

Text Extractions from Documents into Elasticsearch

Team: sddec22-19

Client: Buildertrend

Advisor: Goce Trajcevski

Team Members: Bruce Bitwayiki, Jared Hayashi,
Rushal Sohal, Tiffany Mayberry

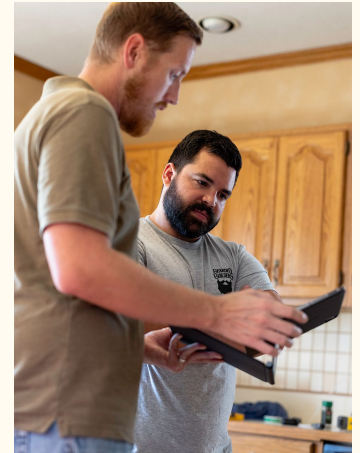
<http://sddec22-19.sd.ece.iastate.edu/>

Overview



The Client

- Industry-Leader in Construction Management Software
- An all-in-one platform: set a budget, communicate with clients, schedule work and keep tabs on job site progress
 - Presales, Project Management, Financial Tools, Customer Management
- 27,000+ construction companies use Buildertrend

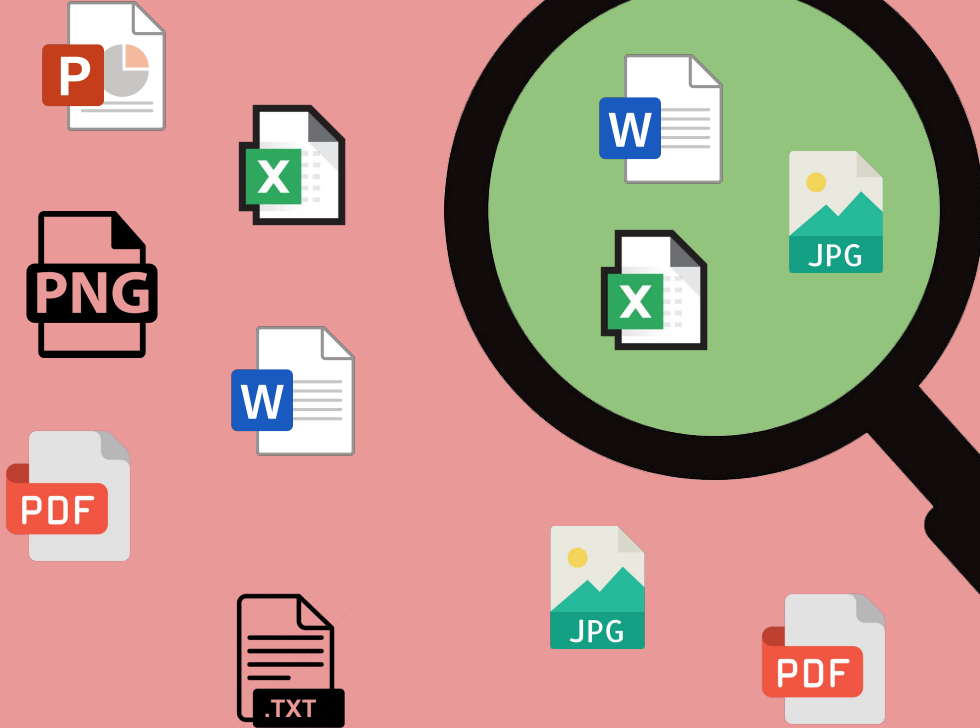


The Problem



Functional Requirements

ElasticSearch



Extract File Content & Metadata

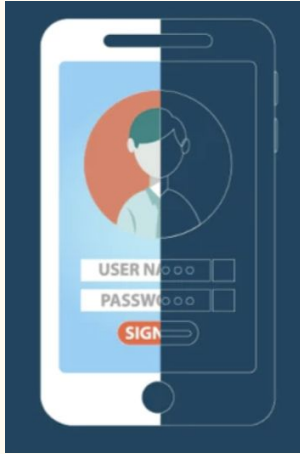


Search for Files Containing:
"5213 Pines Dr."

