# AIE(Upload and Image Extraction) Architecture

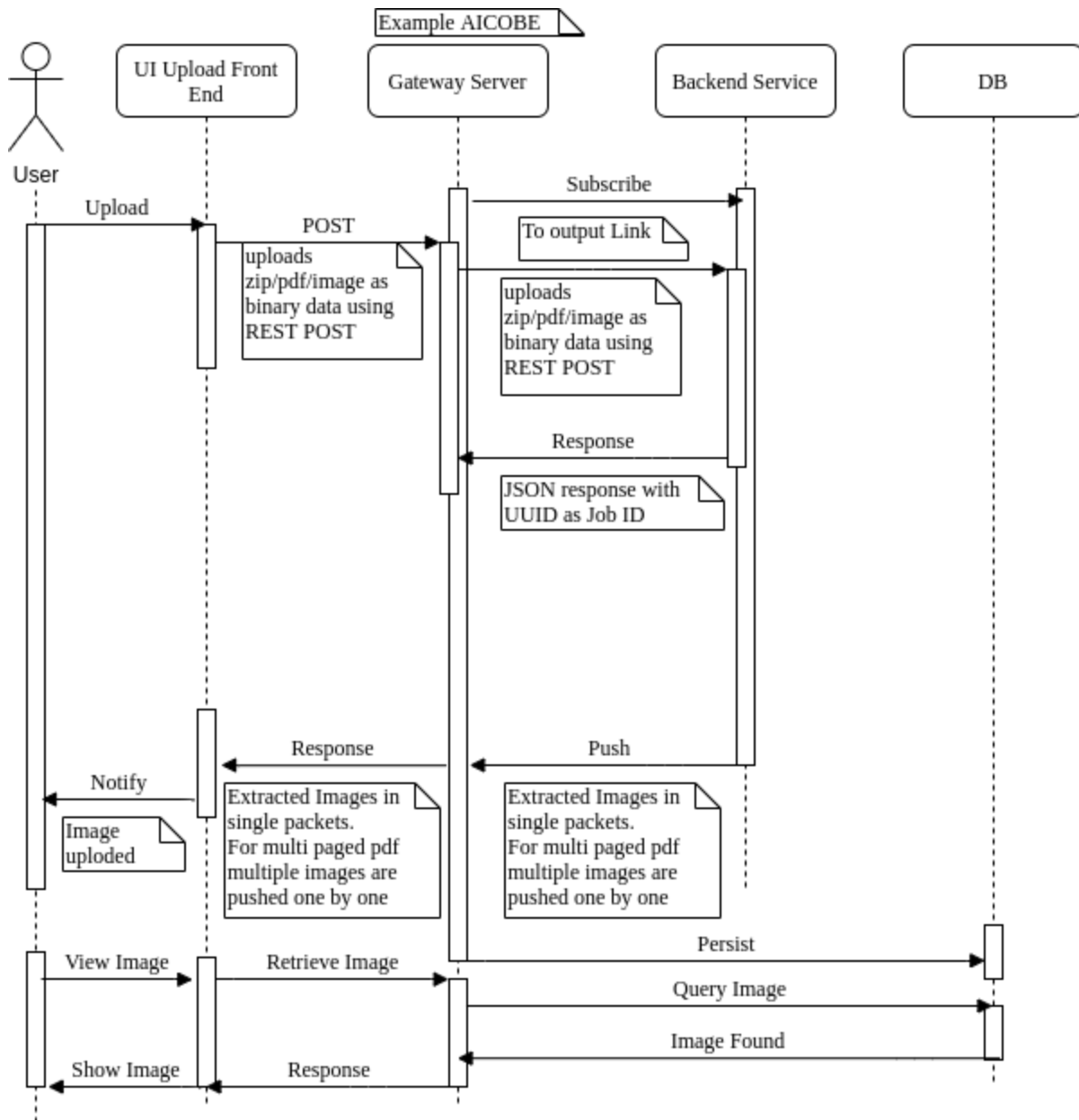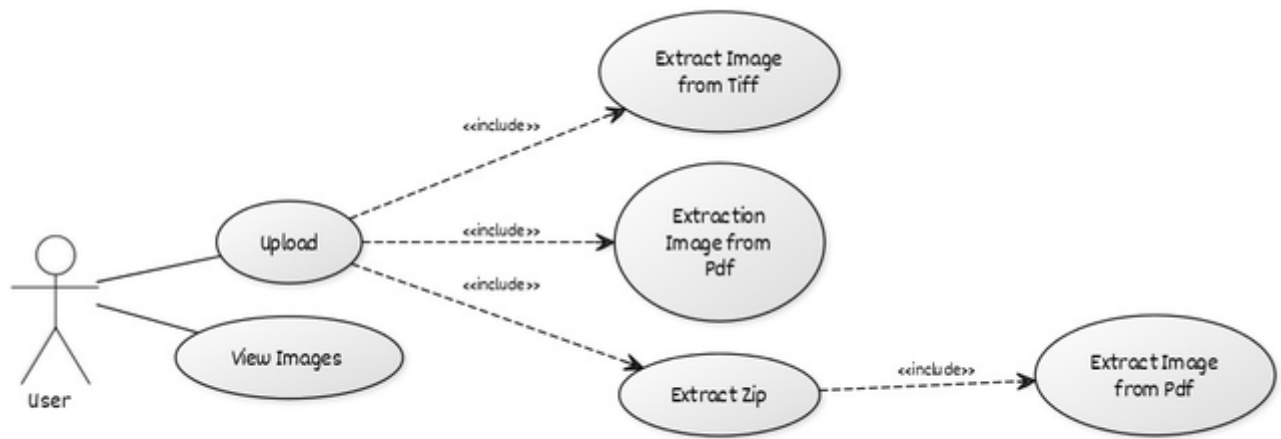| Architecture review date | 26 May 2020 |
|---|---|
| Refactor review date | 05 Aug 2020 |
| Project lead | @ Gregor Blichmann |
| On this page | <ul><li>Overview</li><li>Workflow Sequence</li><li>Use Case</li><li>Architecture Design Options<ul><li>Option 1</li><li>Option 2</li><li>Option 3</li></ul></li><li>Architecture issues</li><li>Hybrid Architecture: (Option 1 Option 2)</li><li>Data Flow</li><li>Data Model of Extraction Output<ul><li>Success Message:</li><li>Failure Message:</li></ul></li><li>Comparison : Hazelcast vs Redis</li><li>Stakeholders</li><li>Refactoring Points</li><li>Strategy Interaction Diagram</li></ul> |

