



DATA EXTRACTION PLATFORM

# DUMONT

Connectors & Data Pipeline

D O C U M E N T A T I O N

v 2026.1

March 27, 2026



viglet

[viglet.org](http://viglet.org)

# Dumont DEP

## Data Extraction Platform

---

PRODUCT	Dumont DEP — Data Extraction Platform
VERSION	2026.1
GENERATED	March 27, 2026
WEBSITE	<a href="https://viglet.org">viglet.org</a>
SOURCE	<a href="https://github.com/openviglet/dumont">github.com/openviglet/dumont</a>
LICENSE	Apache License 2.0

Dumont DEP is a data extraction platform that connects content sources — websites, databases, file systems, AEM, and WordPress — to search engines through a reliable, asynchronous processing pipeline. This document covers architecture, installation, configuration, connectors, and the developer guide.

# Table of Contents

---

What is Dumont DEP?	1
Core Concepts	5
Viglet Dumont DEP: Installation Guide	13
Configuration Reference	25
Dumont DEP — Architecture Overview	33
Indexing Plugins	40
Connectors Overview	50
Web Crawler Connector	57
Database Connector	63
FileSystem Connector	68
AEM Connector	74
AEM Event Listener Setup	88
WordPress Connector	96
Developer Guide	101
REST API Reference	106
Extending the AEM Connector	118
Extending the Database Connector	132

# What is Dumont DEP?

---

**Viglet Dumont DEP** is an open-source data extraction platform. It connects your content — wherever it lives — to **Viglet Turing ES** for indexing and search.

Think of Dumont DEP as the bridge between your content sources and the search engine. It crawls websites, queries databases, scans file systems, reads AEM repositories, and pulls from WordPress — then delivers every document to Turing ES through a reliable, asynchronous pipeline.

---

## What can you do with Dumont DEP?

---

### Crawl websites

Point the **Web Crawler** at a starting URL and let it discover and extract content from your entire site — pages, articles, documentation. It handles authentication, URL filtering, locale detection, and incremental updates.

### Index databases

The **Database Connector** runs any SQL query against Oracle, PostgreSQL, MariaDB, MySQL, or any JDBC-compatible database and turns each row into a searchable document. Complex joins, custom transformations, and batch processing are all supported.

### Scan file systems

The **FileSystem Connector** walks a directory tree, extracts text from PDFs, Word documents, spreadsheets, presentations, and images (via OCR), and indexes everything with full metadata — file size, extension, modification date.

### Connect to AEM

The **AEM Connector** indexes content from Adobe Experience Manager author and publish instances — supporting content fragments, delta tracking, locale mapping, and custom extension points for your specific AEM setup.

## Integrate with WordPress

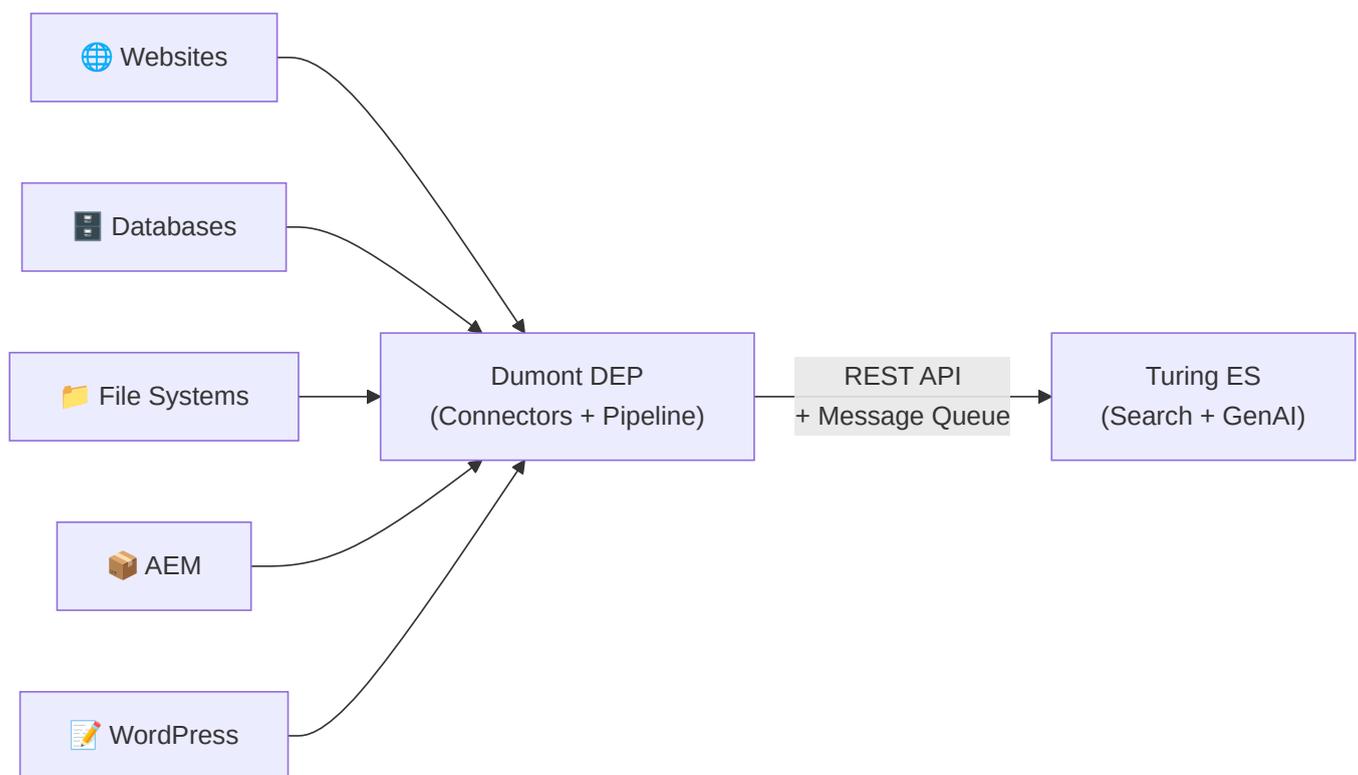
The **WordPress Connector** indexes posts, pages, and custom content types from WordPress installations directly into Turing ES.

## Send to any search engine

Beyond Turing ES, Dumont DEP can deliver content directly to **Apache Solr** or **Elasticsearch** via pluggable indexing adapters – no Turing ES required.

---

## How it works at a glance



Content flows from its original source through Dumont DEP connectors, into an internal message queue for reliable processing, and out to the search engine – where it becomes immediately searchable.

---

## Key concepts

These are the main building blocks you will work with in Dumont DEP. You do not need to understand all of them before getting started – come back to each one as you need it.

CONCEPT	WHAT IT IS	LEARN MORE
<b>Connector</b>	A component that extracts content from a specific source type (web, database, files, AEM, WordPress)	<a href="#">Connectors</a>
<b>Indexing Plugin</b>	The output adapter that delivers documents to a search engine — Turing ES (default), Apache Solr, or Elasticsearch	<a href="#">Indexing Plugins</a>
<b>Job Item</b>	A single document being processed through the pipeline — with fields, metadata, and an action (index, de-index, update)	<a href="#">Core Concepts</a>
<b>Batch Processor</b>	Groups individual job items into configurable batches (default: 50) for efficient queue delivery	<a href="#">Core Concepts</a>
<b>Message Queue</b>	Apache Artemis (embedded) — decouples connector extraction from indexing delivery for reliability and throughput	<a href="#">Architecture</a>
<b>Processing Strategy</b>	Priority-based rules that determine how each document is handled — index, re-index, de-index, ignore, or skip unchanged	<a href="#">Core Concepts</a>
<b>Indexing Rules</b>	Regex-based filters that skip documents matching specific attribute patterns before they reach the queue	<a href="#">Core Concepts</a>

---

## Where to go next

Not sure where to start? Here is a suggested path depending on what you want to do:

**I want to understand how Dumont DEP works** → Read [Core Concepts](#) first, then [Architecture](#).

**I want to install Dumont DEP** → Go to the [Installation Guide](#).

**I want to crawl a website** → Go to [Web Crawler Connector](#).

**I want to index a database** → Go to [Database Connector](#).

**I want to index files from a directory** → Go to [FileSystem Connector](#).

**I want to connect AEM to Turing ES** → Go to [AEM Connector](#).

**I want to index WordPress content** → Go to [WordPress Connector](#).

**I want to send content to Solr or Elasticsearch directly** → Go to [Indexing Plugins](#).

**I want to contribute or build from source** → Go to the [Developer Guide](#).

---

# Core Concepts

---

This page explains the fundamental concepts of Dumont DEP in plain terms. No configuration files, no code – just the mental model you need before diving into the technical documentation.

---

## Connectors

---

A **Connector** is a component that knows how to extract content from a specific type of source. Each connector:

- Connects** to a content source (a website, a database, a file folder, an AEM instance)

- Extracts** documents according to its configuration

- Emits** each document as a **Job Item** into the processing pipeline

Dumont DEP ships with five connectors:

CONNECTOR	SOURCE	HOW IT WORKS
<b>Web Crawler</b>	Websites	Recursively follows links from a starting URL, extracts page content via HTML parsing
<b>Database</b>	JDBC databases	Executes SQL queries and maps each result row to a document
<b>FileSystem</b>	Local/network directories	Walks directory trees, extracts text from files via Apache Tika
<b>AEM</b>	Adobe Experience Manager	Reads content from AEM author/publish instances via the JCR API
<b>WordPress</b>	WordPress sites	Pulls posts, pages, and custom content types from WordPress installations

Each connector implements the `DumConnectorPlugin` interface and provides three operations: **crawl** (full extraction), **indexAll** (re-index a source), and **indexById** (index specific documents by ID).

---

## Job Items

---

A **Job Item** is a single document moving through the pipeline. It contains:

**Fields** — key-value pairs representing the document's content (title, text, URL, date, author, and any custom fields)

**Action** — what to do with this document: **INDEX** (add or update) or **DELETE** (remove from the index)

**Provider** — which connector produced this item

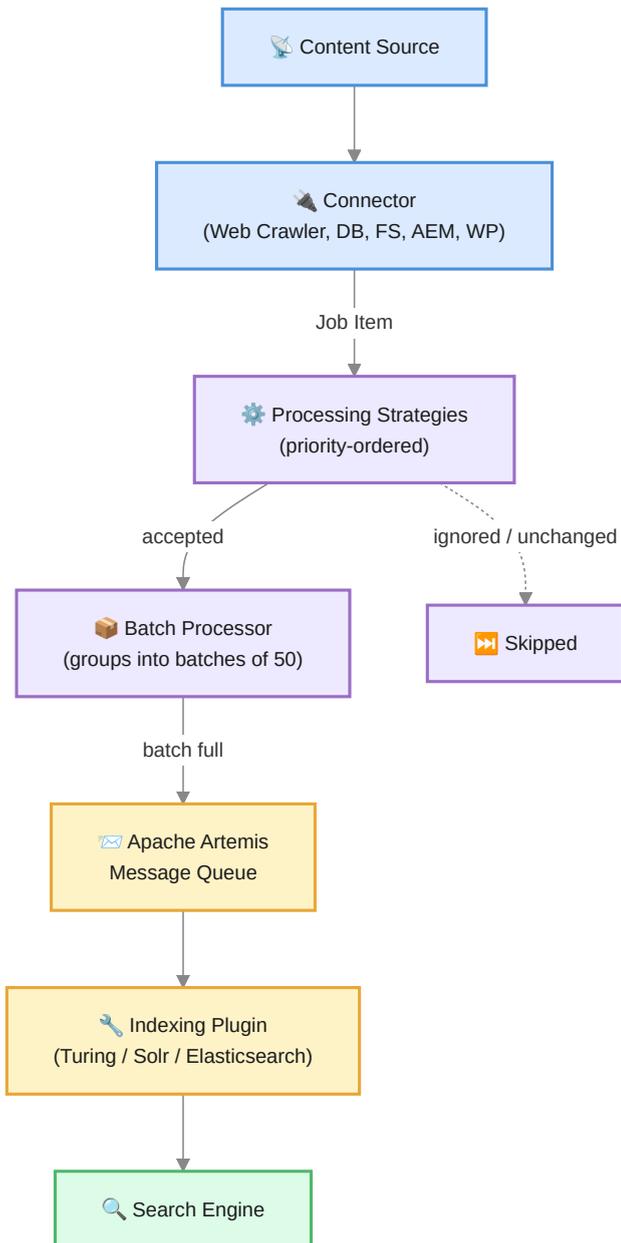
**Locale** — the language/country of the content (e.g., **en\_US**, **pt\_BR**)

Job Items are the universal data format inside Dumont DEP. Every connector produces them, every strategy evaluates them, and every indexing plugin consumes them.

---

## The Processing Pipeline

When a connector extracts a document, it does not send it directly to the search engine. Instead, the document passes through a multi-stage pipeline designed for reliability, efficiency, and flexibility.



### Stage 1 – Extraction

The connector reads content from the source and produces Job Items. Each item includes all the fields needed for indexing.

## Stage 2 — Strategy Evaluation

Every Job Item passes through a chain of **Processing Strategies**, evaluated in priority order. Each strategy decides whether the item should be indexed, re-indexed, de-indexed, ignored, or skipped:

PRIORITY	STRATEGY	WHAT IT DOES
10	<b>De-index</b>	Removes documents marked for deletion
20	<b>Ignore (Indexing Rules)</b>	Skips documents matching regex-based ignore rules
30	<b>Index</b>	Indexes new documents (not seen before)
40	<b>Re-index</b>	Updates documents whose content has changed (checksum comparison)
50	<b>Unchanged</b>	Skips documents that haven't changed since the last run

The strategy chain uses **checksum-based change detection** — each document's content is hashed, and the hash is compared against the last indexed version. Only documents that have actually changed are re-indexed.

## Stage 3 — Batching

Accepted Job Items are collected by the **Batch Processor** into groups (default: 50 items per batch). This reduces the number of messages sent to the queue and improves indexing throughput.

When a batch reaches its configured size, or when the connector signals it has finished, the batch is flushed to the message queue.

## Stage 4 — Queue

The batch is sent to **Apache Artemis** (embedded JMS message queue). The queue decouples extraction from delivery — if the search engine is temporarily unavailable, messages wait in the queue and are delivered when the connection is restored.

Queue messages are persisted to disk ( `store/queue/` ) and survive application restarts.

## Stage 5 – Delivery

The **Indexing Plugin** consumes messages from the queue and delivers them to the configured search engine. Dumont DEP supports three output targets:

PLUGIN	TARGET	DESCRIPTION
<b>Turing</b> (default)	Viglet Turing ES	Uses the Turing Java SDK to deliver documents via REST API
<b>Solr</b>	Apache Solr	Uses SolrJ to add documents directly to a Solr collection
<b>Elasticsearch</b>	Elasticsearch	Uses the Elasticsearch Java Client for bulk indexing

---

## Indexing Rules

---

**Indexing Rules** allow you to filter content during extraction — before it enters the pipeline. A rule defines:

An **attribute** (a field name in the document)

A **rule type** (currently `IGNORE`)

One or more **values** (regex patterns)

When a document's attribute matches any of the values, the document is skipped entirely. For example, a rule with `attribute = template` and `values = [error-page, redirect]` will prevent any document with those templates from being indexed.

Indexing Rules are configured per source in the admin console and are evaluated by the `IgnoreIndexingRuleStrategy` at priority 20 — before the index/re-index/unchanged strategies run.

---

## Change Detection

---

Dumont DEP tracks every document it has processed using a persistent indexing database. For each document, it stores:

**Object ID** — the unique identifier from the source

**Checksum** — a CRC32 hash of the document's content

**Timestamp** — when the document was last indexed

**Status** — the current state (preparing, indexed, de-indexed, etc.)

On subsequent runs, the connector compares the new checksum against the stored one. If they match, the document is **unchanged** and skipped. If they differ, the document is **re-indexed**. If a previously indexed document is no longer present in the source, it is **de-indexed**.

This mechanism enables efficient **incremental indexing** — only changed content is sent to the search engine, regardless of how large the source is.

---

## Indexing Status Values

Every document's journey through the pipeline is tracked with a status code:

STATUS	MEANING
<code>PREPARE_INDEX</code>	Preparing to index the document
<code>PREPARE_UNCHANGED</code>	No changes detected since last indexing
<code>PREPARE_REINDEX</code>	Preparing a re-indexation (content changed)
<code>PREPARE_FORCED_REINDEX</code>	Forced re-indexation triggered
<code>RECEIVED_AND_SENT_TO_TURING</code>	Document received and forwarded to the search engine
<code>SENT_TO_QUEUE</code>	Document placed in the Artemis processing queue
<code>RECEIVED_FROM_QUEUE</code>	Document consumed from the queue by the indexing plugin
<code>INDEXED</code>	Document successfully indexed
<code>FINISHED</code>	Operation finished
<code>DEINDEXED</code>	Document removed from the index
<code>NOT_PROCESSED</code> / <code>IGNORED</code>	Document skipped due to an Indexing Rule or strategy decision

## Ready to go deeper?

---

I WANT TO...	GO TO
Understand the full system architecture	<a href="#">Architecture</a>
Install Dumont DEP	<a href="#">Installation Guide</a>
Configure a Web Crawler	<a href="#">Web Crawler Connector</a>
Index a database	<a href="#">Database Connector</a>
Understand the indexing plugins	<a href="#">Indexing Plugins</a>
See the full configuration reference	<a href="#">Configuration Reference</a>

---

# Viglet Dumont DEP: Installation Guide

Viglet Dumont DEP is an open-source data extraction platform that connects content sources to search engines. The connector application (`dumont-connector.jar`) provides the pipeline infrastructure — message queue, batch processing, strategies, and indexing plugins — while **connector plugins** (AEM, Web Crawler) provide the actual content extraction capabilities.

