

# ElasticSearch 附件內容搜尋 with Java

陳宗霆

# 起源



瀏覽記錄：

您好, 陳宗彥 (還未登入)

登入

- [首頁](#)
- [Wiki規範](#)
- [程設科合作管理會議](#)

## SA 專區

- [News](#)
- [SA WIKI專區](#)
- [Scrum經驗分享](#)

## 投資程式設計科

- [News](#)
- [新人專區](#)

歡迎 加入

Base: JSPWiki

站內文章搜尋

頁面內容 附件 (7) 頁面資訊

更多...

## Developer Wiki

<< 使用規範 >>

1. 未授權使用者可瀏覽：首頁、最新消息。  
另外開放以下頁面可供瀏覽：程式開發、投資系統、CSR系統、越南系統、上版系統、CM共用。
2. 下載附件時若發生無法下載的狀況，請改用firefox下載

## Hot News

2010-02-21

# 法務需求

- 希望搜尋公文、函示等文件內容  
(政府發布法條異動通常都是doc, pdf)

想找某個關鍵字的相关法令，都得一個一個點開檔案ctrl+F確認

# Idea



檔案



AAAAA.....  
Where am I  
....  
....  
....

字串

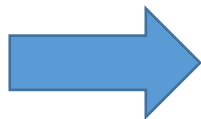


elasticsearch

# JSPWiki



檔案



AAAAA.....  
Where am I  
....  
....  
....

字串

透過Java poi,  
PDF Reader 等  
套件分別處理  
各種類檔案



elasticsearch

用Lucene而非  
ElasticSearch

# Idea



檔案



AAAAA.....  
Where am I  
....  
....  
....

字串



elasticsearch

改用Tika !!!

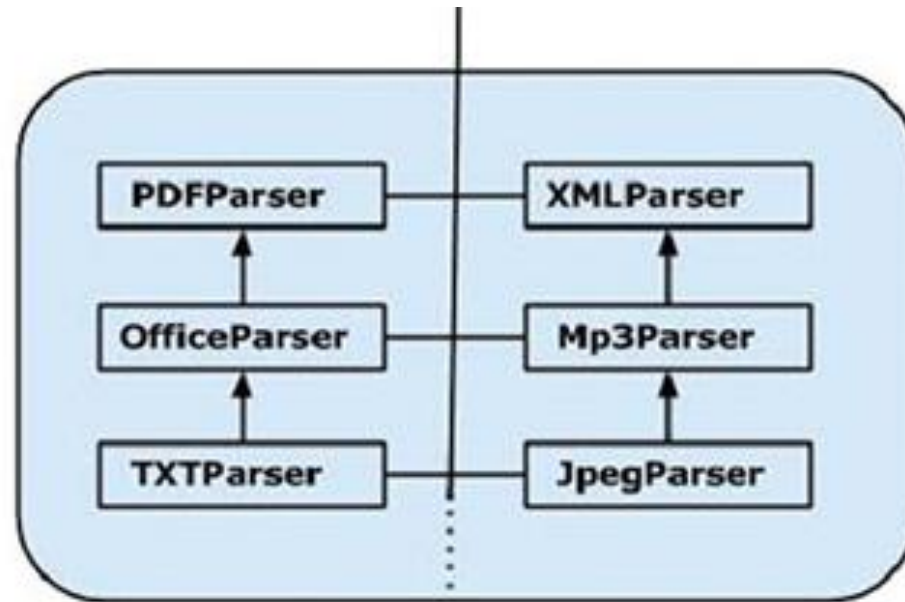
# Apache Tika

功能：自動檢視檔案類型，若有支援則解析出檔案資訊、內容

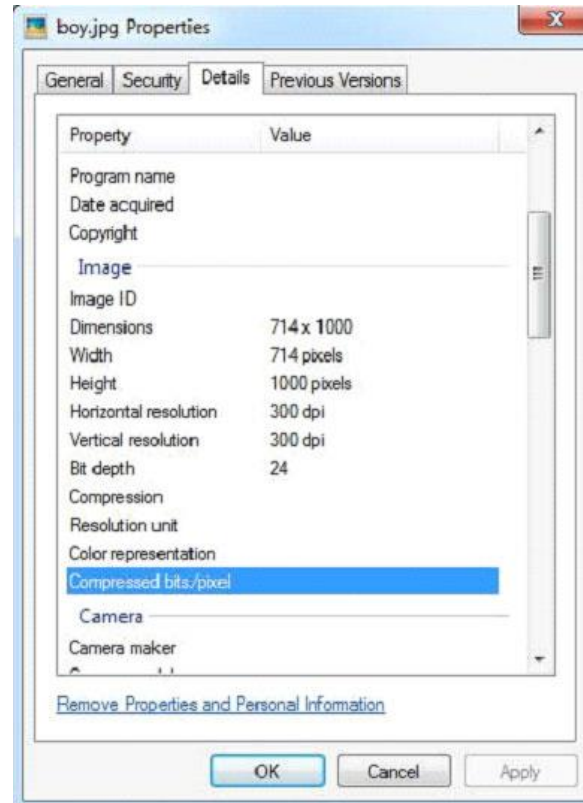
Input: File, FileStream

Output: String

Java Based



# Apache Tika



Contents of the document:

Metadata of the document:

IPTC-NAA record: 92 bytes binary data

Number of Components: 3

Image Height: 1000 pixels

Resolution Units: inch

Data Precision: 8 bits

tiff:BitsPerSample: 8

Compression Type: Baseline

Component 1: Y component: Quantization table 0, Sampling factors 1 horiz/1 vert

Component 2: Cb component: Quantization table 1, Sampling factors 1 horiz/1 vert

tiff:ImageLength: 1000

Component 3: Cr component: Quantization table 1, Sampling factors 1 horiz/1 vert

X Resolution: 300 dots

tiff:ImageWidth: 714

Application Record Version: 4

Image Width: 714 pixels

Original Transmission Reference:

53616c7465645f5fd22a84941585d89cc735d889c9d5ac58a01faf2c92ee3c6f9bcb38359bbe1eef

Y Resolution: 300 dots



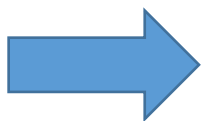
# Elastic Search Plugin

- <https://www.elastic.co/guide/en/elasticsearch/plugins/master/ingest-attachment.html>
- Take a rest & install it !!!
- Input: Base64 of File

# Idea



檔案



234gfdst34fse  
wbgdfgtDSFfd  
SDAF

Base64



AAAAA.....  
Where am I  
....  
....  
....

字串



elasticsearch

ingest attachment plugin

Let's start practice~~

# Environment

- Java 1.7 or higher
- UTF-8
- Library (given pom.xml)

-----

- Elasticsearch 5.6.10 or higher
- UTF-8

# Method

- Create Pipeline
- Put File to ElasticSearch
- Query
- Delete
- .....