



MASTER SCIENCE ET TECHNOLOGIE DU LOGICIEL

---

Rapport Projet DAAR

# Indexation de CVs dans Elasticsearch

---

*Réalisé par:*

Azzoug Koceila

Sofiane Braneci

2021/2022

# Table des matières

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Plan de développement</b>	<b>1</b>
<b>3</b>	<b>Indexation des documents</b>	<b>1</b>
3.1	Parsing . . . . .	1
3.2	Stockage et indexation . . . . .	2
<b>4</b>	<b>API</b>	<b>3</b>
<b>5</b>	<b>Quelque tests d'ajout et recherche</b>	<b>3</b>
<b>6</b>	<b>Conclusion</b>	<b>5</b>

# 1 Introduction

L'informatique moderne et l'internet ont rendu accessible une grande quantité d'informations. La capacité à rechercher efficacement ces informations est fondamentale pour les traiter. À l'époque où nous vivons, la quantité de données générées chaque jour est exponentielle, et le temps nécessaire pour retrouver l'information pertinente est proportionnel à la quantité des données dans lesquelles nous effectuons la recherche. Toutes les applications modernes doivent être conçues dans un souci de performance afin de répondre aux besoins des usagers, et la recherche est un processus clé pour répondre à chaque requête. Presque toutes ces applications utilisent une sorte de solution de gestion et de stockage des données, et la plupart d'entre elles s'appuient sur des systèmes de bases de données comme MySQL, MongoDB. En revanche, ces dernières années, Elasticsearch a trouvé sa place dans les plus grandes entreprises mondiales comme Facebook, Netflix, Deezer, Microsoft ...etc. Cela est dû à ses performances en recherche et analyse de données car il opère en quasi temps réel et sur de grands volumes de données, elle est devenue la solution de recherche la plus populaire. Les indexs sont utilisés pour accélérer la localisation des données. Dans ce projet nous utiliserons Elasticsearch comme solution pour indexer des CV, Elasticsearch est un moteur de recherche open source hautement évolutif. Bien qu'il ait commencé comme un moteur de recherche de texte, il est en train d'évoluer comme un moteur analytique, qui peut supporter non seulement la recherche mais aussi des agrégations complexes.

## 2 Plan de développement

Le but du projet étant d'indexer des CVs sur Elastic search dans le but de faire des recherches efficaces et optimisées. La démarche entreprise consiste à :

— Indexation de documents

1. Parsing des documents en string et base64
2. Stockage des documents sous forme d'un JSON contenant

```
1 {  
2   "filename": "filename-UUID.pdf",  
3   "textContent": "file as String",  
4   "content": "file as Base64 string"  
5 }
```

— Recherche : recherche de tous les CVs contenant une compétence

## 3 Indexation des documents

### 3.1 Parsing

Le besoin de cette partie était d'indexer un CV d'une manière à faire une recherche efficace et de garder le document de manière à pouvoir l'instaurer au moment

voulu ( modèle et contenu).

La classe *PDFParser* s'occupe de tous le parsing.

- `parse(byte [] bytes)` : ici, on utilise la bibliothèque **PDFBox** qui charge le fichier initialement sous forme d'un tableau de Bytes puis invoque **PDFTextStripper** qui extrait le contenu en string.
- `encodePdfBase64(InputStream fileInputStream)` : La requete d'ajout d'un document passe le CV sous forme de InputStream. Cette méthode transforme les bytes de ce stream en une chaine de caractères encodé en Base64.
- `decode(String strEncoded,String saveRep)` : cette méthode à pour but d'instaurer le contenu original du document passée ici en paramètre `strEncoded` (encodé en base64 )et le sauvegarde dans le répertoire de destination `saveRep`.

## 3.2 Stockage et indexation

Pour Pouvoir interagir avec Elasticsearch, on a besoin d'un client `RestHighLevelClient` a travers le quel on peut emettre des requetes et recevoir des réponses.

La classe `Config` annoté par **@Configuration** contenant un **Bean** se charge de créer ce client et le rajoute au contexte d'exécution au lancement de l'application. Les étapes ci-dessous résume les fonctionnalités principales implémentées dans `IndexHandler` ;

### 1. Création d'index :

La méthode **createIndex** qui prend en paramètre le nom souhaité pour l'index s'occupe de cette fonctionnalité, On utilise le client crée précédemment pour vérifier dans un premier temps l'existence de l'index. Si ce dernier n'existe pas, on crée une requete de type **CreateIndexRequest** ayant comme paramètres le nombre de shards ainsi que le nombre de replicas.

### 2. Ajout de document : Cette fonctionnalité est assuré par la méthode `handleUpload` prenant en paramètres un objet de type `MultiPartFile` qui permet de récupérer le fichier en Bytes ou `InputStream`.

- (a) Parser le fichier en string en utilisant `PDFParser`
- (b) Encoder le fichier en Base64 en utilisant `PDFParser`
- (c) On construit le document a insérer contenant les trois attributs `filename`, `textContent`, `content`.
- (d) On crée la requete `IndexRequest` avec notre index et le document crée précédemment
- (e) le client se charge d'envoyer la réponse et de collecter la réponse **Index-Response**

### 3. Recherche : La recherche s'effectue au niveau de la méthode `search` prenant en paramètres un mot clé 'skill'.

- (a) Construire une requête de type `MatchQuery` afin de rechercher tous les documents qui match `skill` dans leur attribut 'textContent'
- (b) On envoie la requete a travers notre client, on récupère la réponse sous forme de `SearchResponse`

- (c) A partir de la réponse, on extrait les SearchHit( contenant les attributs dont on a besoin ).
- (d) On parcourt la réponse( qui correspond a un tableau) pour construire une liste de CV( CV : structure contenant le nom du fichier, l'encodage en base64).
- (e) On décode chaque CV en enregistrant le résultat( fichier PDF) dans le répertoire Docs et on renvoie au client les noms des CVs résultants.