## ⚠ THE CONNECTOR JAR ALONE DOES NOT CRAWL CONTENT

`dumont-connector.jar` is the pipeline engine. To actually extract content from sources, you must add a **connector plugin JAR** (AEM or Web Crawler) to the classpath via `-Dloader.path`. Without a plugin, the connector has no data source to crawl.

## How It Works

```
dumont-connector.jar    ← Pipeline engine (queue, strategies, indexing)
├─ libs/                ← Plugin directory (added via -Dloader.path)
│  └─ aem-plugin.jar    ← AEM connector plugin
│  └─ web-crawler-plugin.jar ← Web Crawler connector plugin
│  └─ my-custom-plugin.jar ← Your custom extensions (optional)
```

The connector JAR is built with Spring Boot's **ZIP layout** (PropertiesLauncher), which allows loading additional JARs from external directories at runtime. When you pass `-Dloader.path=libs`, Spring Boot scans the `libs/` directory and adds all JARs to the application classpath — making the connector plugins available to the pipeline.

## Installing Java

Dumont DEP requires Java.

### ⓘ NOTE

Dumont DEP only supports **Java 21** or later.

Set the `JAVA_HOME` environment variable and add `$JAVA_HOME/bin` to your `PATH`.

---

## Hardware Requirements

MINIMUM	RECOMMENDED
Single 1.5 GHz processor, 1 GB available memory	Dual 2 GHz processors, 2 GB available memory

In addition:

500 MB of disk space for software

Additional disk space proportional to the content being crawled (Web Crawler caches visited URLs)

---

## Installing Dumont DEP

### Option 1 — Docker (fastest)

```
docker pull openviglet/dumont:2026.1
docker run -p 30130:30130 openviglet/dumont:2026.1
```

### Option 2 — Download JARs

Download `dumont-connector.jar` and the connector plugin JARs from the [releases page](#):

```
mkdir -p /appl/viglet/dumont/server/libs

# Main connector engine
cp dumont-connector.jar /appl/viglet/dumont/server/

# Connector plugins (choose one or both)
cp aem-plugin.jar /appl/viglet/dumont/server/libs/
cp web-crawler-plugin.jar /appl/viglet/dumont/server/libs/
```

## Option 3 – Build from source

```
git clone https://github.com/openviglet/dumont.git
cd dumont
mvn clean install

# Copy the connector engine
cp connector/connector-app/target/dumont-connector.jar
/appl/viglet/dumont/server/

# Copy connector plugins
mkdir -p /appl/viglet/dumont/server/libs
cp aem/aem-plugin/target/aem-plugin.jar /appl/viglet/dumont/server/libs/
cp web-crawler/wc-plugin/target/web-crawler-plugin.jar
/appl/viglet/dumont/server/libs/
```

## Starting with a Connector Plugin

The key to running Dumont DEP is the `-Dloader.path` JVM property. It tells Spring Boot's `PropertiesLauncher` where to find the connector plugin JARs.

### With AEM Plugin

```
java -Dloader.path=libs \  
-Dturing.url=http://localhost:2700 \  
-Dturing.apiKey=<YOUR_TURING_API_KEY> \  
-jar dumont-connector.jar
```

### With Web Crawler Plugin

```
java -Dloader.path=libs \  
-Dturing.url=http://localhost:2700 \  
-Dturing.apiKey=<YOUR_TURING_API_KEY> \  
-jar dumont-connector.jar
```

#### ONE CONNECTOR PLUGIN PER JVM

Currently, only **one connector plugin** can be active per JVM instance. To run both AEM and Web Crawler, start two separate instances of `dumont-connector.jar` — each with its own `libs/` directory containing the appropriate plugin, and each on a different port.

### With Custom AEM Extensions

If you have a custom AEM extension (like the `aem-plugin-sample`), add its JAR to the same `libs/` directory:

```
ls libs/
aem-plugin.jar
aem-plugin-sample.jar    # Your custom extensions

java -Dloader.path=libs \
     -Dturing.url=http://localhost:2700 \
     -Dturing.apiKey=<YOUR_TURING_API_KEY> \
     -jar dumont-connector.jar
```

The configuration JSON references extension classes by fully-qualified name (e.g., `com.viglet.dumont.connector.aem.sample.ext.DumAemExtSampleModelJson`). These classes are found on the classpath because the extension JAR is in the `libs/` directory.

### **LOADER.PATH ACCEPTS DIRECTORIES OR INDIVIDUAL JARS**

You can point to a directory or a specific JAR: `-Dloader.path=libs` or `-Dloader.path=/path/to/aem-plugin.jar`

## Database Configuration

By default, Dumont DEP uses **H2** as its embedded database for tracking indexing state. For production, replace with PostgreSQL.

### H2 (default — development only)

No configuration needed. The database is created automatically at `./store/db/dumontDB`.

### PostgreSQL (recommended for production)

Create a database and user:

```
CREATE USER dumont PASSWORD '<password>' WITH CREATEUSER;
CREATE DATABASE dumont_db;
ALTER DATABASE dumont_db OWNER TO dumont;
```

Configure via JVM properties:

```
java -Dloader.path=libs \  
-Dspring.datasource.url=jdbc:postgresql://localhost:5432/dumont_db \  
-Dspring.datasource.username=dumont \  
-Dspring.datasource.password=<password> \  
-Dspring.datasource.driver-class-name=org.postgresql.Driver \  
-Dturing.url=http://localhost:2700 \  
-Dturing.apiKey=<YOUR_TURING_API_KEY> \  
-jar dumont-connector.jar
```

## Connecting to Turing ES

Dumont DEP sends content to Turing ES by default. Configure the connection:

```
turing:  
url: http://localhost:2700  
apiKey: <YOUR_TURING_API_KEY>
```

Create an API Token in **Turing ES → Administration → API Tokens** and use it as the `apiKey` value.

### API TOKEN REQUIRED

Dumont DEP cannot send content to Turing ES without a valid API Token. Create one in the Turing ES admin console before starting Dumont DEP.

## Creating a Linux Service

Each connector plugin runs as a separate service. Create a dedicated directory per connector with its own configuration.

### Directory Layout (AEM example)

```
/appl/viglet/dumont/aem/  
├─ dumont-connector.jar  
├─ dumont-connector.properties    # Spring Boot external config  
└─ libs/  
    └─ aem-plugin.jar            # Connector plugin
```

### External Properties File

Create `/appl/viglet/dumont/aem/dumont-connector.properties`:

```
turing.url=http://localhost:2700  
turing.apiKey=<YOUR_TURING_API_KEY>  
server.port=30130
```

### Systemd Service File

As root, create `/etc/systemd/system/dumont-aem.service`:

```
[Unit]
Description=Viglet Dumont DEP (AEM)
After=syslog.target

[Service]
User=turing
Group=turing
WorkingDirectory=/appl/viglet/dumont/aem
ExecStart=/appl/java/jdk21/bin/java -Xmx512m -Xms512m \
  -Dloader.path=/appl/viglet/dumont/aem/libs/ \
  -jar /appl/viglet/dumont/aem/dumont-connector.jar \
  --spring.config.additional-location=file:/appl/viglet/dumont/aem/dumont-
connector.properties
SuccessExitStatus=143

[Install]
WantedBy=multi-user.target
```

Enable and start the service:

```
systemctl daemon-reload
systemctl enable dumont-aem.service
systemctl start dumont-aem.service
```

## Running Multiple Connectors

Since only one connector plugin is supported per JVM, run separate services for each connector — each with its own directory, config, and port:

```
/appl/viglet/dumont/
├── aem/                                # AEM connector (port 30130)
│   ├── dumont-connector.jar
│   ├── dumont-connector.properties
│   └── libs/
│       └── aem-plugin.jar
└── webcrawler/                          # Web Crawler connector (port 30131)
    ├── dumont-connector.jar
    ├── dumont-connector.properties
    └── libs/
        └── web-crawler-plugin.jar
```

Create a separate systemd service for each (e.g., `dumont-aem.service`, `dumont-webcrawler.service`), each pointing to its own `WorkingDirectory` and properties file with a different `server.port`.

## Standalone Connectors (Database & FileSystem)

The **Database** and **FileSystem** connectors are **standalone CLI tools** — they do not run as plugins inside the connector application. Instead, they connect to a running Dumont DEP instance via its REST API.

### Database Connector

```
java -cp dumont-db-indexer.jar com.viglet.dumont.connector.db.DumDbImportTool \
  --server http://localhost:30130 \
  --api-key <API_KEY> \
  --driver org.mariadb.jdbc.Driver \
  --connect "jdbc:mariadb://localhost:3306/products" \
  --query "SELECT id, name, description FROM products" \
  --site ProductCatalog \
  --locale en_US
```

## FileSystem Connector

```
java -cp dumont-filesystem-indexer.jar
com.viglet.dumont.filesystem.DumFSImportTool \
  --source-dir /mnt/shared/documents \
  --server http://localhost:30130 \
  --api-key <API_KEY> \
  --site InternalDocs \
  --locale en_US
```

These tools run independently and can be scheduled via cron jobs or CI/CD pipelines.

---

## Docker Compose

A full stack with Dumont DEP, Turing ES, Solr, and MariaDB:

```
services:
  dumont:
    image: openviglet/dumont:2026.1
    ports:
      - "30130:30130"
    environment:
      TURING_URL: http://turing:2700
      TURING_APIKEY: <YOUR_API_KEY>

  turing:
    image: openviglet/turing:2026.1
    ports:
      - "2700:2700"

  solr:
    image: solr:9
    ports:
      - "8983:8983"

  mariadb:
    image: mariadb:11
    environment:
      MARIADB_ROOT_PASSWORD: root
      MARIADB_DATABASE: turing
      MARIADB_USER: turing
      MARIADB_PASSWORD: turing
```

## Verifying the Installation

After starting Dumont DEP, verify it is running:

```
curl http://localhost:30130/api/v2/connector/status
```

A successful response confirms the service is up and ready to accept connector configurations.

## Directory Structure Summary

A complete production deployment with two connectors:

```
/appl/viglet/dumont/
├── aem/                                # AEM connector instance
│   ├── dumont-connector.jar           # Pipeline engine
│   ├── dumont-connector.properties   # Turing URL, API key, port
│   ├── libs/
│   │   └── aem-plugin.jar            # AEM connector plugin
│   └── store/                         # Auto-created at runtime
│       ├── db/                       # H2 indexing database
│       ├── queue/                    # Artemis persistent queue
│       └── logs/                     # Application logs
└── webcrawler/                       # Web Crawler connector instance
    ├── dumont-connector.jar
    ├── dumont-connector.properties
    ├── libs/
    │   └── web-crawler-plugin.jar
    └── store/
```

## Related Pages

PAGE	DESCRIPTION
<a href="#">Configuration Reference</a>	All application.yaml properties
<a href="#">Architecture</a>	System components and data flow
<a href="#">Connectors Overview</a>	Available connectors and how to configure them

# Configuration Reference

---

This page documents every significant property in the Dumont DEP `application.yaml` file. Any property can be overridden via environment variables, JVM system properties (`-Dkey=value`), or a separate `application-production.yaml` file.

## OVERRIDE PATTERN

To override a property at runtime, use the environment variable convention: replace `.` and `-` with `_` and uppercase everything. For example, `dumont.indexing.provider` becomes `DUMONT_INDEXING_PROVIDER=solr`.

## Full Default Configuration

```
spring:
  datasource:
    url:
jdbc:h2:file:./store/db/dumontDB;DATABASE_TO_UPPER=false;CASE_INSENSITIVE_IDENTIFIERS=true
    username: sa
    password: ""
    driver-class-name: org.h2.Driver
  h2:
    console:
      enabled: false
      path: /h2
      settings:
        web-allow-others: true
  jpa:
    properties:
      hibernate:
        "[globally_quoted_identifiers]": true
        "[enable_lazy_load_no_trans]": true
    hibernate:
      ddl-auto: update
      show-sql: false
  artemis:
    mode: embedded
    broker-url: localhost:61616
    embedded:
      enabled: true
      persistent: true
      data-directory: store/queue
      queues: connector-indexing.queue
  jms:
    template:
      default-destination: connector-indexing.queue

server:
  port: ${PORT:30130}

dumont:
  indexing:
    provider: turing
    job:
      batch-size: 50
  solr:
    url: http://localhost:8983/solr
    collection: dumont
  elasticsearch:
```

```
url: http://localhost:9200
index: dumont
username: ~
password: ~

turing:
url: http://localhost:2700
apiKey: ""

aem.querybuilder: false
aem.querybuilder.parallelism: 10

logging:
file:
name: store/logs/dum-connector.log
level:
com:
viglet: INFO
```

## Property Reference

### Server

PROPERTY	DEFAULT	DESCRIPTION
<code>server.port</code>	<code>30130</code>	HTTP port for the Dumont DEP application

### Database

PROPERTY	DEFAULT	DESCRIPTION
<code>spring.datasource.url</code>	<code>jdbc:h2:file:./store/db/dumontDB</code>	JDBC URL for the indexing state database
<code>spring.datasource.username</code>	<code>sa</code>	Database username
<code>spring.datasource.password</code>	<code>""</code>	Database password
<code>spring.datasource.driver-class-name</code>	<code>org.h2.Driver</code>	JDBC driver class
<code>spring.h2.console.enabled</code>	<code>false</code>	Enable H2 web console — <b>do not enable in production</b>

## Message Queue (Apache Artemis)

PROPERTY	DEFAULT	DESCRIPTION
<code>spring.artemis.mode</code>	<code>embedded</code>	<code>embedded</code> (default) or <code>native</code> (external broker)
<code>spring.artemis.broker-url</code>	<code>localhost:61616</code>	Broker URL when using <code>native</code> mode
<code>spring.artemis.embedded.enabled</code>	<code>true</code>	Enable the embedded broker
<code>spring.artemis.embedded.persistent</code>	<code>true</code>	Persist queue messages to disk (survives restarts)
<code>spring.artemis.embedded.data-directory</code>	<code>store/queue</code>	Directory for persisted queue data
<code>spring.artemis.embedded.queues</code>	<code>connector-indexing.queue</code>	Queue name created on startup

## Indexing Configuration

PROPERTY	DEFAULT	DESCRIPTION
<code>dumont.indexing.provider</code>	<code>turing</code>	Output target: <code>turing</code> , <code>solr</code> , or <code>elasticsearch</code>
<code>dumont.indexing.job.batch-size</code>	<code>50</code>	Number of Job Items per batch before queue delivery

## Turing ES Connection

PROPERTY	DEFAULT	DESCRIPTION
<code>turing.url</code>	<code>http://localhost:2700</code>	Base URL of the Turing ES instance
<code>turing.apiKey</code>	<code>""</code>	API Token for authenticating with Turing ES

## Solr Direct Connection

PROPERTY	DEFAULT	DESCRIPTION
<code>dumont.indexing.solr.url</code>	<code>http://localhost:8983/solr</code>	Apache Solr base URL
<code>dumont.indexing.solr.collection</code>	<code>dumont</code>	Solr collection name

## Elasticsearch Direct Connection

PROPERTY	DEFAULT	DESCRIPTION
<code>dumont.indexing.elasticsearch.url</code>	<code>http://localhost:9200</code>	Elasticsearch base URL
<code>dumont.indexing.elasticsearch.index</code>	<code>dumont</code>	Elasticsearch index name
<code>dumont.indexing.elasticsearch.username</code>	<i>(none)</i>	Optional authentication username
<code>dumont.indexing.elasticsearch.password</code>	<i>(none)</i>	Optional authentication password

## AEM QueryBuilder

PROPERTY	DEFAULT	DESCRIPTION
<code>dumont.aem.querybuilder</code>	<code>false</code>	Enable QueryBuilder-based content discovery instead of tree traversal during full indexing
<code>dumont.aem.querybuilder.parallelism</code>	<code>10</code>	Number of parallel threads for processing discovered paths

## Logging

PROPERTY	DEFAULT	DESCRIPTION
<code>logging.file.name</code>	<code>store/logs/dum-connector.log</code>	Log file path
<code>logging.level.com.viglet</code>	<code>INFO</code>	Log level for Dumont DEP application code

## Common Production Overrides

---

A minimal production override:

```
spring:
  datasource:
    url: jdbc:postgresql://db-host:5432/dumont_db
    username: dumont
    password: strong_password
    driver-class-name: org.postgresql.Driver
  h2:
    console:
      enabled: false

dumont:
  indexing:
    provider: turing

turing:
  url: https://search.yourcompany.com
  apiKey: your-turing-api-token

server:
  port: 30130
```

# Dumont DEP — Architecture Overview

---

## Introduction

---

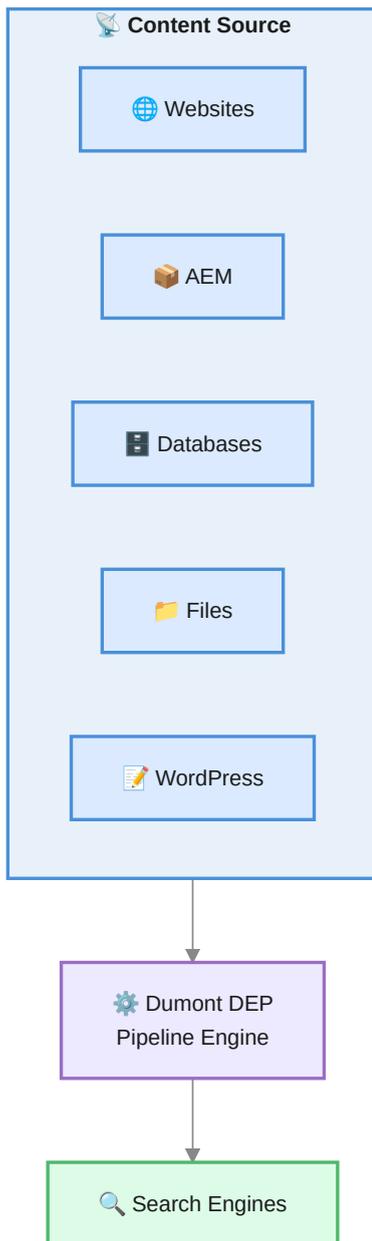
Viglet Dumont DEP is a modular data extraction platform that runs connectors independently and delivers indexed documents to a search engine via an asynchronous message queue. It is designed as a companion application to Viglet Turing ES, but can also deliver content directly to Apache Solr or Elasticsearch.