Filters:

Only search .docx,
.jpg, .xlsx files

Non-Functional Requirements



Easy to use UI



Return result within 10
seconds



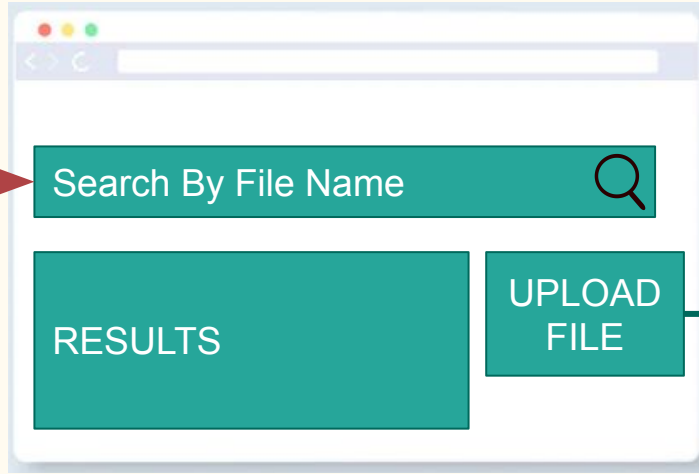
Easily modifiable and
scalable for future
functionality

Design



Buildertrends Current Solution

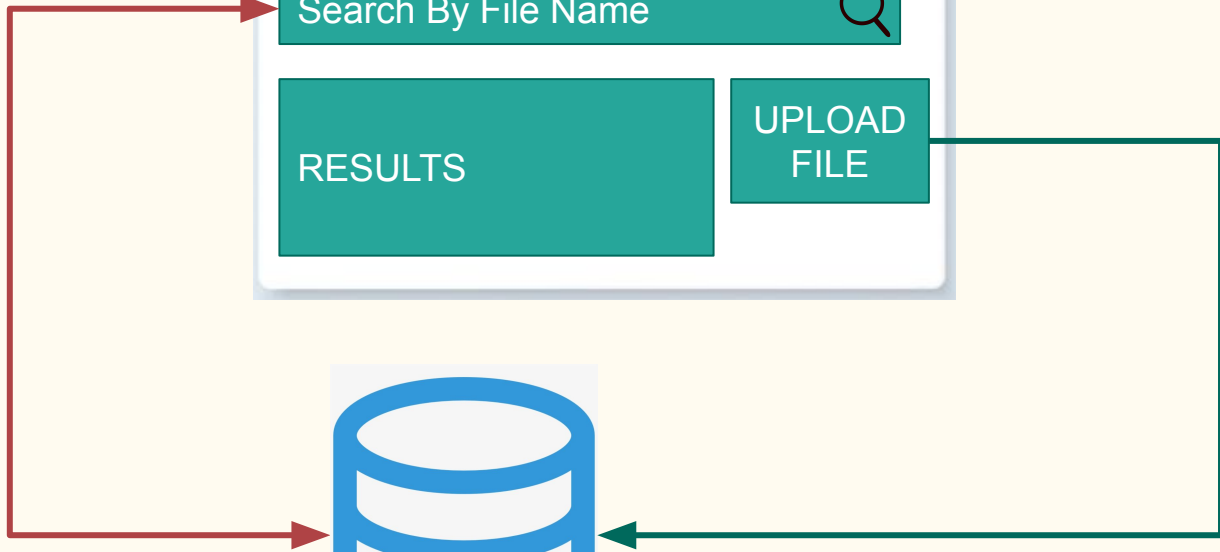
Fetch & Return
Results



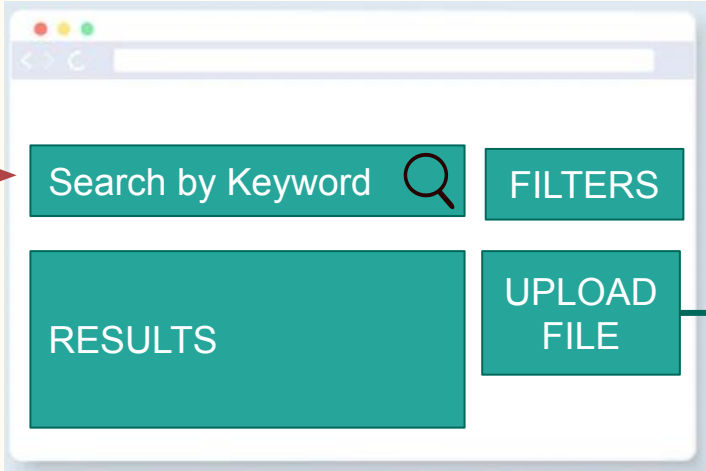
Upload
File



DATABASE



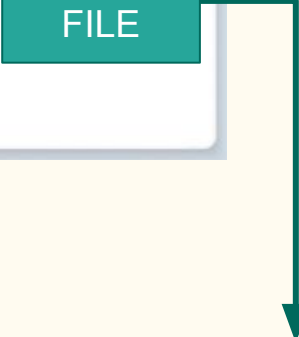
Basic Overview of Our Project









Fetch & Return Results



EXTRACT TEXT & INSERT INTO DATABASE



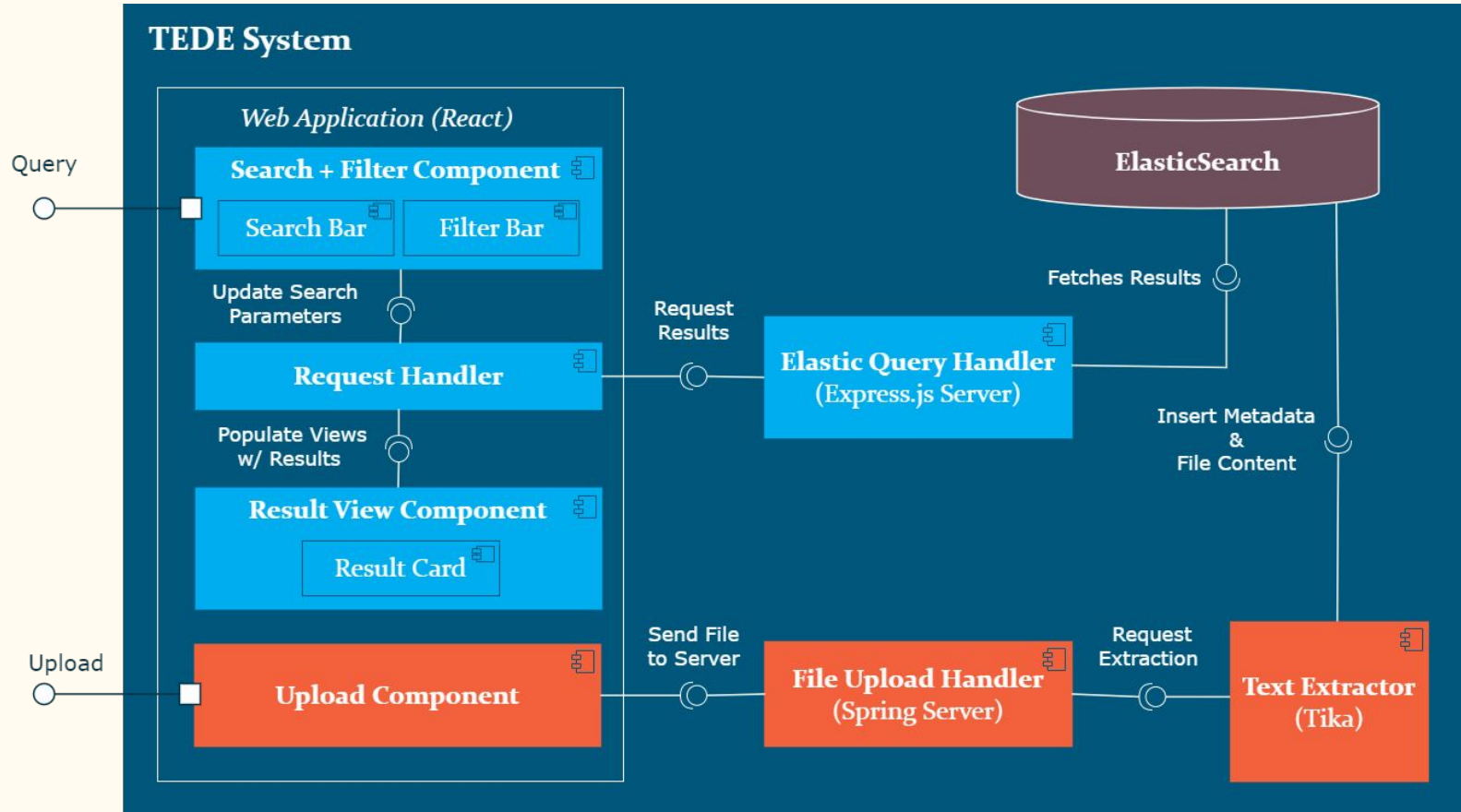
Technology Stack

Web Application	Database	Text Extractor
 	 	 