📋 Overview

This is our discussion on 25 May 2020 on white board where we put all our different design approaches for the end to end workflow. Below every designs are drawn in detail with explanations.

## 📈 Workflow Sequence

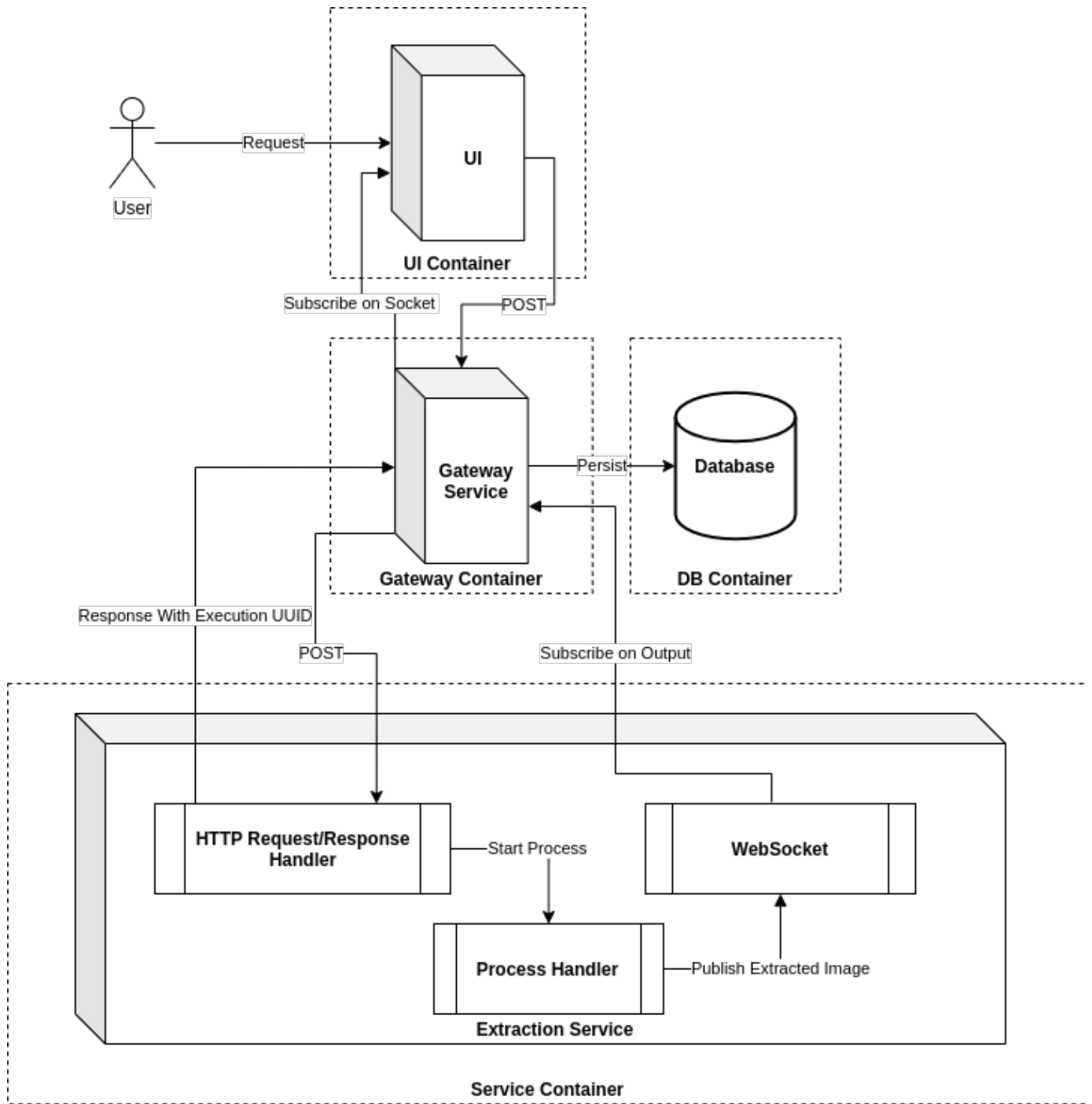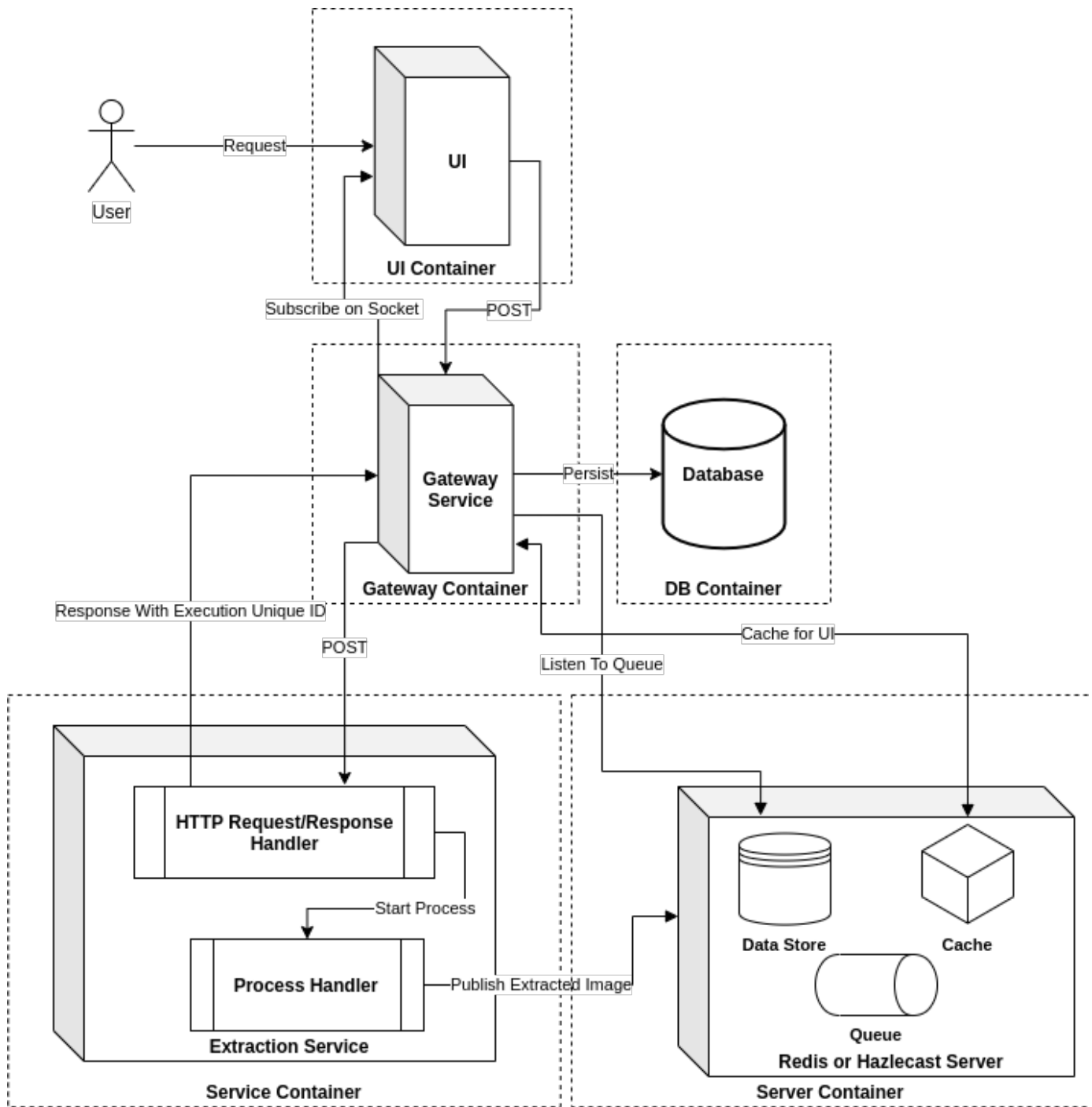Overall Idea on the Step by Step Sequence of Image Extraction Workflow