---

## System Overview

---

The system has three layers: content sources on the left, the Dumont DEP pipeline in the center, and search engines on the right.

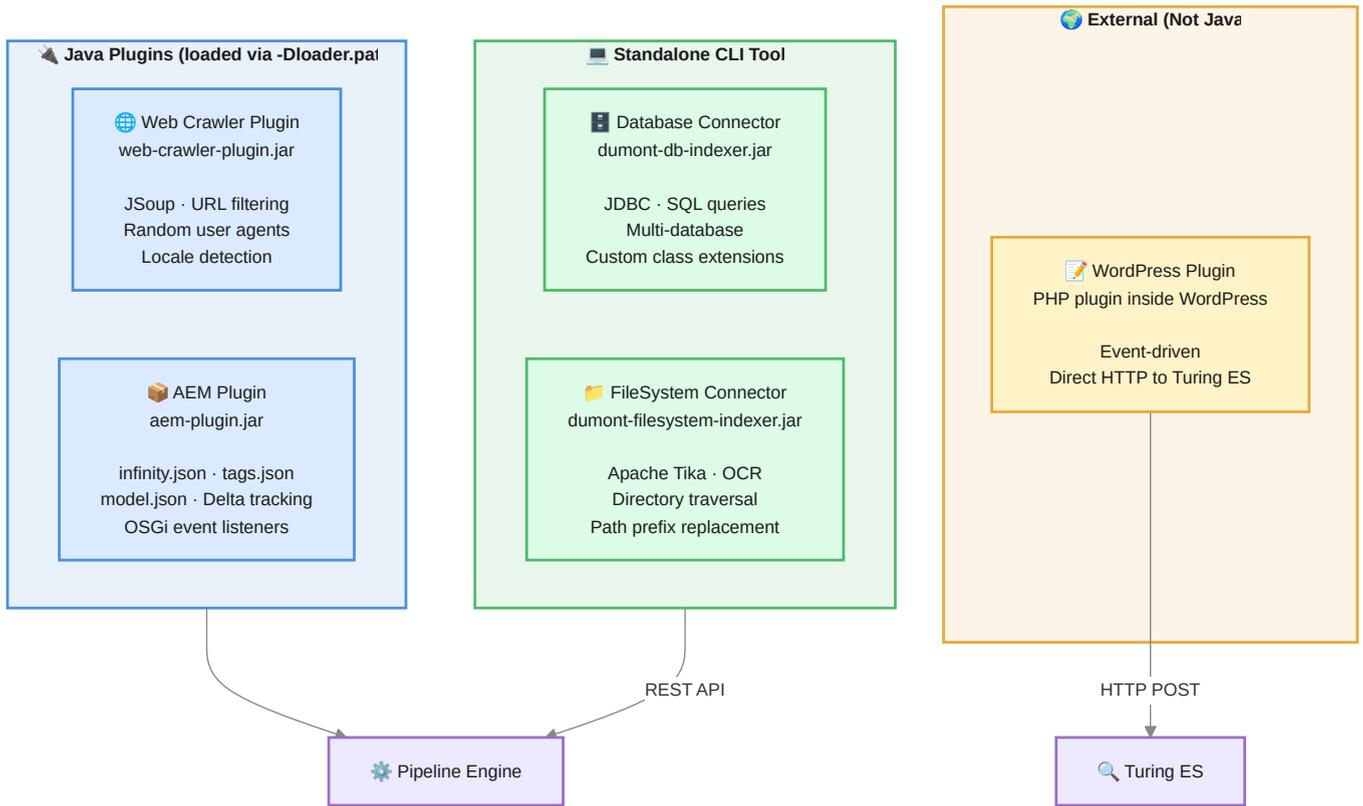


Each numbered block is detailed in its own diagram below.

---

## ① Connectors — How Content Enters the Pipeline

Connectors extract content and feed it into the pipeline. They come in three forms: Java plugins, standalone CLI tools, and the WordPress PHP plugin.



TYPE	CONNECTORS	HOW THEY CONNECT
Java Plugins	Web Crawler, AEM	Loaded into <code>dumont-connector.jar</code> via <code>-Dloader.path</code> — one plugin per JVM
Standalone CLI	Database, FileSystem	Separate JARs that connect to a running Dumont DEP instance via REST API
PHP Plugin	WordPress	Installed inside WordPress — sends content directly to Turing ES, bypasses Dumont

For details on each connector, see [Connectors Overview](#).

## ② Pipeline Engine – How Content Is Processed

Once a connector produces a Job Item, it passes through a multi-stage pipeline before reaching the search engine.

STAGE	COMPONENT	WHAT IT DOES
①	<b>Job Item</b>	A single document with fields, an action (INDEX / DELETE), and a locale
②	<b>Processing Strategies</b>	Priority chain: deindex (P10) → ignore rules (P20) → index new (P30) → reindex changed (P40) → skip unchanged (P50)
③	<b>Batch Processor + Queue</b>	Groups items into batches of 50, sends to Apache Artemis persistent queue
④	<b>Indexing Plugin</b>	Delivers to Turing ES (default), Apache Solr, or Elasticsearch

The **Indexing DB** stores checksums and status for every processed document — enabling incremental indexing on subsequent runs.

For the conceptual explanation of each stage, see [Core Concepts — The Processing Pipeline](#).

## ③ Data Flow — Indexing Sequence

The complete sequence from content source to search engine:

## Internal Module Structure

MODULE	PACKAGE	RESPONSIBILITY
Connector Core	<code>connector</code>	Plugin interface, session management, and REST API controllers
Processing Strategies	<code>connector/strategy</code>	Priority-based chain — index, re-index, de-index, ignore, or skip
Batch Processor	<code>connector/batch</code>	Thread-safe buffer that groups Job Items before queue delivery
Message Queue	<code>connector/queue</code>	JMS listener on Apache Artemis — delegates to indexing plugins
Indexing Plugins	<code>connector/indexing</code>	Output adapters: Turing ES, Apache Solr, Elasticsearch
Web Crawler	<code>web-crawler</code>	JSoup, URL filtering, authentication, locale detection
Database	<code>db</code>	JDBC queries, batch chunking, multi-database support
FileSystem	<code>filesystem</code>	Apache Tika text extraction, OCR, metadata mapping
AEM	<code>aem</code>	infinity.json, tags, model.json, delta tracking, custom extensions
WordPress	<code>wordpress</code>	PHP plugin — event-driven indexing inside WordPress
Commons	<code>commons</code>	Shared models, interfaces, utilities
AEM Commons	<code>aem-commons</code>	AEM extension interfaces (published to Maven Central)
DB Commons	<code>db-commons</code>	DB extension interface (published to Maven Central)

## Technology Stack

LAYER	TECHNOLOGY	NOTES
Runtime	Java 21	Minimum supported version
Framework	Spring Boot 4.0.3	JMS, caching, async, scheduling
Message Broker	Apache Artemis	Embedded, persistent queues
Database	H2 (dev) / PostgreSQL (prod)	Indexing state, checksums, config
HTML Parsing	JSoup 1.22.1	Web Crawler
Text Extraction	Apache Tika 3.2.3	FileSystem – PDF, DOCX, images (OCR)
Search Clients	Turing Java SDK, SolrJ 10.0.0, ES Client 9.3.2	One active per deployment
Build	Apache Maven	Multi-module project

## Deployment Topologies

### Development

```
Dumont DEP (H2 embedded + Artemis embedded)  
→ Turing ES (http://localhost:2700)
```

### Production

```
Dumont DEP + PostgreSQL  
→ Turing ES + Apache Solr
```

## Direct Indexing (without Turing ES)

Dumont DEP + PostgreSQL  
→ Apache Solr (direct via SolrJ)  
→ Elasticsearch (direct via ES Client)

## Related Pages

PAGE	DESCRIPTION
<a href="#">Core Concepts</a>	Pipeline stages, strategies, and change detection
<a href="#">Connectors Overview</a>	All connectors and deployment types
<a href="#">Indexing Plugins</a>	Turing ES, Solr, Elasticsearch output targets
<a href="#">Installation Guide</a>	Setup with <code>-Dloader.path</code> and systemd
<a href="#">Turing ES — Architecture</a>	Search-side architecture — indexing reception, search flow, and deployment topologies

# Indexing Plugins

An **Indexing Plugin** is the output adapter that delivers processed documents from the Dumont DEP pipeline to a search engine. Dumont DEP supports three targets — you choose one per deployment.

## Available Plugins

PLUGIN	TARGET	CLIENT LIBRARY	USE CASE
<b>Turing</b> (default)	Viglet Turing ES	Turing Java SDK 2026.1.17	Full enterprise search with GenAI, facets, spotlights
<b>Solr</b>	Apache Solr	SolrJ 10.0.0	Direct Solr indexing without Turing ES
<b>Elasticsearch</b>	Elasticsearch	ES Java Client 9.3.2	Direct Elasticsearch indexing without Turing ES

## Turing ES Plugin (default)

The default plugin delivers documents to Turing ES via its REST API, using the official Turing Java SDK.

## Configuration

```
dumont:  
  indexing:  
    provider: turing  
  
turing:  
  url: http://localhost:2700  
  apiKey: your-turing-api-token
```

PROPERTY	DESCRIPTION
<code>turing.url</code>	Base URL of the Turing ES instance
<code>turing.apiKey</code>	API Token created in Turing ES → Administration → API Tokens

## How It Works

Receives a batch of Job Items from the message queue

Creates a `TurSNServer` instance using the Turing Java SDK

Calls `TurSNJobUtils.importItems()` to submit the batch

Turing ES validates each document against the target SN Site configuration

Documents are queued internally in Turing ES for Solr indexing

### API TOKEN REQUIRED

The Turing plugin cannot deliver content without a valid API Token. Create one in **Turing ES → Administration → API Tokens** before starting Dumont DEP.

## Apache Solr Plugin

The Solr plugin delivers documents directly to an Apache Solr collection, bypassing Turing ES entirely. Use this when you want Dumont DEP as a pure data extraction tool without Turing ES features.

### Configuration

```
dumont:
  indexing:
    provider: solr
    solr:
      url: http://localhost:8983/solr
      collection: my-collection
```

PROPERTY	DESCRIPTION
<code>dumont.indexing.solr.url</code>	Apache Solr base URL
<code>dumont.indexing.solr.collection</code>	Target Solr collection name

### How It Works

Receives a batch of Job Items from the message queue

Converts each Job Item into a `SolrInputDocument`

Adds all documents to the Solr collection via SolrJ

Commits the changes

The Solr client is cleaned up automatically when the application shuts down (`@PreDestroy`).

## Elasticsearch Plugin

The Elasticsearch plugin delivers documents directly to an Elasticsearch index using bulk requests.

## Configuration

```
dumont:
  indexing:
    provider: elasticsearch
    elasticsearch:
      url: http://localhost:9200
      index: my-index
      username: ~
      password: ~
```

PROPERTY	DESCRIPTION
<code>dumont.indexing.elasticsearch.url</code>	Elasticsearch base URL
<code>dumont.indexing.elasticsearch.index</code>	Target index name
<code>dumont.indexing.elasticsearch.username</code>	Optional authentication username
<code>dumont.indexing.elasticsearch.password</code>	Optional authentication password

## How It Works

Receives a batch of Job Items from the message queue

Builds a bulk request containing all documents

Submits the bulk request to Elasticsearch

Logs any per-document errors from the bulk response

Authentication is optional — leave `username` and `password` empty for unauthenticated clusters.

## Why Use Turing ES Instead of Solr or Elasticsearch Directly?

---

The Solr and Elasticsearch plugins deliver raw documents to the search engine — your application is responsible for everything else: building queries, rendering facets, managing spotlights, handling locales, and building the search UI.

**Turing ES adds an entire enterprise search platform on top of the search engine.** Here's what you get by choosing the Turing plugin over direct Solr or Elasticsearch indexing:

## Search Experience Layer

CAPABILITY	WITH TURING ES	DIRECT SOLR / ELASTICSEARCH
<b>Faceted navigation</b>	Configured per site — facet types, sort, AND/OR operators, secondary facets, custom facets with ranges	Build manually with query params
<b>Spotlights</b>	Curated results pinned to search terms, injected at configured positions	Not available — build from scratch
<b>Targeting Rules</b>	Filter results by user profile (department, role, country) at query time	Not available
<b>Merge Providers</b>	Combine documents from two connectors into one enriched result using a join key	Not available
<b>Spell check</b>	Built-in with auto-correction mode	Configure Solr/ES suggester manually
<b>Autocomplete</b>	Ready-to-use endpoint	Configure Solr/ES suggester manually
<b>More Like This</b>	One toggle per site	Configure MLT handler manually
<b>Result ranking</b>	Boost rules with conditions and weights via admin UI	Write boost queries manually
<b>Highlighting</b>	Configurable HTML tags per site	Configure highlight params manually
<b>Self-describing JSON</b>	Response includes pre-built links for pagination, facet filters, locale switching — the front-end is a pure rendering layer	Build all query logic client-side

## Generative AI & RAG

CAPABILITY	WITH TURING ES	DIRECT SOLR / ELASTICSEARCH
<b>RAG (Retrieval-Augmented Generation)</b>	Documents are automatically embedded as vectors during indexing — users ask questions in natural language and get grounded answers	Not available
<b>AI Agents</b>	Compose assistants that combine LLM + search + web browsing + code execution + MCP tools	Not available
<b>Chat interface</b>	Built-in UI with direct LLM, Semantic Navigation, and AI Agent tabs	Not available
<b>Knowledge Base</b>	Upload files to MinIO — automatically indexed as vector embeddings for RAG	Not available
<b>LLM providers</b>	Anthropic Claude, OpenAI, Azure OpenAI, Google Gemini, Ollama — configured via admin UI	Not available
<b>Tool calling</b>	27 native tools across 7 categories available to AI Agents	Not available
<b>Token usage monitoring</b>	Dashboard showing LLM consumption by model, day, and month	Not available

## Administration & Operations

CAPABILITY	WITH TURING ES	DIRECT SOLR / ELASTICSEARCH
<b>Admin console</b>	Browser-based React UI for all configuration	Solr Admin UI (limited) / Kibana (separate)
<b>Multi-site</b>	Multiple SN Sites on one instance, each with independent fields, facets, and AI settings	Manage collections/indices manually
<b>Multi-language</b>	Locale-aware indexing and search with per-locale Solr cores	Configure manually per core/index
<b>Connector management</b>	Integration page with monitoring, stats, double-check, and indexing manager	Not available
<b>Search metrics</b>	Top search terms by day/week/month/all-time per site	Not available out of the box
<b>Application logs</b>	MongoDB-backed log viewer in the admin console	External log tools required
<b>Security (SSO)</b>	Keycloak OAuth2/OIDC with SAML, LDAP, social login, and MFA	Configure separately per product

## Integration

CAPABILITY	WITH TURING ES	DIRECT SOLR / ELASTICSEARCH
<b>REST API</b>	Self-describing search response with pre-built navigation links	Raw query/response
<b>GraphQL</b>	Built-in endpoint	Not available
<b>Java SDK</b>	Official typed client on Maven Central	Use SolrJ / ES Client directly
<b>JavaScript SDK</b>	Official <a href="#">@viglet/turing-sdk</a> on npm	No official SDK

## When to Use Direct Solr or Elasticsearch

The Solr and Elasticsearch plugins are appropriate when:

- You already have a search infrastructure and only need Dumont DEP as a **data extraction tool**
- You have your own search UI and query layer built on top of Solr/Elasticsearch
- You don't need the GenAI, faceted navigation, or admin console features
- You're integrating with a third-party system that requires direct Solr/Elasticsearch access

For all other cases — especially when building a new search experience — **Turing ES is the recommended target** because it provides a complete, ready-to-use enterprise search platform with GenAI capabilities on top of the same Apache Solr engine.