## 4 API

Notre API est implémenté par CVController, elle contient deux **Endpoints**, upload et seach/skill qui servent respectivement a ajouter un CV et chercher un mot clé parmi tous les documents indexés.

- "/upload" : prend en charge les requetes HTTP POST contenant le fichier à indexer, la méthode correspondante a ce mapping invoque la service handleUpload de indexHandler.
- "/search/skill : prend en charge les requetes HTTP GET invoquant le service search en lui passant 'skill' en paramètre, la réponse étant une liste de String correspondant au noms des fichiers résultants.

## 5 Quelques tests d'ajout et recherche

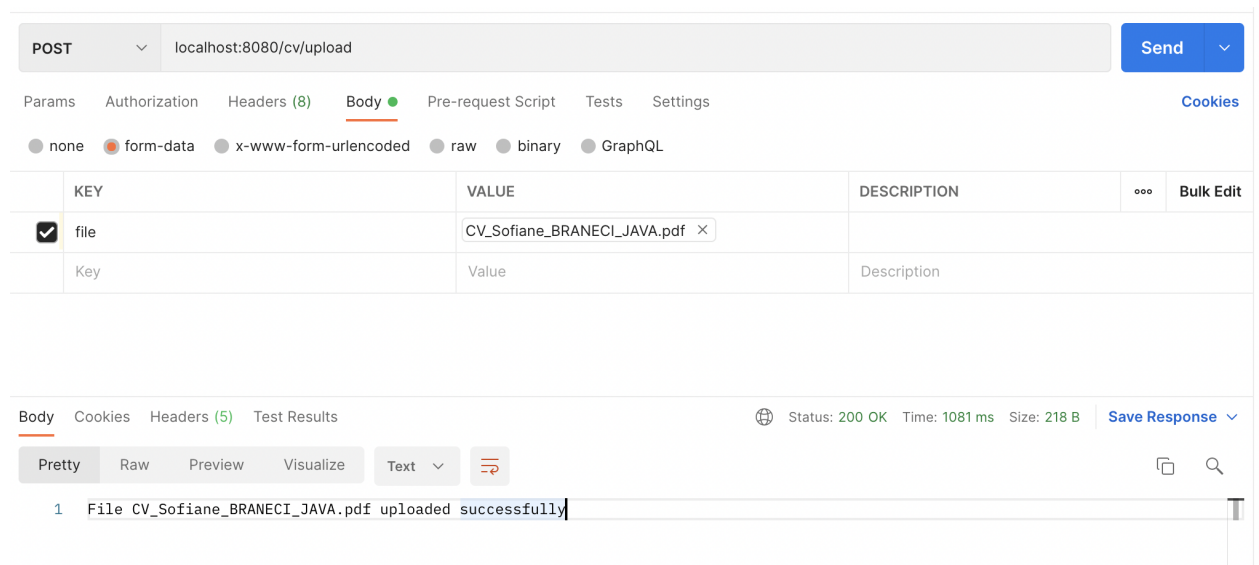


FIGURE 1 – Ajout d'un cv et réponse du serveur.

```

L87 | | | | "_source" : {
L88 | | | | | "filename" : "CV_Sofiane_BRANECI_JAVAae4be832-9501-4a3d-a35d
      | | | | | -a851ca74f1b4.pdf",
L89 | | | | | "textContent" : """/
L90 | sofiane|
L91 | braneci
L92 | développeur java junior à la recherche d'un
L93 | stage de 6 mois en développement logiciel
L94 | projet s r é a l i s é s
L95 | model checking: cas global properties
L96 | sorbonne université paris
L97 | mise en oeuvre d'un nouvel outil pour la vérification
L98 | de programmes
L99 | entièrement implémenté en java
200 | optimisation du solver existant en rajoutant des
201 | stratégies

```

FIGURE 2 – Preuve d'ajout sur elasticsearch.

GET localhost:8080/cv/search/java Send

Params Authorization Headers (8) Body Pre-request Script Tests Settings Cookies

none form-data x-www-form-urlencoded raw binary GraphQL

	KEY	VALUE	DESCRIPTION	...	Bulk Edit
<input checked="" type="checkbox"/>	file	bert-de-vriess-cv.pdf ×			
	Key	Value	Description		

Body Cookies Headers (5) Test Results Status: 200 OK Time: 172 ms Size: 360 B Save Response

Pretty Raw Preview Visualize JSON

```

1 | [
2 |   "CV_Sofiane_BRANECI_JAVAae4be832-9501-4a3d-a35d-a851ca74f1b4.pdf",
3 |   "navneet-pandeys-resumee43bcbf7-6cb7-4f81-a3e5-4596e763e087.pdf",
4 |   "a-customised-curve-cv449f7164-0406-4f16-8494-ed729ecea7ae.pdf"
5 | ]

```

FIGURE 3 – Recherche CV contenant JAVA.

## 6 Conclusion

L'indexation des données est une pratique très efficace dans le domaine de recherche de l'information, elle nous permet entre autres de réduire notre espace de recherche, ce qui nous fait gagner énormément de temps, particulièrement quand on manipule une masse importante de données, c'est ce que nous propose aujourd'hui Elasticsearch, stocker, rechercher et analyser d'énormes volumes de données rapidement et en temps quasi réel pour retourner des réponses en quelques millisecondes.