Implementation

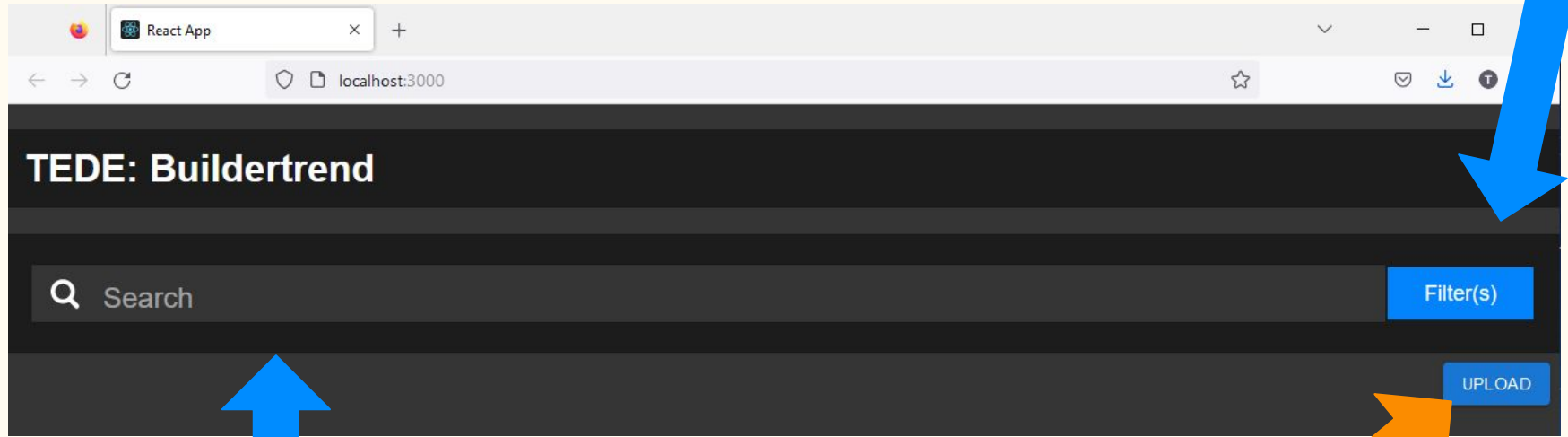
—

Component Diagram



Frontend Start Page

Filter Bar Component
Shows available filters



Search Bar Component
*Sends inputs to Request Handler
Component (on enter)*

Upload Component
*Gives user the ability to upload
supported file types*

Frontend Expanded Filter Bar

React App

localhost:3000

TEDE: Buildertrend

Search Filter(s)

Selected Filters: pdf docx Jonah /Documents/Jonah/2022

File Type(s): pdf txt docx pptx xlsx png jpg

Author: Jonah Add Clear

Path: /Documents/Jonah/2022 Add Clear

Upload Date: Add Clear

File Created Date: Add Clear

Last Updated Date: Add Clear

Match Case:





UPLOAD

Frontend Results View

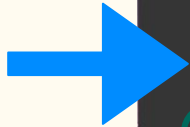
TEDE: Buildertrend

Q 123 Filter(s) UPLOAD

Found 4 Result(s) for: "123"

	really-big-contract.xlsx /Users/brucepro/Downloads
2 matches: keyword = '123' fileType = 'Microsoft Excel document (.xlsx)'	
	really-big-contract2.xlsx /Users/brucepro/Downloads
2 matches: keyword = '123' fileType = 'Microsoft Excel document (.xlsx)'	
	construction-contract.docx /Users/brucepro/Downloads
2 matches: keyword = '123' fileType = 'Microsoft Word document (.docx)'	
	construction-contract.pdf /Users/brucepro/Downloads
2 matches: keyword = '123' fileType = 'PDF (.pdf)'	