---

## Switching Plugins

Change the active plugin by setting `dumont.indexing.provider`:

```
# Via JVM property
java -Ddumont.indexing.provider=solr -jar viglet-dumont.jar

# Via environment variable
DUMONT_INDEXING_PROVIDER=elasticsearch
java -jar viglet-dumont.jar
```

Only one plugin is active per deployment. All connectors share the same output target.

---

## Related Pages

---

PAGE	DESCRIPTION
<a href="#">Configuration Reference</a>	All application.yaml properties
<a href="#">Architecture</a>	Where indexing plugins fit in the pipeline
<a href="#">Turing ES – REST API</a>	Turing ES indexing API reference
<a href="#">Turing ES – Semantic Navigation</a>	SN Sites, facets, spotlights, and search configuration
<a href="#">Turing ES – RAG</a>	Retrieval-Augmented Generation and vector embeddings

---

# Connectors Overview

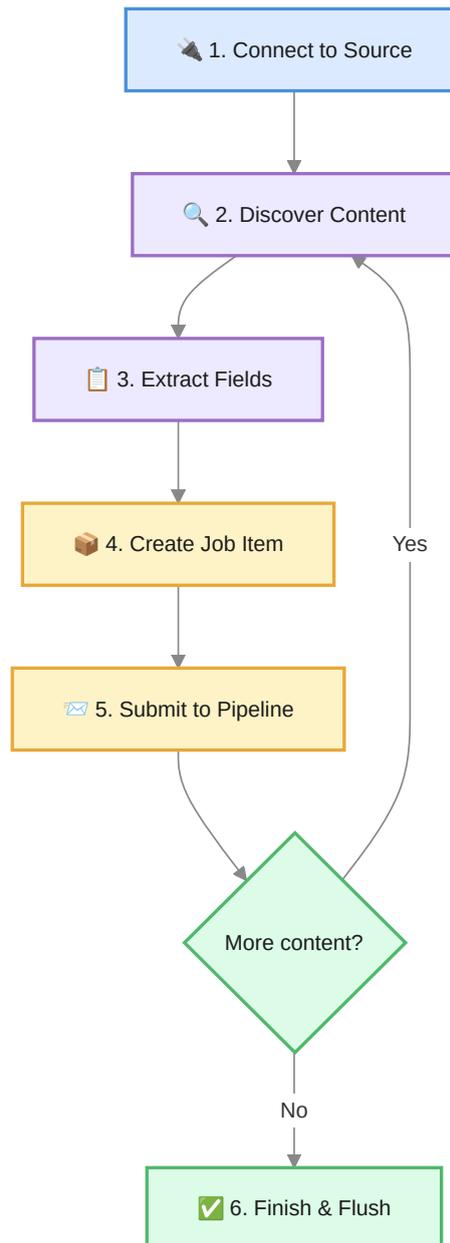
Connectors are the components that extract content from external sources and feed it into the Dumont DEP processing pipeline. Each connector specializes in a specific type of content source and knows how to navigate, extract, and map content into Job Items that the pipeline can process.

## Available Connectors

CONNECTOR	SOURCE TYPE	DEPLOYMENT	ARTIFACT
<b>Web Crawler</b>	Websites	<b>Java Plugin</b> — loaded via <code>-Dloader.path</code>	<code>web-crawler-plugin.jar</code>
<b>AEM</b>	Adobe Experience Manager	<b>Java Plugin</b> — loaded via <code>-Dloader.path</code>	<code>aem-plugin.jar</code>
<b>Database</b>	JDBC databases	<b>Standalone Java CLI</b> — runs independently	<code>dumont-db-indexer.jar</code>
<b>FileSystem</b>	Local/network directories	<b>Standalone Java CLI</b> — runs independently	<code>dumont-filesystem-indexer.jar</code>
<b>WordPress</b>	WordPress sites	<b>PHP Plugin</b> — installed inside WordPress	<code>viglet-turing-for-wordpress/</code>

## How Connectors Work

Every connector follows the same lifecycle:



**Connect** — Establish a connection to the content source (HTTP, JDBC, file handle, JCR session)

**Discover** — Find content to process (follow links, execute query, list files, traverse nodes)

**Extract** — Pull field values from each content item (title, text, URL, date, custom fields)

**Create** — Build a Job Item with the extracted fields, an action (INDEX/DELETE), and metadata

**Submit** — Pass the Job Item into the processing pipeline (strategies → batch → queue)

**Finish** — Flush any remaining items in the batch processor and signal completion

## Connector Interface

All connectors implement the `DumConnectorPlugin` interface:

METHOD	DESCRIPTION
<code>crawl()</code>	Full extraction — discover and process all content from the source
<code>indexAll(source)</code>	Re-index all content from a specific source
<code>indexById(source, contentIds)</code>	Index specific documents by their IDs
<code>getProviderName()</code>	Returns the connector's identifier (e.g., <code>web-crawler</code> , <code>database</code> )

## Connector Plugins vs. Standalone Tools

Dumont DEP connectors are distributed in two forms:

### Connector Plugins (AEM, Web Crawler)

The **AEM** and **Web Crawler** connectors are **plugin JARs** that run inside the `dumont-connector.jar` pipeline. They must be placed in a `libs/` directory and loaded via Spring Boot's `-Dloader.path`:

```
# Directory layout
dumont-connector.jar
libs/
├─ aem-plugin.jar
└─ web-crawler-plugin.jar

# Launch with plugins on the classpath
java -Dloader.path=libs -jar dumont-connector.jar
```

#### DUMONT-CONNECTOR.JAR ALONE DOES NOT CRAWL

The connector JAR provides only the pipeline infrastructure (queue, strategies, indexing). Without a plugin JAR on the classpath, there is no data source to extract content from. You must add exactly **one** connector plugin via `-Dloader.path`.

#### ONE PLUGIN PER JVM INSTANCE

Only **one connector plugin** can be loaded per JVM instance. To run multiple connectors (e.g., AEM and Web Crawler), start separate `dumont-connector.jar` instances — each with its own plugin and port.

### Standalone CLI Tools (Database, FileSystem)

The **Database** and **FileSystem** connectors are **standalone command-line tools** — separate JARs that run independently and connect to a running Dumont DEP instance via REST API:

```
# Database import (standalone JAR)
java -cp dumont-db-indexer.jar com.viglet.dumont.connector.db.DumDbImportTool \
  --server http://localhost:30130 \
  --api-key <API_KEY> \
  --driver org.mariadb.jdbc.Driver \
  --connect "jdbc:mariadb://localhost:3306/products" \
  --query "SELECT id, name, description, price FROM products" \
  --site ProductCatalog \
  --locale en_US

# FileSystem import (standalone JAR)
java -cp dumont-filesystem-indexer.jar \
  com.viglet.dumont.filesystem.DumFSImportTool \
  --source-dir /mnt/shared/documents \
  --server http://localhost:30130 \
  --api-key <API_KEY> \
  --site InternalDocs
```

These tools can be scheduled via cron jobs or CI/CD pipelines.

---

## Managing Connectors via the Turing ES Console

The AEM and Web Crawler connector plugins can be managed through the **Turing ES Admin Console**. To connect a running `dumont-connector.jar` instance to the Turing ES UI:

Open the Turing ES Admin Console

Navigate to **Enterprise Search → Integration**

Click **New** to create a new integration instance

Set the **Integration Type** (AEM or Web Crawler)

Set the **Endpoint** to the URL of your Dumont DEP connector instance (e.g.,

`http://localhost:30130`)

Enable the integration

Once connected, the Turing ES console provides a graphical interface for:

Configuring sources, content types, and field mappings

Triggering full indexing and re-indexing operations

Monitoring indexing progress in real time

Viewing indexing statistics and status

Running double-check consistency validation

For full details on the Integration UI – including monitoring, indexing stats, and double-check – see the [Turing ES Integration documentation](#).

For AEM-specific configuration (sources, content types, author/publish, delta tracking, locales, indexing rules) see [Turing ES AEM Connector documentation](#).

---

## Common Configuration Pattern

Every connector needs at least these pieces of information:

SETTING	DESCRIPTION
<b>Source</b>	Where to read content (URL, connection string, directory path, AEM endpoint)
<b>Credentials</b>	Authentication (username/password, API key, or none)
<b>Target SN Site</b>	The Turing ES Semantic Navigation Site that will receive the content
<b>Locale</b>	The language/country code for the content (e.g., <code>en_US</code> )
<b>Field Mapping</b>	How source fields map to search index fields

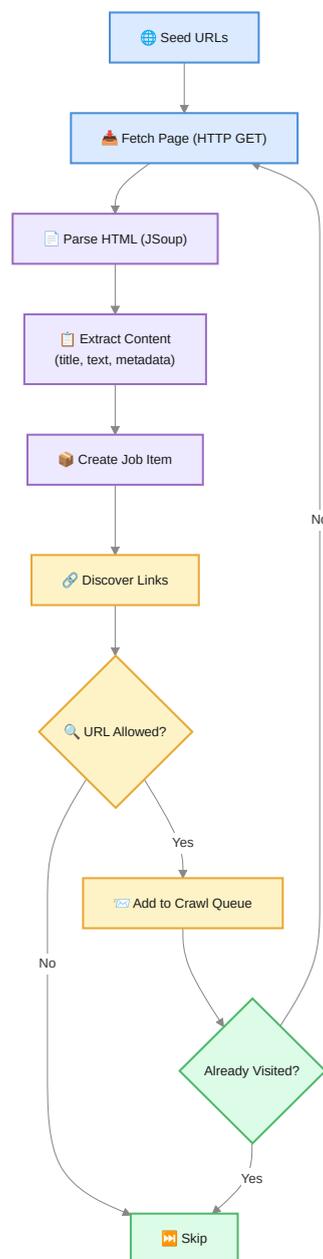
## How Each Connector Is Configured

CONNECTOR	CONFIGURATION METHOD	WHERE
<b>Web Crawler</b>	Turing ES Admin Console	<a href="#">Turing ES → Integration</a>
<b>AEM</b>	JSON file + Turing ES Admin Console	<code>export/</code> directory + <a href="#">Turing ES → Integration</a>
<b>Database</b>	CLI parameters	Command-line — see <a href="#">Database Connector</a>
<b>FileSystem</b>	CLI parameters	Command-line — see <a href="#">FileSystem Connector</a>
<b>WordPress</b>	WordPress Admin UI	WordPress → Settings → Viglet Dumont

# Web Crawler Connector

The Web Crawler connector recursively discovers and extracts content from websites. Starting from one or more seed URLs, it follows links, extracts page content using JSoup HTML parsing, and delivers each page as a searchable document.

## How It Works



Start with one or more **seed URLs** (starting points for the crawl)

Fetch each page via HTTP GET with a randomized user agent

Parse the HTML response with JSoup  
Extract content fields based on configured **attribute mappings**  
Create a Job Item and submit it to the pipeline  
Discover all links on the page  
Filter links against **allow/deny patterns** and file extension rules  
Add new, unvisited URLs to the crawl queue  
Repeat until no more URLs remain

---

## Key Features

FEATURE	DESCRIPTION
<b>Recursive crawling</b>	Follows links from seed URLs to discover all reachable pages
<b>URL filtering</b>	Allow and deny patterns (regex) control which URLs are crawled
<b>File extension filtering</b>	Configurable list of file extensions to include or exclude
<b>Authentication</b>	Basic HTTP authentication for protected sites
<b>Random user agents</b>	Generates random browser user-agent strings to avoid blocking
<b>Visited URL tracking</b>	CRC32 checksums prevent re-visiting the same URL in a single crawl
<b>Locale detection</b>	Extensible interface ( <code>DumWCExtLocaleInterface</code> ) for detecting content language
<b>URL normalization</b>	Normalizes URLs to avoid duplicate crawling of the same page
<b>Incremental indexing</b>	Checksum-based change detection — only re-indexes pages that changed

---

# Configuration

---

## Source Settings

FIELD	DESCRIPTION
<b>Name</b>	Identifier for this crawl source
<b>Starting URLs</b>	One or more seed URLs where the crawl begins
<b>Endpoint</b>	Base URL of the website
<b>Username / Password</b>	Optional HTTP Basic authentication credentials

## URL Filtering

FIELD	DESCRIPTION
<b>Allow URL Patterns</b>	Regex patterns — only URLs matching these patterns are crawled
<b>Deny URL Patterns</b>	Regex patterns — URLs matching these patterns are skipped
<b>Allowed File Extensions</b>	List of file extensions to include (e.g., <code>.html</code> , <code>.php</code> , <code>.aspx</code> )

## Attribute Mapping

Define which parts of the HTML page map to which fields in the search index:

FIELD	DESCRIPTION
<b>Title</b>	CSS selector or extraction rule for the page title
<b>Text</b>	CSS selector for the main content body
<b>URL</b>	The page URL (extracted automatically)
<b>Date</b>	CSS selector or meta tag for the publication date
<b>Custom fields</b>	Additional attribute mappings for any HTML element or meta tag

---

## Example: Crawling a Documentation Site

Crawl a documentation site starting from the homepage, only following links within the `/docs/` path:

SETTING	VALUE
Starting URL	<code>https://docs.example.com/</code>
Allow Pattern	<code>https://docs\example\.com/docs/.*</code>
Deny Pattern	<code>.*\.(css js png jpg gif svg)\$</code>
Target SN Site	<code>Documentation</code>
Locale	<code>en_US</code>

The crawler will:

- Start at the homepage
- Follow all links matching `/docs/` paths
- Skip links to stylesheets, scripts, and images
- Extract each page's title, body text, and URL
- Send each page to the `Documentation` SN Site in Turing ES

---

## Locale Detection

The Web Crawler supports automatic locale detection via the `DumWCExtLocaleInterface` extension point. Implement this interface to determine the locale of each page based on:

- URL path patterns (e.g., `/en/`, `/pt-br/`)
- HTML `lang` attribute
- Content-Language HTTP header
- Custom logic specific to your site structure

If no locale extension is configured, the default locale from the source configuration is used for all pages.

## Limitations

---

The crawler follows **HTML links only** — it does not execute JavaScript. Single-page applications (SPAs) that render content via JavaScript are not supported without a pre-rendering solution.

**robots.txt** — The crawler does not currently enforce robots.txt directives. Ensure you have permission to crawl the target site.

**Rate limiting** — There is no built-in rate limiter. For large sites, monitor server load and consider adding delays between requests.

---

## Related Pages

---

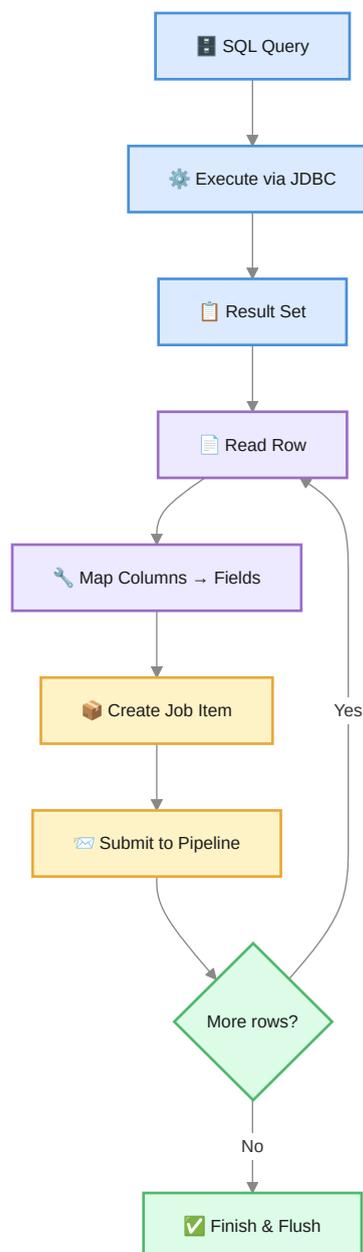
PAGE	DESCRIPTION
<a href="#">Connectors Overview</a>	All available connectors
<a href="#">Core Concepts</a>	Pipeline, strategies, and change detection
<a href="#">Turing ES — Integration</a>	How Turing ES receives content from connectors

---

# Database Connector

The Database Connector extracts content from any JDBC-compatible database by executing SQL queries and mapping each result row to a searchable document. It supports Oracle, PostgreSQL, MariaDB, MySQL, and any database with a JDBC driver.