**Example AICOBE**

Sequence diagram:

- **User** (actor)
- **UI Upload Front End**
- **Gateway Server**
- **Backend Service**
- **DB**

Messages and notes:

- User → UI Upload Front End: **Upload**
- UI Upload Front End → Gateway Server: **POST**
  - Note: *uploads zip/pdf/image as binary data using REST POST*
- Gateway Server → Backend Service: **Subscribe**
  - Note: *To output Link*
- Gateway Server → Backend Service: (upload)
  - Note: *uploads zip/pdf/image as binary data using REST POST*
- Backend Service → Gateway Server: **Response**
  - Note: *JSON response with UUID as Job ID*
- Backend Service → Gateway Server: **Push**
  - Note: *Extracted Images in single packets. For multi paged pdf multiple images are pushed one by one*
- Gateway Server → UI Upload Front End: **Response**
  - Note: *Extracted Images in single packets. For multi paged pdf multiple images are pushed one by one*
- UI Upload Front End → User: **Notify**
  - Note: *Image uploded*
- Gateway Server → DB: **Persist**
- User → UI Upload Front End: **View Image**
- UI Upload Front End → Gateway Server: **Retrieve Image**
- Gateway Server → DB: **Query Image**
- DB → Gateway Server: **Image Found**
- Gateway Server → UI Upload Front End: **Response**
- UI Upload Front End → User: **Show Image**

**📊 Use Case**

CREATED WITH YUML

🗺 Architecture Design Options

**Option 1**

User

Request

**UI**

UI Container

Subscribe on Socket

POST

**Gateway Service**

Persist

**Database**

Gateway Container

DB Container

Response With Execution UUID

POST

Subscribe on Output

**HTTP Request/Response Handler**

Start Process

**WebSocket**

**Process Handler**

Publish Extracted Image

**Extraction Service**

**Service Container**

**Option 2**

**Option 3**

## 🚩 Architecture issues

| Architecture Options | Implementation Complexity | Maintenance | External Dependency | Notes |
|---|---|---|---|---|
| Option 1 | HIGH | MEDIUM | LOW | |
| Option 2 | MEDIUM | HIGH | HIGH | |
| Option 3 | MEDIUM | HIGH | HIGH | |