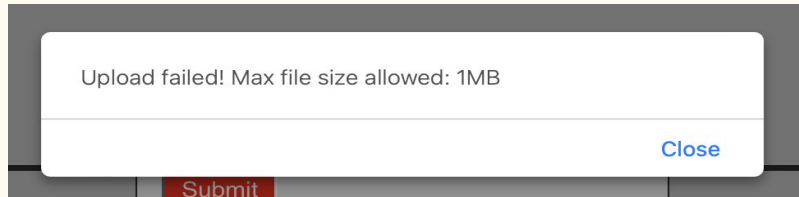
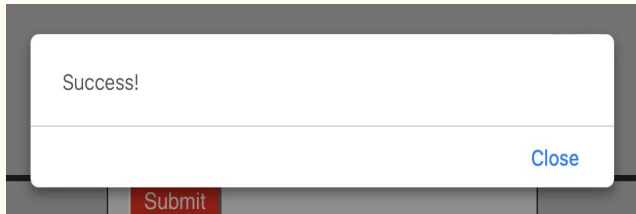
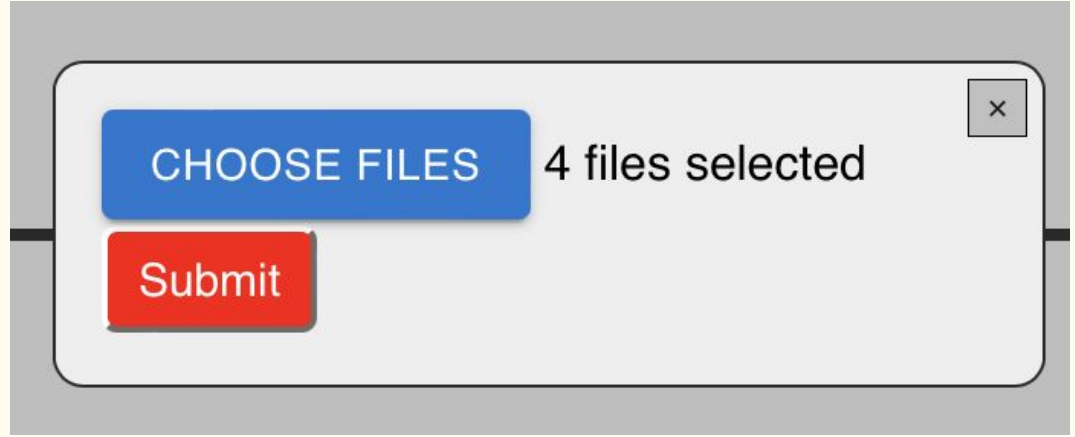
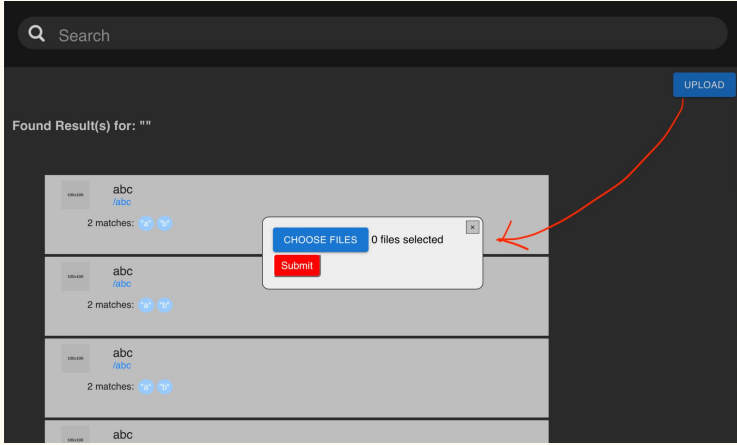
Result Card