## How It Works



Execute the configured **SQL query** against the database

Iterate through the result set

Map each **column** to a search index **field** based on the attribute mapping

Create a Job Item with the mapped fields

Submit to the pipeline in configurable **chunks**

Repeat until all rows are processed

---

## Key Features

FEATURE	DESCRIPTION
<b>Any JDBC database</b>	Oracle, PostgreSQL, MariaDB, MySQL, and any JDBC-compatible source
<b>SQL flexibility</b>	Use any SELECT query — joins, aggregations, subqueries, functions
<b>External SQL files</b>	Load queries from files using the <code>file://</code> protocol
<b>Batch processing</b>	Configurable chunk size for memory-efficient processing
<b>Max content size</b>	Limits per-document content size (default: 5 MB) to prevent oversized documents
<b>Standalone CLI</b>	Run imports from the command line without the full Dumont DEP application
<b>Custom extensions</b>	<code>DumDbExtCustomImpl</code> interface for custom row processing logic
<b>Locale support</b>	Configurable default locale for all extracted documents

---

## CLI Parameters

The standalone Database Connector accepts the following command-line parameters:

PARAMETER	REQUIRED	DEFAULT	DESCRIPTION
<code>--driver</code>	Yes	—	JDBC driver class name
<code>--connect</code>	Yes	—	JDBC connection string
<code>--query</code> / <code>-q</code>	Yes	—	SQL query or <code>file://path/to/query.sql</code>
<code>--site</code>	Yes	—	Target Semantic Navigation Site name
<code>--server</code> / <code>-s</code>	Yes	—	Dumont DEP server URL
<code>--api-key</code> / <code>-a</code>	Yes	—	API key for authentication
<code>--locale</code>	No	<code>en_US</code>	Default locale for documents
<code>--chunk</code> / <code>-z</code>	No	<i>(varies)</i>	Batch size for processing
<code>--max-content-size</code>	No	<code>5</code>	Maximum content size in MB
<code>--deindex-before-importing</code>	No	<code>false</code>	Remove all existing documents before import

## Example: Indexing a Product Catalog

### SQL Query

```
SELECT
  p.id,
  p.name AS title,
  p.description AS text,
  p.price,
  c.name AS category,
  p.updated_at AS date,
  CONCAT('https://shop.example.com/products/', p.slug) AS url
FROM products p
JOIN categories c ON p.category_id = c.id
WHERE p.active = 1
```

### Standalone Import

```
java -cp dumont-db.jar com.viglet.dumont.db.DumDbImportTool \
  --server http://localhost:30130 \
  --api-key your-api-key \
  --driver org.mariadb.jdbc.Driver \
  --connect "jdbc:mariadb://localhost:3306/shop?user=reader&password=secret" \
  --query "file:///opt/queries/products.sql" \
  --site ProductCatalog \
  --locale en_US \
  --chunk 100
```

This will:

Connect to the MariaDB database

Execute the SQL query from the external file

Map each row to a document (columns become fields)

Send documents in batches of 100 to the `ProductCatalog` SN Site

## Supported Databases

DATABASE	DRIVER	CONNECTION STRING EXAMPLE
MariaDB	<code>org.mariadb.jdbc.Driver</code>	<code>jdbc:mariadb://host:3306/dbname</code>
MySQL	<code>com.mysql.cj.jdbc.Driver</code>	<code>jdbc:mysql://host:3306/dbname</code>
PostgreSQL	<code>org.postgresql.Driver</code>	<code>jdbc:postgresql://host:5432/dbname</code>
Oracle	<code>oracle.jdbc.OracleDriver</code>	<code>jdbc:oracle:thin:@host:1521/service</code>

## Tips

**Use aliases** in your SQL query to match the field names expected by Turing ES (e.g., `SELECT name AS title`).

**External SQL files** (`file://` protocol) keep complex queries maintainable and version-controlled.

**De-index before importing** (`--deindex-before-importing`) is useful for full refreshes — it removes stale documents that no longer exist in the database.

**Chunk size** affects memory usage: larger chunks are faster but use more memory. Start with 100 and adjust based on document size.

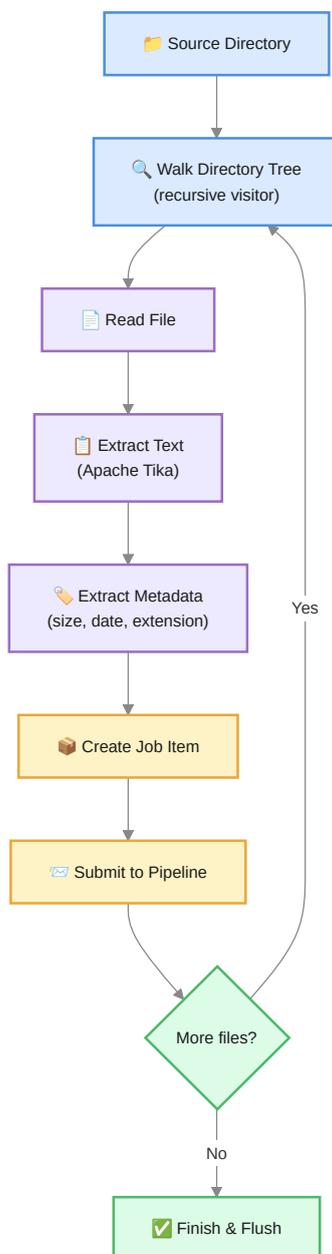
### CUSTOMIZING THE DATABASE CONNECTOR

Need to transform or enrich rows during import? See [Extending the Database Connector](#) for the `DumDbExtCustomImpl` interface and step-by-step guide.

# FileSystem Connector

The FileSystem Connector walks a directory tree, extracts text and metadata from files using Apache Tika, and indexes everything as searchable documents. It supports PDFs, Word documents, spreadsheets, presentations, plain text, HTML, and images (with OCR).

## How It Works



Start at the configured **source directory**

Recursively traverse all subdirectories

- For each file, extract text content using **Apache Tika**
  - Collect file metadata (size, modification date, extension, MIME type)
  - Create a Job Item with the extracted content and metadata
  - Submit to the pipeline in configurable batches
  - Repeat until all files are processed
- 

## Key Features

FEATURE	DESCRIPTION
<b>Recursive traversal</b>	Walks the full directory tree using Java's file visitor pattern
<b>Apache Tika integration</b>	Text extraction from 1000+ file formats
<b>OCR support</b>	Extract text from images and scanned PDFs (requires Tesseract)
<b>File metadata</b>	Captures file size, extension, modification date, and MIME type
<b>Prefix replacement</b>	Replace file path prefixes with custom URLs (e.g., replace local path with a web URL)
<b>Standalone CLI</b>	Run imports from the command line independently
<b>Configurable batch size</b>	Control memory usage with chunk-based processing

---

## Supported File Formats

Apache Tika supports text extraction from a broad range of formats:

CATEGORY	FORMATS
<b>Documents</b>	PDF, DOCX, DOC, ODT, RTF, EPUB
<b>Spreadsheets</b>	XLSX, XLS, ODS, CSV
<b>Presentations</b>	PPTX, PPT, ODP
<b>Web / Markup</b>	HTML, XHTML, XML
<b>Plain text</b>	TXT, LOG, Markdown
<b>Email</b>	EML, MSG, MBOX
<b>Images (with OCR)</b>	PNG, JPEG, TIFF, BMP, GIF

Files that Tika cannot extract text from are silently skipped.

---

## CLI Parameters

PARAMETER	REQUIRED	DEFAULT	DESCRIPTION
<code>--source-dir</code> / <code>-d</code>	Yes	—	Root directory to scan
<code>--server</code> / <code>-s</code>	Yes	—	Dumont DEP server URL
<code>--api-key</code> / <code>-a</code>	Yes	—	API key for authentication
<code>--site</code>	Yes	—	Target Semantic Navigation Site name
<code>--type</code> / <code>-t</code>	No	<code>Static</code> <code>File</code>	Content type label for all documents
<code>--locale</code>	No	<code>en_US</code>	Default locale
<code>--chunk</code> / <code>-z</code>	No	<code>100</code>	Batch size
<code>--file-size-field</code>	No	—	Field name to store file size
<code>--file-extension-field</code>	No	—	Field name to store file extension
<code>--prefix-from-replace</code>	No	—	Path prefix to replace (e.g., <code>/mnt/docs</code> )
<code>--prefix-to-replace</code>	No	—	Replacement prefix (e.g., <code>https://docs.example.com</code> )

## Example: Indexing a Document Repository

```
java -cp dumont-fs.jar com.viglet.dumont.filesystem.DumFSImportTool \  
  --source-dir /mnt/shared/documents \  
  --server http://localhost:30130 \  
  --api-key your-api-key \  
  --site InternalDocs \  
  --locale en_US \  
  --chunk 50 \  
  --file-size-field fileSize \  
  --file-extension-field fileExtension \  
  --prefix-from-replace /mnt/shared/documents \  
  --prefix-to-replace https://intranet.example.com/docs
```

This will:

- Scan all files under `/mnt/shared/documents` recursively
- Extract text from each file using Apache Tika
- Replace the local path with a web URL for each document
- Index everything into the `InternalDocs` SN Site
- Store file size and extension as searchable/facetable fields

## Path Prefix Replacement

The `--prefix-from-replace` and `--prefix-to-replace` parameters transform file paths into web-accessible URLs. This is essential when the files are served by a web server:

LOCAL PATH	AFTER REPLACEMENT
<code>/mnt/shared/documents/reports/q1-2026.pdf</code>	<code>https://intranet.example.com/docs/reports/q1-2026.pdf</code>
<code>/mnt/shared/documents/policies/security.docx</code>	<code>https://intranet.example.com/docs/policies/security.docx</code>

The transformed path becomes the document's **URL** field in the search index, allowing users to click through from search results to the actual file.

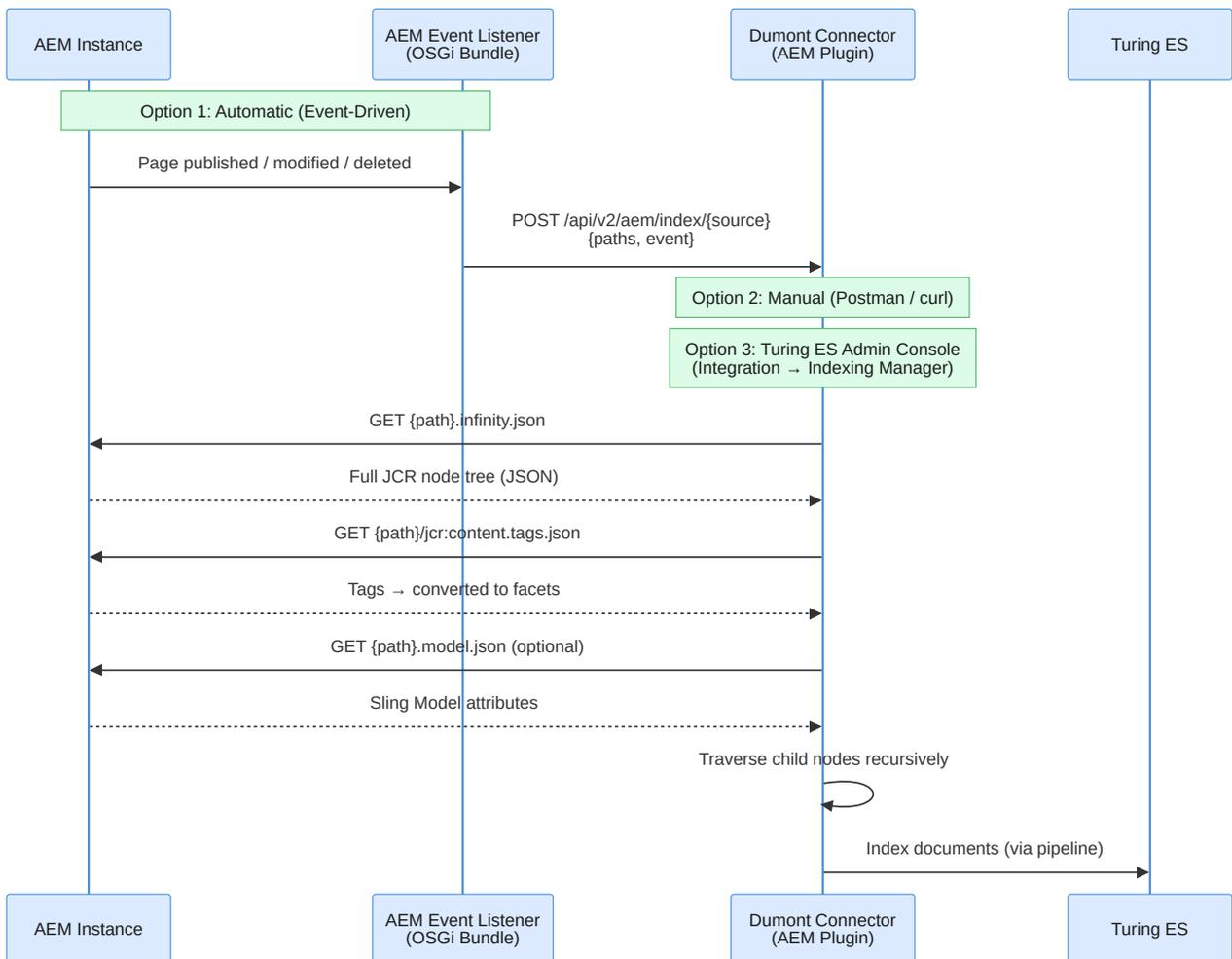
---

# AEM Connector

The AEM Connector indexes content from Adobe Experience Manager (AEM) author and publish instances. It consists of two components: an **AEM server-side bundle** (OSGi event listeners installed inside AEM) and a **connector plugin** (Java JAR loaded into `dumont-connector.jar`).

## How It Works

The AEM connector receives indexing requests, then accesses AEM to traverse the content tree, fetch page data, extract tags as facets, and optionally call `.model.json` for custom attributes.



## Three Ways to Trigger Indexing

METHOD	HOW	WHEN TO USE
<b>AEM Event Listeners</b>	Install the <code>aem-server</code> OSGi bundle inside AEM — it automatically sends indexing requests when content is published, modified, or deleted	Production — real-time content sync
<b>Manual API Call</b>	Send a POST request to <code>/api/v2/aem/index/{source}</code> with a JSON payload containing paths and event type	Development, testing, one-off re-indexing
<b>Turing ES Admin Console</b>	Use <b>Enterprise Search → Integration → Indexing Manager</b> to select paths and trigger indexing/deindexing/publishing operations	Operations — selective re-indexing via UI

## Content Discovery Strategies

The AEM connector supports two strategies for discovering content during a full **Index All** operation:

STRATEGY	PROPERTY	HOW IT DISCOVERS CONTENT
<b>Tree Traversal</b> <i>(default)</i>	<code>dumont.aem.query</code> <code>builder=false</code>	Recursively walks the content tree from the root path using <code>infinity.json</code> — follows parent→child relationships
<b>QueryBuilder</b>	<code>dumont.aem.query</code> <code>builder=true</code>	Uses AEM's native QueryBuilder API ( <code>/bin/querybuilder.json</code> ) to find all content matching the configured content type in paginated batches

### QueryBuilder Discovery

When enabled, the connector queries the QueryBuilder endpoint instead of walking the tree:

```
GET http://localhost:4502/bin/querybuilder.json?
path=/content/wknd&type=cq:Page&p.hits=slim&p.limit=500&p.offset=0
```

Each page of discovered paths is immediately processed in parallel using a configurable thread pool – the full path list is never held in memory.

**Enable it** by adding the following properties to `application.yaml` or as JVM arguments:

```
dumont:  
  aem.querybuilder: true  
  aem.querybuilder.parallelism: 10 # number of parallel threads (default)
```

Or via command-line:

```
java -Ddumont.aem.querybuilder=true -Ddumont.aem.querybuilder.parallelism=10 -  
jar dumont-connector.jar
```

### CONTENT TYPE REQUIRED

QueryBuilder discovery requires a **Content Type** (e.g., `cq:Page`) configured in the AEM source. Without it, the command logs a warning and skips processing.

### WHEN TO USE QUERYBUILDER

QueryBuilder is recommended for **large content trees** where recursive traversal is slow. It reduces the number of HTTP round-trips to AEM by discovering all paths in bulk, then fetching content in parallel. For small sites (< 1 000 pages), the default tree traversal is typically sufficient.

## The Indexing Flow (Step by Step)

When the connector receives an indexing request (from any of the three triggers), it processes each path as follows:

### 1. Fetch the Content Node

The connector calls AEM's `infinity.json` endpoint to get the full JCR node tree:

```
GET http://localhost:4502/content/wknd/us/en/my-page.infinity.json
```

This returns the complete node hierarchy as JSON — all properties, child nodes, and metadata. The connector filters out internal nodes (prefixed with `jcr:`, `rep:`, `cq:`).