## 🧬 Hybrid Architecture: (Option 1 ➕ Option 2)

User

Request

UI

UI Container

Subscribe on Socket

POST

Gateway Service

Gateway Container

Persist

Database

DB Container

Response With Execution UUID

POST

Subscribe on Output

POST

Socket Usage or HazelCast Usage

Listens to Output Queue

HTTP Request/Response Handler

Start Process

WebSocket

Process Handler

Publish Extracted Image

Extraction Service

HTTP Request/Response Handler

Start Process

Process Handler

Publish Extracted Image

Extraction Service

Data Store

Cache

Queue

Hazlecast Server

Service Container

Service Container

Server Container

## 💾 Data Flow

ML guy

0. Select file to upload

AICODA

File upload form

REST API

Collections (jobs) view

WS-Client

1. Send request with ZIP file as body

6. Forward extracted files one by one to the frontend

7. Display upload progress in the UI

AICOBE

REST API as proxy

**Data Handler**
1. Creates collections
2. Creates data objects
3. Creates GridFs objects
4. Modify WS message for the frontend representation

WS-Server

2. Forward the request to extraction service

Image Extraction Service

Extract from zip, convert to png etc.

3. Push extracted files one by one

WS-Client

WS-Server

5. Save created objects into DB

MongoDb

## 🟠 Data Model of Extraction Output

## Success Message:

```
{
  "jobId": "0ef90871-c72b-411e-9e71-d8fc22b967ce",
  "documentName": "temp_0ef90871-c72b-411e-9e71-d8fc22b967ce.pdf",
  "pageFileName": "temp_0ef90871-c72b-411e-9e71-d8fc22b967ce_0.png",
  "pageCount": 32,
  "pageIndex": 0,
  "pageWidth": 800,
  "pageHeight: 200,
  "pageContent": "iVBORw0KGgoAAAANSUhEUgAABnUAAAkiCAIAAACTwI40...",
//Base64 png image string
  "documentCount":1,
  "documentIndex":0,
  "documentPath":"folderOne/nested",
  "meta": {}
}
```

**Failure Message:**

```
{
  "jobId": "0ef90871-c72b-411e-9e71-d8fc22b967ce",
  "documentName": "temp_0ef90871-c72b-411e-9e71-d8fc22b967ce.pdf",
  "pageFileName": "",
  "pageCount": 32,
  "pageIndex": 2,
  "pageContent": "", //Base64 png image string - no content since failed
  "documentCount":1,
  "documentIndex":0,
  "documentPath":"folderOne/nested",
  "failureReason": "NullPointerException"
}
```

## ⚒ Comparison : Hazelcast vs Redis

| Name | Hazelcast | Redis |
| --- | --- | --- |
| Description | A widely adopted in-memory data grid | In-memory data structure store, used as database, cache and message broker |
| Primary database model | Key-value store | Key-value store |
| Secondary database models | Document store | Document store<br><br>Graph DBMS<br><br>Search engine<br><br>Time Series DBMS |

| | | |
|---|---|---|
| Website | hazelcast.com | redis.io |
| Technical documentation | hazelcast.org/-imdg/-docs | redis.io/-documentation |
| Developer | Hazelcast | Salvatore Sanfilippo |
| Initial release | 2008 | 2009 |
| Current release | 4.0, February 2020 | 6.0.1, May 2020 |
| License | Open Source | Open Source |
| Cloud-based only | no | no |
| DBaaS offerings (sponsored links) | | |
| Implementation language | Java | C |
| Server operating systems | All OS with a Java VM | BSD<br>Linux<br>OS X<br>Windows |
| Data scheme | schema-free | schema-free |
| Typing | yes | partial |
| XML support | yes | no |
| Secondary indexes | yes | yes |
| SQL | SQL-like query language | no |
| APIs and other access methods | JCache<br>JPA<br>Memcached protocol<br>RESTful HTTP API | proprietary protocol |