Result View



Upload Component



Post Request

Elastic Query Handler

Post Response

```
1 {
2   ... "contents": "*contract*",
3   ... "types": [
4     ... "pdf",
5     ... "txt"
6   ],
7   ... "authors": [
8     ... "brucepro"
9   ],
10  ... "uploadDates": [],
11  ... "fileCreatedDate": [],
12  ... "lastUpdatedDate": [],
13  ... "containingFolder": []
14 }
15
```

Elasticsearch

```
1 {
2   ... "matches": [
3     ... {
4       ... "id": "1",
5       ... "fileName": "really-big-contract2.txt",
6       ... "filePath": "/Users/brucepro/Downloads",
7       ... "author": "brucepro",
8       ... "fileType": "Text (.txt)",
9       ... "contents": [
10        ... "This is another contract for XXX company."
11      ]
12    },
13    ... {
14      ... "id": "2",
15      ... "fileName": "another_sort_of_file.pdf",
16      ... "filePath": "/Documents",
17      ... "author": "brucepro",
18      ... "fileType": "PDF (.pdf)",
19      ... "contents": [
20        ... "THE_FILE_CONTENT. contract."
21      ]
22    }
23  ]
24 }
25
```

Post Request

Elastic Query Handler

Post Response

Elasticsearch Query

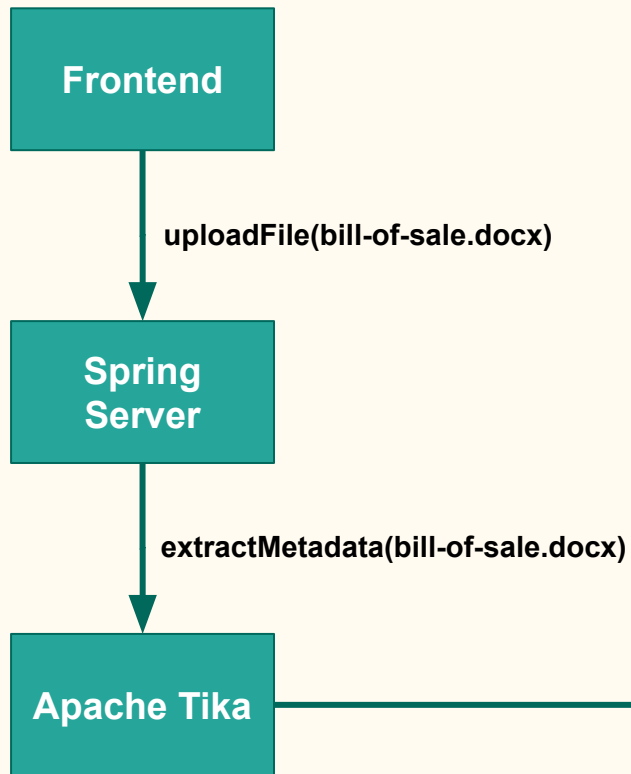
```
{
  index: 'documents',
  query: {
    bool: {
      "must": [
        { "query_string": {
          "default_field": "contents",
          "query": "*contract*"
        }}
      ],
      "filter": [
        { "bool": { "should": [
          { "match": { "type": "PDF (.pdf)" } },
          { "match": { "type": "Text (.txt)" } }
        ] } },
        { "bool": { "should": [
          { "match_phrase": { "owner": "brucepro" } }
        ] } }
      ]
    }
  }
}
```



Elasticsearch Response

```
{
  query: params.body.query,
  hits: { hits: [
    {
      _id: "1",
      _source: {
        "type": "Text (.txt)",
        "title": "really-big-contract2.txt",
        "path": "/Users/brucepro/Downloads",
        "owner": "brucepro",
        "contents": "This is another contract for XXX company."
      }
    },
    {
      _id: "2",
      _source: {
        "type": "PDF (.pdf)",
        "title": "another_sort_of_file.pdf",
        "path": "/Documents",
        "owner": "brucepro",
        "contents": "THE FILE_CONTENT. contract."
      }
    }
  ] }
}
```