## 2. Extract Tags as Facets

For each page, the connector fetches tags:

```
GET http://localhost:4502/content/wknd/us/en/my-page/jcr:content.tags.json
```

Tags are **automatically converted to facets** in the search index — no manual configuration needed. Each tag becomes a filterable value in the Turing ES facet panel.

## 3. Fetch Model JSON (Optional — Requires Configuration)

The connector does **not** call `.model.json` by default. To enable it, you must configure a `DumAemExtContentInterface` implementation in the `models` section of the export JSON:

```
"models": [  
  {  
    "type": "cq:Page",  
    "className":  
    "com.viglet.dumont.connector.aem.sample.ext.DumAemExtSampleModelJson",  
    "targetAttrs": [ ... ]  
  }  
]
```

When this is configured, the extension class fetches the Sling Model exporter:

```
GET http://localhost:4502/content/wknd/us/en/my-page.model.json
```

This returns structured content from AEM's Sling Models — useful for extracting custom attributes like content fragment paths, component data, or experience fragment references. See [Extending the AEM Connector](#) for how to implement `DumAemExtContentInterface` and the full JSON configuration reference.

## 4. Traverse the Content Tree

Starting from the configured **root path** (e.g., `/content/wknd`), the connector recursively traverses all child nodes that match the configured **content type** (e.g., `cq:Page`). Each matching node goes through steps 1–3 above.

## 5. Map Attributes and Index

For each page, the connector:

- Applies **attribute mappings** from the configuration (global attributes + model-specific source→target mappings)

- Runs **custom extension classes** (if configured via `className`)

- Creates a **Job Item** for both author and publish environments (if enabled)

- Sends the Job Item through the Dumont DEP pipeline to Turing ES

---

## AEM Server-Side Bundle (Event Listeners)

The `aem-server` module is an **OSGi bundle installed inside AEM**. It provides event listeners that automatically notify the Dumont connector when content changes.

### Events Captured

EVENT LISTENER	AEM EVENT	DUMONT ACTION
<code>DumAemPageEventHandler</code>	Page created / modified	<b>INDEXING</b>
<code>DumAemPageReplicationEventHandler</code>	Page activated (published)	<b>PUBLISHING</b>
<code>DumAemPageReplicationEventHandler</code>	Page deactivated (unpublished)	<b>UNPUBLISHING</b>
<code>DumAemResourceEventHandler</code>	DAM asset created / modified	<b>INDEXING</b>

### OSGi Configuration

The event listeners are configured in AEM via **OSGi Configuration** (AEM → Web Console → Configuration):

SETTING	DESCRIPTION
<b>Enabled</b>	Toggle to enable/disable automatic indexing
<b>Host</b>	Dumont connector URL (e.g., <code>http://dumont-server:30130</code> )
<b>Config Name</b>	Source name configured in the Dumont connector

### HTTP Payload

When an event fires, the bundle sends:

```
POST http://dumont-server:30130/api/v2/aem/index/{configName}
Content-Type: application/json

{
  "paths": ["/content/wknd/us/en/my-page"],
  "event": "INDEXING"
}
```

Event types: `INDEXING`, `DEINDEXING`, `PUBLISHING`, `UNPUBLISHING`.

---

## Manual API Triggering

You can trigger indexing manually via HTTP (Postman, curl, etc.):

### Index Specific Paths

```
curl -X POST http://localhost:30130/api/v2/aem/index/WKND \
-H "Content-Type: application/json" \
-d '{
  "paths": ["/content/wknd/us/en/about"],
  "event": "INDEXING",
  "recursive": true
}'
```

### Deindex Specific Paths

```
curl -X POST http://localhost:30130/api/v2/aem/index/WKND \
-H "Content-Type: application/json" \
-d '{
  "paths": ["/content/wknd/us/en/old-page"],
  "event": "DEINDEXING"
}'
```

## Request Body Fields

FIELD	TYPE	DEFAULT	DESCRIPTION
<code>paths</code>	string[]	<i>(required)</i>	AEM content paths to process
<code>event</code>	string	<code>INDEXING</code>	<code>INDEXING</code> , <code>DEINDEXING</code> , <code>PUBLISHING</code> , or <code>UNPUBLISHING</code>
<code>recursive</code>	boolean	<code>false</code>	Traverse child nodes recursively
<code>attribute</code>	string	<code>ID</code>	<code>ID</code> (path-based) or <code>URL</code> (URL-based)

## Source Configuration

Each AEM source defines connection details, content scope, author/publish environments, locale mappings, and delta tracking. Sources are configured in the **Turing ES Admin Console** under **Enterprise Search → Integration → [your AEM instance] → Sources**.

For the JSON configuration file used by custom extensions (attributes, models, locale paths), see [Extending the AEM Connector](#).

### General

FIELD	DESCRIPTION
Name	Source identifier
Endpoint	URL of the AEM instance (e.g., <code>http://localhost:4502</code> )
Username / Password	Credentials for authenticated access to the AEM instance

### Root Path

Defines the root content path within the AEM repository from which content is traversed (e.g., `/content/wknd`). All child nodes matching the configured content type are indexed recursively from this path.

### Content Types

FIELD	DESCRIPTION
Content Type	Primary content type to be indexed (e.g., <code>cq:Page</code> )
Sub Type	Optional sub-type filter within the content type

### Delta Tracking

Controls incremental indexing — how the connector detects which content has changed since the last run.

FIELD	DESCRIPTION
Once Pattern	Pattern used to identify content that should only be indexed once
Delta Class	Fully-qualified Java class name responsible for detecting changed content since the last run (see <a href="#">Extending AEM</a> for custom implementations)

## Author / Publish

Configures which AEM environments are indexed and how they map to Turing ES Semantic Navigation Sites.

FIELD	DESCRIPTION
Author	Enable indexing from the AEM author environment
Publish	Enable indexing from the AEM publish environment
SN Site (Author)	Semantic Navigation Site that receives author content
SN Site (Publish)	Semantic Navigation Site that receives publish content
URL Prefix (Author)	URL prefix prepended to document paths in the author index
URL Prefix (Publish)	URL prefix prepended to document paths in the publish index

## Locales

Maps content language codes to repository paths.

FIELD	DESCRIPTION
Default Locale	Locale used when no language-specific path is matched
Locale Class	Fully-qualified Java class name responsible for resolving document locale (see <a href="#">Extending AEM</a> for custom implementations)
Locale → Path	Dynamic list mapping each locale code (e.g., <code>en_US</code> ) to its root path in the repository

## Source Actions

Each source has two action buttons available in the Turing ES admin console:

**Index All** — triggers a full indexing run for all content in this source

**Reindex All** — forces a full reindexation, replacing all previously indexed content

---

## Indexing Rules

Indexing Rules allow you to filter content during indexing — for example, to exclude error pages or draft content before it reaches the search index. Rules are configured in the **Turing ES Admin Console** under **Enterprise Search → Integration → [your AEM instance] → Indexing Rules**.

FIELD	DESCRIPTION
Name	Rule identifier (required)
Description	Purpose of this rule
Source	The source this rule applies to
Attribute	Document field to evaluate (e.g., <code>template</code> )
Rule Type	How the rule is applied — currently supports <b>IGNORE</b> (skip documents that match)
Values	Dynamic list of values that trigger the rule (add or remove entries)

**Example:** A rule with `Attribute = template`, `Rule Type = IGNORE`, and `Values = [error-page]` will prevent any document with `template:error-page` from being indexed.

## Indexing Manager

The Indexing Manager provides a stepper form in the **Turing ES Admin Console** for targeting specific documents for manual operations.

OPERATION	DESCRIPTION	COLOUR
INDEXING	Index specific content	Blue
DEINDEXING	Remove specific content from the index	Red
PUBLISHING	Publish content	Green
UNPUBLISHING	Unpublish content	Orange

Each operation step allows you to:

Select the **Source** to operate on

Choose the **attribute** to identify documents: **ID** or **URL**

Enter one or more specific values (IDs or URLs)

Expand **Advanced Settings** to toggle **Recursive** mode, which traverses child content in hierarchical repositories

## Concurrency

The connector supports two execution modes:

MODE	WHEN	BEHAVIOR
<b>Exclusive</b>	Full crawl ( <code>indexAll</code> )	Only one full crawl per source at a time
<b>Standalone</b>	Specific paths (event-driven / manual)	Multiple concurrent updates allowed

Reactive (parallel) processing can be enabled for large sites:

```
dumont.reactive.indexing=true
dumont.reactive.parallelism=10
```

When using **QueryBuilder discovery**, parallelism is controlled separately via

`dumont.aem.querybuilder.parallelism` (default `10`). See [Content Discovery Strategies](#) above.

## CUSTOMIZING THE AEM CONNECTOR

Need custom attribute extractors, delta date logic, or content processors? See [Extending the AEM Connector](#) for the full extension system, configuration JSON reference, and step-by-step guide.

## Related Pages

PAGE	DESCRIPTION
<a href="#">AEM Event Listener</a>	Install the OSGi event listener bundle inside AEM for real-time indexing
<a href="#">Extending the AEM Connector</a>	Custom attribute extractors, content processors, and configuration JSON reference
<a href="#">Turing ES – Integration</a>	General integration management – monitoring, indexing stats, and system information
<a href="#">Turing ES – AEM Connector</a>	AEM integration overview from the Turing ES perspective
<a href="#">Turing ES – Semantic Navigation</a>	Configure the SN Sites that receive indexed content

# AEM Event Listener Setup

This guide explains how to install the Dumont AEM event listener bundle inside your AEM as a Cloud Service instance. Once installed, AEM automatically notifies the Dumont connector whenever content is published, modified, or deleted – enabling real-time search index synchronization.

## What Gets Installed

The Dumont AEM Server module (`aem-server`) deploys three artifacts into AEM:

ARTIFACT	TYPE	PURPOSE
<code>aem-dumont.structure</code>	Content package	Creates the <code>/apps/turing</code> repository structure
<code>aem-dumont.config</code>	Content package	Deploys the OSGi configuration (host, configName, enabled)
<code>aem-dumont.core</code>	OSGi bundle	Event handlers that listen for page/asset changes and send HTTP requests to Dumont

## Prerequisites

An AEM as a Cloud Service environment (dev, stage, or prod)

A running Dumont connector instance with the AEM plugin (accessible from AEM)

Access to your AEM project's Git repository (for Cloud Manager pipeline deployment)

Java 11+ and Maven 3.6+ for local builds

### **i** AEM AS A CLOUD SERVICE DEPLOYMENT

AEM as a Cloud Service does not support direct bundle installation via the Felix Console. All code must be deployed through a **Cloud Manager pipeline** from a Git repository. See [Adobe's deployment documentation](#) for details.

## Option 1: Add to an Existing AEM Project

If you already have an AEM as a Cloud Service project (created from the [AEM Project Archetype](#)), you can integrate the Dumont event listeners into it.

### Step 1 – Copy the Source Files

From the [Dumont repository](#), copy these directories into your AEM project:

```
From: dumont/aem/aem-server/core/src/main/java/com/viglet/turing/aem/server/  
├── config/  
│   └── DumAemIndexerConfig.java  
└── core/  
    ├── events/  
    │   ├── DumAemPageEventHandler.java  
    │   ├── DumAemPageReplicationEventHandler.java  
    │   ├── DumAemResourceEventHandler.java  
    │   └── beans/  
    │       ├── DumAemEvent.java  
    │       └── DumAemPayload.java  
    │   └── utils/  
    │       └── DumAemEventUtils.java  
    └── services/  
        ├── DumAemIndexerService.java  
        └── DumAemIndexerServiceImpl.java
```

```
To: your-aem-project/core/src/main/java/com/viglet/turing/aem/server/
```

### Step 2 – Add Dependencies to Your Core Module

Add these dependencies to your project's `core/pom.xml`:

```
<!-- Lombok (if not already present) -->
<dependency>
  <groupId>org.projectlombok</groupId>
  <artifactId>lombok</artifactId>
  <version>1.18.34</version>
  <scope>provided</scope>
</dependency>

<!-- Apache Commons Collections -->
<dependency>
  <groupId>org.apache.commons</groupId>
  <artifactId>commons-collections4</artifactId>
  <version>4.5.0-M3</version>
  <scope>provided</scope>
</dependency>
```

The following dependencies are already provided by AEM SDK and don't need to be added: Jackson, Apache HTTP Client, Sling API, OSGi annotations, AEM Replication API, AEM WCM API.

### Step 3 — Add the OSGi Configuration

Create the configuration file for your AEM environment:

```
your-aem-project/ui.config/src/main/content/jcr_root/apps/<your-
project>/osgiconfig/config/
└─
com.viglet.turing.aem.server.core.services.TurAemIndexerServiceImpl.cfg.json
```

Contents:

```
{
  "enabled": true,
  "host": "${env:DUMONT_CONNECTOR_HOST;default=http://localhost:30130}",
  "configName": "${env:DUMONT_CONFIG_NAME;default=WKND}"
}
```

PROPERTY	DESCRIPTION
<code>enabled</code>	<code>true</code> to activate event-driven indexing, <code>false</code> to disable
<code>host</code>	URL of the Dumont connector instance (with AEM plugin)
<code>configName</code>	AEM source name configured in the Dumont connector

### CLOUD MANAGER ENVIRONMENT VARIABLES

Use `[$[env:VAR_NAME;default=value]]` syntax to configure different values per environment (dev, stage, prod) without changing code. Set the variables in **Cloud Manager → Environments → Environment Variables**.

See [Adobe's environment variables documentation](#).

## Step 4 — Configure the Service User

The replication event handler requires a service user mapping. Add this to your project's `ui.config`:

```
your-aem-project/ui.config/src/main/content/jcr_root/apps/<your-project>/osgiconfig/config/
└─ org.apache.sling.serviceusermapping.impl.ServiceUserMapperImpl.amended-dumont-indexer.cfg.json
```

```
{
  "user.mapping": [
    "aem-dumont.core:dumont-indexer=[content-reader-service]"
  ]
}
```

## Step 5 — Deploy via Cloud Manager

Commit and push your changes to the Git repository

In **Cloud Manager**, run the deployment pipeline for your target environment

After deployment, verify the bundle is active in AEM → **Tools → Operations → Web Console → Bundles** (search for "dumont")



## Option 2: Create a New AEM Project from Archetype

If you don't have an existing AEM project, create one from the [AEM Project Archetype](#):

```
mvn -B org.apache.maven.plugins:maven-archetype-plugin:3.2.1:generate \
  -D archetypeGroupId=com.adobe.aem \
  -D archetypeArtifactId=aem-project-archetype \
  -D archetypeVersion=49 \
  -D appTitle="My Site" \
  -D appId="mysite" \
  -D groupId="com.mycompany" \
  -D aemVersion="cloud"
```

Then follow Steps 1–5 from Option 1 above.

For the full archetype reference, see [Adobe's archetype documentation](#).

---

## Option 3: Use the Standalone Dumont AEM Server Module

You can also build and deploy the `aem-server` module directly from the Dumont repository:

```
git clone https://github.com/opencviglet/dumont.git
cd dumont/aem/aem-server
mvn clean install -PautoInstallPackage -Daem.host=localhost -Daem.port=4502
```

### LOCAL SDK ONLY

The `-PautoInstallPackage` profile installs directly to a local AEM SDK instance. For AEM as a Cloud Service, you must deploy through a Cloud Manager pipeline (Option 1 or 2).

## OSGi Configuration Properties

After deployment, the configuration can be managed in **AEM → Tools → Operations → Web Console → Configuration** (search for "Dumont"):

PROPERTY	DEFAULT	DESCRIPTION
Enabled	false	Toggle event-driven indexing on/off
Host	http://localhost:30130	Dumont connector URL
Config Name	WKND	Source name in the Dumont connector

## Cloud Manager Environment Variables

For AEM as a Cloud Service, set these variables in Cloud Manager:

VARIABLE	ENVIRONMENT	VALUE
DUMONT_CONNECTOR_HOST	Dev	http://dev-dumont:30130
DUMONT_CONNECTOR_HOST	Stage	http://stage-dumont:30130
DUMONT_CONNECTOR_HOST	Prod	http://prod-dumont:30130
DUMONT_CONFIG_NAME	All	WKND (or your source name)

## Verifying the Installation

**Check bundle status:** AEM → Tools → Operations → Web Console → Bundles → search "dumont" — should be **Active**

**Check configuration:** Web Console → Configuration → search "TurAemIndexer" — verify host and configName

**Test:** Publish a page in AEM and check the Dumont connector logs for incoming requests

**Monitor:** In Turing ES → Integration → Monitoring, filter by your AEM source to see indexing events

---

## Related Pages

---

PAGE	DESCRIPTION
<a href="#">AEM Connector</a>	How the connector processes AEM content (infinity.json, tags, model.json)
<a href="#">Extending the AEM Connector</a>	Custom extensions, configuration JSON, and Maven dependencies
<a href="#">Installation Guide</a>	Deploying Dumont connector with the AEM plugin
<a href="#">Turing ES – AEM Integration</a>	Managing AEM indexing via the Turing ES admin console

# WordPress Connector

The WordPress Connector is a **PHP plugin installed directly inside WordPress** — not a Java connector plugin. It hooks into WordPress publish, update, and delete events to keep the search index synchronized with your content in real time.