| | | |
|---|---|---|
| Supported programming languages | .Net<br>C#<br>C++<br>Clojure<br>Go<br>Java<br>JavaScript (Node.js)<br>Python<br>Scala | C<br>C#<br>C++<br>Clojure<br>Crystal<br>D<br>Dart<br>Elixir<br>Erlang<br>Fancy<br>Go<br>Haskell<br>Haxe<br>Java<br>JavaScript (Node.js)<br>Lisp<br>Lua<br>MatLab<br>Objective-C<br>OCaml<br>Pascal<br>Perl<br>PHP<br>Prolog<br>Pure Data<br>Python<br>R<br>Rebol<br>Ruby<br>Rust<br>Scala<br>Scheme<br>Smalltalk<br>Swift<br>Tcl<br>Visual Basic |
| Server-side scripts | yes | Lua |
| Triggers | yes | no |
| Partitioning methods | Sharding | Sharding |
| Replication methods | yes | Master-slave replication<br><br>Multi-master replication |
| MapReduce | yes | no |
| Consistency concepts | Immediate Consistency or Eventual Consistency selectable by user | Eventual Consistency<br>Strong eventual consistency with CRDTs |
| Foreign keys | no | no |
| Transaction concepts | one or two-phase-commit; repeatable reads; read committed | Optimistic locking, atomic execution of commands blocks and scripts |
| Concurrency | yes | yes |
| Durability | yes | yes |
| In-memory capabilities | yes | yes |
| User concepts | Role-based access control | Simple password-based access control |

## 👥 Stakeholders

| Name | Role |
|---|---|

| | | |
|---|---|---|
| @ Sankha Sil  @ Alexander Naumenko | | Software Engineer |
| @ Sergej Tschigraj | | UI Engineer |
| @ Gregor Blichmann | | Product Owner |

## 🧰 Refactoring Points

Currently Ixtraction service is running in https://ixtract.ai4bd.org to check one can use https://ixtract.ai4bd.org/image/extraction/ping link.

There lots of areas of improvement regarding code quality, code design and memory leaks. Some are listed below.

1. ~~Controller has overdesign of filetype check 💡 Remove~~
   a. ~~fileType request header. Check content type.~~
2. ~~Controller endpoints are not meaningful. 💡 Remove~~
   a. ~~Get user endpoint~~
   b. ~~detailed request header.~~
   c. ~~as soon as~~ `/info` ~~called remove the entry of hashmap.~~
   d. ~~Change the endpoint urls. Talk to WAPP team and update openAPI.~~
3. ~~Controller has service logic in endpoint of process. 💡 put logic to service~~
4. ~~Service File is huge with all the logic of extraction of image. 💡 Design pattern Strategy pattern,Factory pattern and Helper pattern.~~
   a. ~~First design class diagram to proceed.~~
5. ~~Switch case are used for different type of source type to determine the image extraction logic. 💡 Strategy pattern.~~
   a. ~~Recursive process for archive formats. (upto 3 recursion)~~
6. ~~Too many small methods used. 💡 Put methods into Util class .~~
7. ~~Methods having more than 7 parameters 💡 remove parameter array and put to class.~~
8. ~~Unnecessary Constants used. 💡 put to a common constant or Enum.~~
9. ~~Unit Test coverage is poor. 💡 TDD refactoring , e.g Util test before Util class.~~

1. ~~Methods for Socket/Hazelcast send/push 💡 Observer pattern with using strategy~~
2. Meaningful Logging required.

## 📕 Strategy Interaction Diagram

## ImageExtractionService

+ extractionStrategy : ExtractionStrategyService

---

## ExtractionStrategyService

+ fileType: String
+ archiveStrategy : ArchiveStrategyService
+ docRepository : DocumentRepository

---

«interface»
## Archive Strategy Service
flatten(byte[]) : List<File>

---

«interface»
## Document Strategy Repository
extract(byte[]) : List<File>

---

## ZipArchive Strategy
@Bean("zip")

docRepository: DocumentRepository

+ flatten(byte[]): List<File>

---

## TarArchive Strategy
@Bean("tar")

+ docRepository: DocumentRepository

+ flatten(byte[]): List<File>

---

«interface»
## Notification Strategy Repository
notify(JsonObject) : void

---

## Socket Strategy
@Bean("socket")

+ notify(JsonObject): void

---

## Hazelcast Strategy
@Bean("hazelcast")

+ notify(JsonObject): void

---

## Tiff Strategy
@Bean("tiff")

notificationRepo: NotificationRepository

+ extract(byte[]): List<File>

---

## PDF Strategy
@Bean("pdf")

+ notificationRepo: NotificationRepository

+ extract(byte[]): List<File>

---

## Jpeg Strategy
@Bean("jpeg")

notificationRepo: NotificationRepository

+ extract(byte[]): List<File>

---

## PNG Strategy
@Bean("png")

notificationRepo: NotificationRepository

+ extract(byte[]): List<File>