## **i** THIS IS A WORDPRESS PLUGIN, NOT A JAVA PLUGIN

Unlike the AEM and Web Crawler connectors (which are Java JARs loaded via `-Dloader.path`), the WordPress Connector is a **PHP plugin** that you install in your WordPress `wp-content/plugins/` directory. It communicates directly with Turing ES via HTTP — no `dumont-connector.jar` needed.

Unlike the AEM and Web Crawler connectors, the WordPress plugin is **not part of the Dumont connector application** — it runs entirely inside WordPress and sends content directly to Turing ES via HTTP. No `dumont-connector.jar` is needed.

## How It Works



The plugin integrates with WordPress hooks to automatically push content to the search index:

**On publish** — When a post or page is published, the plugin extracts its content and sends it to Turing ES

**On update** — When published content is modified, the index is updated automatically

**On delete** — When content is trashed or deleted, it is removed from the index

**On status change** — When a published post is changed to draft or private, it is removed from the index

**Manual bulk indexing** — The admin panel provides buttons to index all posts or all pages in batches of 250

There is no scheduled or cron-based indexing — everything is **event-driven** or manually triggered.

## Installation

Copy the `viglet-turing-for-wordpress` plugin folder to your WordPress installation:

```
wp-content/plugins/viglet-turing-for-wordpress/
```

Activate the plugin in **WordPress Admin → Plugins**

Configure the connection in **Settings → Viglet Dumont**

## Configuration

### Server Settings

FIELD	DEFAULT	DESCRIPTION
Host	localhost	Turing ES server hostname
Port	2700	Turing ES server port
Path	/dumont	Base path for the indexing endpoint
Site Name	<i>(required)</i>	Semantic Navigation Site name in Turing ES

### Indexing Options

FIELD	DESCRIPTION
Post types	Which content types to index (posts, pages, custom post types)
Remove on delete	Automatically remove from index when content is deleted
Remove on status change	Remove from index when status changes from published to draft/private
Index comments	Include post comments in the indexed content
Custom fields	Additional WordPress custom fields to include in the index
Excluded IDs	Post IDs to exclude from indexing

## Search Display Options

FIELD	DESCRIPTION
<b>Results per page</b>	Number of search results to display
<b>Show facets</b>	Enable faceted navigation on the search results page
<b>Facet fields</b>	Which fields to show as facets (categories, tags, author, type)
<b>Spellcheck</b>	Enable spelling suggestions for search queries
<b>Max tags</b>	Maximum number of tags to display in facets

---

## Indexed Fields

---

The plugin extracts and sends these fields to Turing ES:

FIELD	SOURCE
<b>id</b>	WordPress post ID
<b>title</b>	Post/page title
<b>content</b>	Full post content
<b>contentnoshortcodes</b>	Content with WordPress shortcodes removed
<b>permalink</b>	Canonical URL
<b>author</b>	Author display name
<b>type</b>	Content type (post, page, custom)
<b>date</b>	Publication date
<b>modified</b>	Last modification date
<b>categories</b>	Category names (multi-valued)
<b>tags</b>	Tag names (multi-valued)
<b>numcomments</b>	Comment count
<b>comments</b>	Comment text (if comment indexing is enabled)
<b>Custom fields</b>	Any configured custom fields

## Multisite Support

The plugin supports WordPress multisite installations. When activated network-wide, it can index content from all blogs in the network, including `blogid`, `blogdomain`, and `blogpath` fields for each document.

## Admin Panel Actions

The plugin adds action buttons to the WordPress admin settings page:

ACTION	DESCRIPTION
<b>Ping</b>	Tests the connection to Turing ES
<b>Index all Posts</b>	Bulk indexes all published posts (batches of 250)
<b>Index all Pages</b>	Bulk indexes all published pages (batches of 250)
<b>Delete all</b>	Removes all WordPress content from the Turing ES index
<b>Optimize</b>	Triggers index optimization on the search engine

## Architecture

The plugin includes its own PHP HTTP client library (`DumontPhpClient`) for communicating with Turing ES:

```
viglet-turing-for-wordpress/  
├─ viglet-turing-for-wordpress.php      # Main plugin (hooks, indexing, search)  
├─ viglet-turing-options-page.php      # Admin settings page  
├─ DumontPhpClient/  
│   └─ Viglet/Dumont/  
│       ├─ Service.php                 # HTTP client for Turing ES  
│       ├─ Document.php                # Document model  
│       ├─ Response.php                # Response parser  
│       └─ HttpTransport/  
│           ├─ Curl.php                 # cURL transport  
│           ├─ CurlNoReuse.php         # cURL without connection reuse  
│           └─ FileGetContents.php     # file_get_contents fallback  
└─ template/  
    ├─ turing4wp_search.php            # Search results template  
    ├─ search.css                      # Search results styling  
    └─ autocomplete.css                # Autocomplete styling
```

# Developer Guide

---

Whether you're building custom extensions, contributing to the project, or integrating Dumont DEP into your CI/CD pipeline, this guide has everything you need.

Dumont DEP is fully open-source at [github.com/openviglet/dumont](https://github.com/openviglet/dumont). All contributions are welcome.

---

## Tech Stack

---

LAYER	TECHNOLOGY
Backend	Java 21 · Spring Boot 4.0.3
Message Queue	Apache Artemis (embedded)
Database	H2 (dev) · PostgreSQL (prod)
HTML Parsing	JSoup 1.22.1
Text Extraction	Apache Tika 3.2.3
Search Clients	Turing Java SDK · SolrJ · ES Client
Build	Apache Maven
CI/CD	GitHub Actions

---

## Setting Up Your Dev Environment

---

### Prerequisites

Java 21 (Temurin recommended)

Maven 3.9+

Git

Turing ES running at <http://localhost:2700> (for end-to-end testing)

## Clone, Build, and Run

```
git clone https://github.com/openviglet/dumont.git
cd dumont
mvn clean install

cd connector/connector-app
mvn spring-boot:run
```

The application starts at <http://localhost:30130>.

---

## Project Structure

```

dumont/
├── commons/                # Shared interfaces, models, utilities
├── aem-commons/          # AEM extension interfaces (published to Maven
Central)
├── spring/                # Spring Boot configuration and JPA
persistence
├── connector/
│   └── connector-app/    # Main pipeline – strategies, batch, queue,
indexing plugins, API
├── web-crawler/
│   └── wc-plugin/        # Web Crawler connector plugin
├── db/
│   └── db-commons/      # DB extension interface (published to Maven
Central)
│       ├── db-app/      # Database connector (standalone CLI)
│       └── db-sample/   # Example custom DB extension
├── filesystem/
│   └── fs-connector/    # FileSystem connector (standalone CLI)
├── aem/
│   ├── aem-plugin/     # AEM connector plugin (loaded via -
Dloader.path)
│   ├── aem-server/     # AEM server-side integration
│   └── aem-plugin-sample/ # Example custom AEM extensions (WKND site)
└── wordpress/          # WordPress PHP plugin

```

## Extension Points

Dumont DEP provides extension points at multiple levels:

## Connector-Level Extensions

EXTENSION	GUIDE	MAVEN ARTIFACT
<b>AEM extensions</b> — custom attribute extractors, content processors, delta date logic	Extending the AEM Connector	<code>com.viglet.dumont:aem-commons:2026.1.19</code>
<b>Database extensions</b> — custom row transformations during SQL import	Extending the Database Connector	<code>com.viglet.dumont:db-commons:2026.1.19</code>

## Platform-Level Extensions

### Creating a Custom Connector

Implement the `DumConnectorPlugin` interface:

```
public interface DumConnectorPlugin {
    void crawl();
    String getProviderName();
    void indexAll(String source);
    void indexById(String source, List<String> contentId);
}
```

### Creating a Custom Indexing Plugin

Implement the `DumIndexingPlugin` interface:

```
public interface DumIndexingPlugin {
    void index(TurSNJobItems turSNJobItems);
    String getProviderName();
}
```

Register your plugin as a Spring `@Component` with `@ConditionalOnProperty(name = "dumont.indexing.provider", havingValue = "your-provider")`.

## Processing Strategy Architecture

---

Strategies are evaluated in priority order for each Job Item:

To add a custom strategy, implement the strategy interface and assign a priority between the existing ones.

---

## REST API

---

ENDPOINT	DESCRIPTION
<code>GET /api/v2/connector/status</code>	Health check
<code>POST /api/v2/connector/indexing/</code>	Submit indexing jobs
<code>GET /api/v2/connector/monitoring/index/{source}</code>	Monitor indexing progress
<code>GET /api/v2/connector/validate/{source}</code>	Validate a content source

For the full API surface, start the application and visit the Swagger UI.

---

## Contributing

---

**Fork** the [openviglet/dumont](#) repository

**Create a branch** for your feature or fix: `git checkout -b feature/my-improvement`

**Commit** with clear, descriptive messages

**Open a Pull Request** — describe what you changed and why

---

# REST API Reference

---

Dumont DEP exposes a REST API for triggering indexing operations, monitoring progress, managing sources, and configuring indexing rules. All endpoints use JSON and support CORS.

---

## Connector API

---

Base path: `/api/v2/connector`

### Health Check

```
GET /api/v2/connector/status
```

Returns `{"status": "ok"}` if the service is running.

---

### Index All Content

```
GET /api/v2/connector/index/{name}/all
```

Triggers a full indexing operation for all content in the specified source.

PARAMETER	LOCATION	DESCRIPTION
<code>name</code>	Path	Source name

Response: `{"status": "sent"}`

---

## Index Specific Content

```
POST /api/v2/connector/index/{name}
```

Indexes specific documents by their IDs.

PARAMETER	LOCATION	DESCRIPTION
<code>name</code>	Path	Source name
<i>(body)</i>	JSON	Array of content IDs: <code>["id-1", "id-2"]</code>

**Response:** `{"status": "sent"}`

---

## Reindex All Content

```
GET /api/v2/connector/reindex/{name}/all
```

Deletes the existing index and reindexes all content from the source.

PARAMETER	LOCATION	DESCRIPTION
<code>name</code>	Path	Source name

**Response:** `{"status": "sent"}`

---

## Reindex Specific Content

```
GET /api/v2/connector/reindex/{name}
```

Deletes and reindexes specific documents.

PARAMETER	LOCATION	DESCRIPTION
<code>name</code>	Path	Source name
<i>(body)</i>	JSON	Array of content IDs

**Response:** `{"status": "sent"}`

---

## Validate Source

```
GET /api/v2/connector/validate/{source}
```

Validates content differences between the source and the search index. Used by the Turing ES Double Check feature.

PARAMETER	LOCATION	DESCRIPTION
<code>source</code>	Path	Source name

**Response:** `DumConnectorValidateDifference` — lists missing and extra documents.

---

## Monitoring API

---

Base path: `/api/v2/connector/monitoring`

### Get Indexing Status

```
GET /api/v2/connector/monitoring/index/{source}
```

Returns indexing status records for a specific source. Returns `404` if no records found.

---

### Get All Monitoring Data

```
GET /api/v2/connector/monitoring/indexing
```

Returns all indexing monitoring data across all sources.

---

### Get Monitoring by Source

```
GET /api/v2/connector/monitoring/indexing/{source}
```

Returns monitoring data filtered by source name.

---

### Search Monitoring Records

```
POST /api/v2/connector/monitoring/indexing/_search
```

Searches indexing monitoring records with pagination and filters.

**Request body:**

```
{
  "source": "WKND",
  "status": "INDEXED",
  "page": 0,
  "size": 20
}
```

## Get Monitoring by Content IDs

```
POST /api/v2/connector/monitoring/indexing
```

**Request body:** Array of content IDs.

## Indexing Statistics

```
GET /api/v2/connector/monitoring/indexing/stats
GET /api/v2/connector/monitoring/indexing/stats/{source}
```

Returns indexing operation statistics (start time, document count, duration, throughput). Returns **404** if no stats found.

## System Information API

```
GET /api/v2/connector/system-info
```

Returns runtime system information consumed by the Turing ES Integration → System Information page:

Application name and version

Java version, vendor, and JVM details

Operating system name and version

Physical memory (total, used, free)

JVM heap memory (total, used, free)

Disk space (total, used, free)

---

## Indexing Rules API

---

Base path: `/api/v2/connector/indexing-rule`

CRUD operations for managing regex-based indexing rules that filter content during extraction.

### List All Rules

```
GET /api/v2/connector/indexing-rule
```

### List Rules by Source

```
GET /api/v2/connector/indexing-rule/source/{source}
```

### Get Rule by ID

```
GET /api/v2/connector/indexing-rule/{id}
```

### Create Rule

```
POST /api/v2/connector/indexing-rule
```

**Request body:**

```
{  
  "name": "Exclude error pages",  
  "source": "WKND",  
  "attribute": "template",  
  "ruleType": "IGNORE",  
  "values": ["error-page", "redirect"]  
}
```

## Update Rule

```
PUT /api/v2/connector/indexing-rule/{id}
```

**Request body:** Same as create.

## Delete Rule

```
DELETE /api/v2/connector/indexing-rule/{id}
```

**Response:** `true` if deleted, `false` if not found.

## Get Empty Structure

```
GET /api/v2/connector/indexing-rule/structure
```

Returns an empty model for reference.

---

## AEM Plugin API

Base path: `/api/v2/aem`

These endpoints are available when the AEM plugin JAR is loaded via `-Dloader.path`.

### Health Check

```
GET /api/v2/aem/status
```

Returns `{"status": "ok"}`.

### Index AEM Paths

```
POST /api/v2/aem/index/{source}
```

Triggers indexing for specific AEM content paths. Includes built-in deduplication — repeated requests for the same paths within 30 seconds are ignored.

PARAMETER	LOCATION	DESCRIPTION
<code>source</code>	Path	AEM source name (e.g., <code>WKND</code> )
<i>(body)</i>	JSON	<code>DumAemPathList</code> object

**Request body:**

```
{
  "paths": ["/content/wknd/us/en/about"],
  "event": "INDEXING",
  "recursive": false,
  "attribute": "ID"
}
```

FIELD	TYPE	DEFAULT	DESCRIPTION
<code>paths</code>	string[]	<i>(required)</i>	AEM content paths
<code>event</code>	string	<code>INDEXING</code>	<code>INDEXING</code> , <code>DEINDEXING</code> , <code>PUBLISHING</code> , or <code>UNPUBLISHING</code>
<code>recursive</code>	boolean	<code>false</code>	Traverse child nodes
<code>attribute</code>	string	<code>ID</code>	<code>ID</code> (path-based) or <code>URL</code> (URL-based)

Response: `{"status": "sent"}`

## AEM Source Management

CRUD operations for AEM source configurations.

### List All Sources

```
GET /api/v2/aem/source
```

### Get Source by ID

```
GET /api/v2/aem/source/{id}
```

## Create Source

```
POST /api/v2/aem/source
```

**Request body:** Full `DumAemSource` JSON (see [Extending the AEM Connector – Configuration JSON](#)).

## Update Source

```
PUT /api/v2/aem/source/{id}
```

Returns `400` if ID mismatch, `404` if not found.

## Delete Source

```
DELETE /api/v2/aem/source/{id}
```

## Index All Content for Source

```
GET /api/v2/aem/source/{id}/indexAll
```

Triggers an asynchronous full indexing operation for the specified AEM source.

## Get Empty Structure

```
GET /api/v2/aem/source/structure
```

Returns an empty AEM source model for reference.

## Related Pages

---

PAGE	DESCRIPTION
<a href="#">AEM Connector</a>	AEM indexing flow, event listeners, and manual triggering
<a href="#">Connectors Overview</a>	All connectors and how they're managed
<a href="#">Turing ES — Integration</a>	Turing ES admin console for managing connectors

# Extending the AEM Connector

---

The AEM connector provides a powerful extension system that lets you customize how content is extracted, transformed, and indexed from Adobe Experience Manager. Extensions are Java classes that implement interfaces from the `aem-commons` library — published on Maven Central.

---

## Maven Dependency

---

To create custom AEM extensions, add `aem-commons` to your project:

```
<!-- Source: https://mvnrepository.com/artifact/com.viglet.dumont/aem-commons -->
-->
<dependency>
  <groupId>com.viglet.dumont</groupId>
  <artifactId>aem-commons</artifactId>
  <version>2026.1.19</version>
  <scope>compile</scope>
</dependency>
```

---

## Extension Interfaces

---

The AEM connector provides four extension interfaces:

INTERFACE	PURPOSE	CONFIG FIELD
<code>DumAemExtAttributeInterface</code>	Custom logic for extracting or transforming individual field values	<code>attributes[].className</code> or <code>sourceAttrs[].className</code>
<code>DumAemExtContentInterface</code>	Extract additional content from AEM pages (e.g., <code>.model.json</code> )	<code>models[].className</code>
<code>DumAemExtDeltaDateInterface</code>	Custom delta date resolution for incremental indexing	<code>sources[].deltaClass</code>
<code>DumAemExtUrlAttributeInterface</code>	Specialized URL handling with ID extraction (extends <code>DumAemExtAttributeInterface</code> )	<code>attributes[].className</code>

## DumAemExtAttributeInterface

The most commonly used extension. Implement this to transform or extract individual attribute values with custom logic.

```

import
com.viglet.dumont.connector.aem.commons.ext.DumAemExtAttributeInterface;
import com.viglet.dumont.connector.aem.commons.DumAemObject;
import com.viglet.dumont.connector.aem.commons.context.DumAemConfiguration;
import com.viglet.dumont.connector.aem.commons.mappers.*;
import com.viglet.turing.client.sn.TurMultiValue;

public class MyCustomAttribute implements DumAemExtAttributeInterface {

    @Override
    public TurMultiValue consume(
        DumAemTargetAttr dumAemTargetAttr,
        DumAemSourceAttr dumAemSourceAttr,
        DumAemObject aemObject,
        DumAemConfiguration dumAemConfiguration
    ) {
        // Example: extract a custom property and transform it
        String rawValue = aemObject.getAttributes()
            .get("myProperty").toString();
        return TurMultiValue.singleItem(rawValue.toUpperCase());
    }
}

```

#### Parameters:

PARAMETER	DESCRIPTION
<code>DumAemTargetAttr</code>	Target search field being populated (name, type)
<code>DumAemSourceAttr</code>	Source AEM property definition (name, className)
<code>DumAemObject</code>	AEM content node: <code>path</code> , <code>title</code> , <code>template</code> , <code>jcrNode</code> , <code>jcrContentNode</code> , <code>lastModified</code> , <code>attributes</code>
<code>DumAemConfiguration</code>	Connection config: <code>url</code> , <code>username</code> , <code>rootPath</code> , <code>authorSNSite</code> , <code>publishSNSite</code> , <code>providerName</code>

#### Built-in implementations:

CLASS	WHAT IT DOES
<code>DumAemExtContentId</code>	Returns the AEM page path as the document ID
<code>DumAemExtContentUrl</code>	Builds the full URL from the page path and URL prefix config
<code>DumAemExtContentTags</code>	Fetches tags from the <code>/jcr:content.tags.json</code> endpoint
<code>DumAemExtCreationDate</code>	Returns the <code>jcr:created</code> date
<code>DumAemExtModificationDate</code>	Returns the <code>cq:lastModified</code> or <code>jcr:lastModified</code> date
<code>DumAemExtPublicationDate</code>	Returns the last replication date
<code>DumAemExtHtml2Text</code>	Converts HTML content to plain text
<code>DumAemExtPageComponents</code>	Extracts text from responsive grid components
<code>DumAemExtTypeName</code>	Returns the content type name
<code>DumAemExtSourceApps</code>	Returns the provider name from configuration
<code>DumAemExtSiteName</code>	Returns the site name

## DumAemExtContentInterface

Implement this to fetch additional content from AEM — for example, calling the `.model.json` Sling Model exporter to get structured data.

```
import com.viglet.dumont.connector.aem.commons.ext.DumAemExtContentInterface;
import com.viglet.dumont.connector.aem.commons.bean.DumAemTargetAttrValueMap;

public class MyModelJsonExtractor implements DumAemExtContentInterface {

    @Override
    public DumAemTargetAttrValueMap consume(
        DumAemObject aemObject,
        DumAemConfiguration dumAemConfiguration
    ) {
        DumAemTargetAttrValueMap result = new DumAemTargetAttrValueMap();

        // Fetch the .model.json endpoint
        String url = dumAemConfiguration.getUrl()
            + aemObject.getPath() + ".model.json";
        // ... HTTP call to get JSON ...

        result.addWithSingleValue("fragmentPath",
            "/content/dam/fragment", false);
        return result;
    }
}
```

Referenced in the `models[].className` field of the configuration JSON.

---

## DumAemExtDeltaDateInterface

Customize how the connector determines the "last modified" date for incremental indexing.

```
import
com.viglet.dumont.connector.aem.commons.ext.DumAemExtDeltaDateInterface;

public class MyDeltaDate implements DumAemExtDeltaDateInterface {

    @Override
    public Date consume(
        DumAemObject aemObject,
        DumAemConfiguration dumAemConfiguration
    ) {
        return Optional.ofNullable(aemObject.getLastModified())
            .map(Calendar::getTime)
            .orElse(null);
    }
}
```

Referenced in the `sources[].deltaClass` field of the configuration JSON.

---

## AEM Configuration JSON

---

The AEM connector is configured via a JSON file that defines sources, attributes, locale mappings, and content models. This file is placed in an `export/` directory and imported at startup.

## Full Example (WKND Site)

```

{
  "sources": [
    {
      "name": "WKND",
      "defaultLocale": "en_US",
      "localeClass":
"com.viglet.dumont.connector.aem.commons.ext.DumAemExtLocale",
      "deltaClass": "com.example.MyDeltaDate",
      "endpoint": "http://localhost:4502",
      "username": "admin",
      "password": "admin",
      "oncePattern": "^/content/wknd/us/en/faqs",
      "rootPath": "/content/wknd",
      "contentType": "cq:Page",
      "author": true,
      "publish": true,
      "authorSNSite": "wknd-author",
      "publishSNSite": "wknd-publish",
      "authorURLPrefix": "http://localhost:4502",
      "publishURLPrefix": "https://wknd.site",
      "localePaths": [
        { "locale": "en_US", "path": "/content/wknd/us/en" },
        { "locale": "es", "path": "/content/wknd/es/es" }
      ],
      "attributes": [
        {
          "name": "id", "type": "STRING", "mandatory": true,
          "className":
"com.viglet.dumont.connector.aem.commons.ext.DumAemExtContentId"
        },
        {
          "name": "title", "type": "TEXT", "mandatory": true,
          "facetName": { "default": "Titles", "pt_BR": "Títulos" }
        },
        {
          "name": "tags", "type": "STRING", "multiValued": true,
          "facet": true, "facetName": { "default": "Tags" },
          "className":
"com.viglet.dumont.connector.aem.commons.ext.DumAemExtContentTags"
        },
        {
          "name": "url", "type": "STRING", "mandatory": true,
          "className":
"com.viglet.dumont.connector.aem.commons.ext.DumAemExtContentUrl"
        }
      ],
    }
  ],
}

```

```
"models": [
  {
    "type": "cq:Page",
    "className": "com.example.MyModelJsonExtractor",
    "targetAttrs": [
      { "name": "title", "sourceAttrs": [{ "name": "jcr:title" }] },
      { "name": "tags", "sourceAttrs": [{ "name": "cq:tags" }] },
      {
        "name": "text",
        "sourceAttrs": [{
          "className":
"com.viglet.dumont.connector.aem.common.ext.DumAemExtPageComponents"
        }]
      }
    ]
  }
]
```

## Source Fields

FIELD	TYPE	DESCRIPTION
<code>name</code>	string	Source identifier (displayed in the admin console)
<code>endpoint</code>	string	AEM instance URL (e.g., <code>http://localhost:4502</code> )
<code>username</code> / <code>password</code>	string	AEM authentication credentials
<code>rootPath</code>	string	Content tree root to crawl (e.g., <code>/content/wknd</code> )
<code>contentType</code>	string	JCR node type to index (e.g., <code>cq:Page</code> )
<code>defaultLocale</code>	string	Fallback locale code (e.g., <code>en_US</code> )
<code>localeClass</code>	string	Class for locale resolution
<code>deltaClass</code>	string	Class implementing <code>DumAemExtDeltaDateInterface</code>
<code>oncePattern</code>	string	Regex — matching paths are indexed only once (never re-indexed)
<code>author</code> / <code>publish</code>	boolean	Enable indexing from author/publish environments
<code>authorSNSite</code> / <code>publishSNSite</code>	string	Turing ES SN Site names for each environment
<code>authorURLPrefix</code> / <code>publishURLPrefix</code>	string	Public URL prefixes for documents

## Locale Paths

```
"localePaths": [
  { "locale": "en_US", "path": "/content/wknd/us/en" },
  { "locale": "es", "path": "/content/wknd/es/es" }
]
```

Content found under each path is tagged with the corresponding locale.

## Attribute Fields

FIELD	TYPE	DESCRIPTION
<code>name</code>	string	Field name in the search index
<code>type</code>	string	<code>STRING</code> , <code>TEXT</code> , or <code>DATE</code>
<code>mandatory</code>	boolean	Whether this field is required
<code>multiValued</code>	boolean	Whether this field holds multiple values
<code>description</code>	string	Human-readable description
<code>facet</code>	boolean	Expose as a facet filter in search results
<code>facetName</code>	object	Localized labels: <code>{ "default": "Tags", "pt_BR": "Etiquetas" }</code>
<code>className</code>	string	Class implementing <code>DumAemExtAttributeInterface</code> – extracts the value instead of reading from JCR

## Model Fields

FIELD	TYPE	DESCRIPTION
<code>type</code>	string	JCR node type this model applies to
<code>className</code>	string	Class implementing <code>DumAemExtContentInterface</code>
<code>targetAttrs[].name</code>	string	Target field (must match <code>attributes[]</code> )
<code>targetAttrs[].sourceAttrs[] .name</code>	string	JCR property to read (e.g., <code>jcr:title</code> )
<code>targetAttrs[].sourceAttrs[] .className</code>	string	Class implementing <code>DumAemExtAttributeInterface</code> for custom extraction

## Creating a Custom AEM Extension

### Step 1 – Create a Maven project

```
<project>
  <groupId>com.example</groupId>
  <artifactId>my-aem-extensions</artifactId>
  <version>1.0.0</version>

  <dependencies>
    <dependency>
      <groupId>com.viglet.dumont</groupId>
      <artifactId>aem-commons</artifactId>
      <version>2026.1.19</version>
    </dependency>
  </dependencies>
</project>
```

### Step 2 – Implement your extension

```
package com.example.ext;

import
com.viglet.dumont.connector.aem.common.ext.DumAemExtAttributeInterface;
import com.viglet.turing.client.sn.TurMultiValue;

public class MyBreadcrumb implements DumAemExtAttributeInterface {
  @Override
  public TurMultiValue consume(DumAemTargetAttr target,
    DumAemSourceAttr source, DumAemObject aemObject,
    DumAemConfiguration config) {
    String path = aemObject.getPath()
      .replace(config.getRootPath(), "");
    return TurMultiValue.singleItem(
      String.join(" > ", path.split("/")));
  }
}
```

## Step 3 – Build and deploy

```
mvn clean package
cp target/my-aem-extensions-1.0.0.jar /apl/viglet/dumont/aem/libs/
```

The `libs/` directory must contain both `aem-plugin.jar` and your extension JAR.

## Step 4 – Reference in the JSON

```
{
  "name": "breadcrumb",
  "type": "STRING",
  "className": "com.example.ext.MyBreadcrumb"
}
```

## How classes are loaded

Extension classes are loaded via `DumCustomClassCache` using `Class.forName()`. Requirements:

**Public no-argument constructor**

**On the classpath** (via `libs/` and `-Dloader.path`)

**Thread-safe** (one instance is shared across all calls)

## Related Pages

PAGE	DESCRIPTION
<a href="#">AEM Connector</a>	AEM connector features, configuration, and locale mapping
<a href="#">Installation Guide</a>	How to deploy plugins with <code>-Dloader.path</code>
<a href="#">Developer Guide</a>	Project structure, build, and contribution guide

# Extending the Database Connector

The Database Connector supports a custom extension point that lets you transform, enrich, or filter database rows before they are indexed. Extensions are Java classes that implement the

`DumDbExtCustomImpl` interface from the `db-commons` library.

## Maven Dependency

```
<!-- Source: https://mvnrepository.com/artifact/com.viglet.dumont/db-commons ->
<dependency>
  <groupId>com.viglet.dumont</groupId>
  <artifactId>db-commons</artifactId>
  <version>2026.1.19</version>
  <scope>compile</scope>
</dependency>
```

## The DumDbExtCustomImpl Interface

```
public interface DumDbExtCustomImpl {
    Map<String, Object> run(Connection connection, Map<String, Object>
attributes);
}
```

PARAMETER	DESCRIPTION
<code>Connection connection</code>	The active JDBC connection — you can execute additional queries to enrich the row
<code>Map&lt;String, Object&gt; attributes</code>	Mutable map of column values from the current row (column aliases become keys)

**Return:** The modified attributes map. The returned map is what gets indexed.

## How It Works

For each row in the SQL result set, the connector:

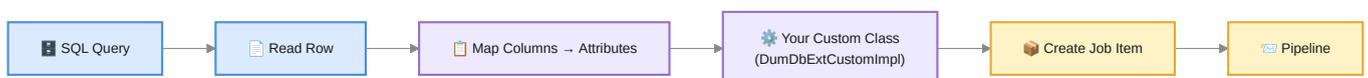
Reads all columns into an attributes map ( `Map<String, Object>` )

Applies value formatting (HTML tag removal, date formatting)

**Calls your custom class** via `run(connection, attributes)`

Processes multi-valued fields

Creates a Job Item and sends it to the pipeline



## Example: Prefixing Titles

A minimal example that prefixes every title with a label:

```
package com.example.ext;

import com.viglet.dumont.connector.db.ext.DumDbExtCustomImpl;
import java.sql.Connection;
import java.util.Map;

public class MyDbCustom implements DumDbExtCustomImpl {

    private static final String TITLE = "title";

    @Override
    public Map<String, Object> run(Connection connection,
                                   Map<String, Object> attributes) {
        if (attributes.containsKey(TITLE)) {
            attributes.replace(TITLE,
                              String.format("Product: %s", attributes.get(TITLE)));
        }
        return attributes;
    }
}
```

## Example: Enriching with a Second Query

Use the JDBC connection to run additional queries and add data:

```
public class EnrichWithCategory implements DumDbExtCustomImpl {

    @Override
    public Map<String, Object> run(Connection connection,
                                   Map<String, Object> attributes) {
        Object categoryId = attributes.get("category_id");
        if (categoryId != null) {
            try (var stmt = connection.prepareStatement(
                "SELECT name FROM categories WHERE id = ?")) {
                stmt.setObject(1, categoryId);
                try (var rs = stmt.executeQuery()) {
                    if (rs.next()) {
                        attributes.put("category", rs.getString("name"));
                    }
                }
            } catch (Exception e) {
                // log and continue
            }
        }
        return attributes;
    }
}
```

## Using the Extension

### Via CLI

Pass the fully-qualified class name with the `--class-name` parameter:

```
java -cp dumont-db-indexer.jar:my-db-extensions.jar \  
  com.viglet.dumont.connector.db.DumDbImportTool \  
  --class-name com.example.ext.MyDbCustom \  
  --server http://localhost:30130 \  
  --api-key <API_KEY> \  
  --driver org.mariadb.jdbc.Driver \  
  --connect "jdbc:mariadb://localhost:3306/shop" \  
  --query "SELECT id, name AS title, price FROM products" \  
  --site ProductCatalog \  
  --locale en_US
```

#### CLASSPATH

Your extension JAR must be on the classpath (`-cp`). Add it alongside the `dumont-db-indexer.jar` separated by `:` (Linux) or `;` (Windows).

# Creating a Custom DB Extension

## Step 1 – Create a Maven project

```
<project>
  <groupId>com.example</groupId>
  <artifactId>my-db-extensions</artifactId>
  <version>1.0.0</version>

  <dependencies>
    <dependency>
      <groupId>com.viglet.dumont</groupId>
      <artifactId>db-commons</artifactId>
      <version>2026.1.19</version>
    </dependency>
  </dependencies>
</project>
```

## Step 2 – Implement the interface

```
package com.example.ext;

import com.viglet.dumont.connector.db.ext.DumDbExtCustomImpl;
import java.sql.Connection;
import java.util.Map;

public class MyDbCustom implements DumDbExtCustomImpl {
    @Override
    public Map<String, Object> run(Connection connection,
                                   Map<String, Object> attributes) {
        // Transform, enrich, or filter attributes
        return attributes;
    }
}
```

## Step 3 – Build

```
mvn clean package
```

## Step 4 — Run with the extension

```
java -cp dumont-db-indexer.jar:target/my-db-extensions-1.0.0.jar \  
  com.viglet.dumont.connector.db.DumDbImportTool \  
  --class-name com.example.ext.MyDbCustom \  
  --server http://localhost:30130 \  
  --api-key <API_KEY> \  
  ...
```

## Requirements

**Public no-argument constructor**

**On the classpath** (via `-cp`)

**Thread-safe** — one instance processes all rows sequentially, but should avoid mutable shared state

## Related Pages

PAGE	DESCRIPTION
<a href="#">Database Connector</a>	CLI parameters, SQL examples, supported databases
<a href="#">Extending the AEM Connector</a>	AEM extension interfaces and configuration JSON
<a href="#">Developer Guide</a>	Project structure, build, and contribution guide