Backend Example



JSON sent to Elasticsearch

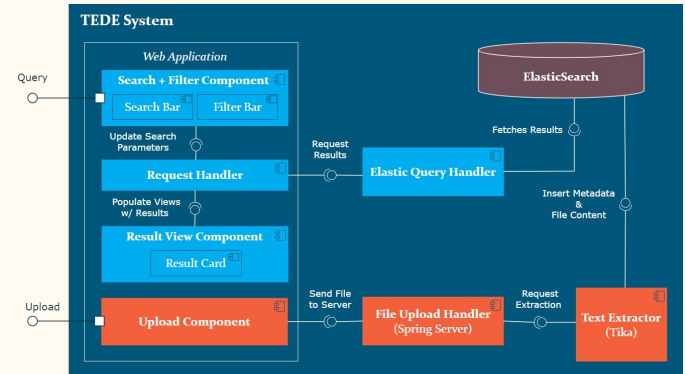
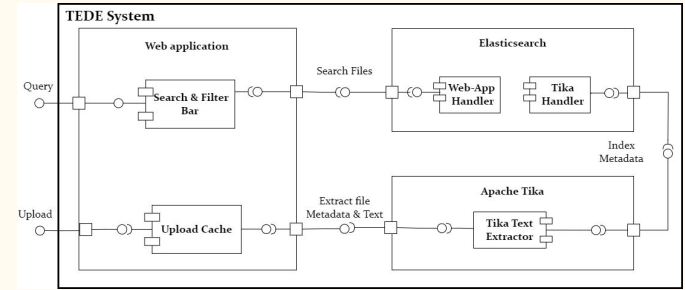
```
{
  "owner": "Microsoft Office User",
  "path": "/home/vm-user/sddec22-19/ExampleDocs/bill-of-sale.docx",
  "creator": "Microsoft Office User",
  "contents": "Bill of Sale THIS BILL OF SALE is executed this day of _____, by _____, (hereinafter Seller) residing at _____, for the benefit of _____ (hereinafter Buyer), residing at _____, located in _____ County, _____ For and in consideration of _____, which has been acknowledged to have been received by Seller. The form of payment used will be _____ and sales tax will not be included as part of the purchase price. The sale and transfer of property is made on an AS IS basis, without any express or implied warranties, with no recourse to the Seller, provided that Seller can issue proof that it has title to the property without any liens or encumbrances. The Buyer has been given the opportunity to inspect, or have inspected, any and all property as defined above. The Buyer agrees to accept all property in its existing state. In witness, the parties execute on this Bill of Sale on _____, Signature of Buyer _____, Signature of Seller _____, Date _____ Additional Notes: 1. Make sure that this Bill of Sale document is completed, and signed by both parties. Once signed, it will go into effect on the effective date specified in the document. 1. The Buyer should be provided with the original document, and a copy should be made and provided to the seller. 1. This document cannot be used to legally buy or sell real estate. ",
  "file-size": "21006 bytes",
  "word-count": "284",
  "title": "bill-of-sale.docx",
  "type": "docx",
  "page-count": "2"
}
```

```
1- {
2   "took" : 801,
3   "timed_out" : false,
4   "_shards" : {
5     "total" : 1,
6     "successful" : 1,
7     "skipped" : 0,
8     "failed" : 0
9   },
10  "hits" : {
11    "total" : {
12      "value" : 1,
13      "relation" : "eq"
14    },
15    "max_score" : 1.0,
16    "hits" : [
17      {
18        "_index" : "documents",
19        "_id" : "20e03493-90bc-49a1-8ecf-03dbf991d99d",
20        "_score" : 1.0,
21        "_ignored" : [
22          "contents.keyword"
23        ],
24        "_source" : {
25          "owner" : "Microsoft Office User",
26          "path" : "/home/vm-user/sddec22-19/ExampleDocs/bill-of-sale.docx",
27          "creator" : "Microsoft Office User",
28          "contents" : "Bill of Sale THIS BILL OF SALE is executed this day of _____, (hereinafter Seller) residing at _____, for the benefit of _____ (hereinafter Buyer), residing at _____, located in _____ County, _____ For and in consideration of _____, which has been acknowledged to have been received by Seller. The form of payment used will be _____ and sales tax will not be included as part of the purchase price. The sale and transfer of property is made on an AS IS basis, without any express or implied warranties, with no recourse to the Seller, provided that Seller can issue proof that it has title to the property without any liens or encumbrances. The Buyer has been given the opportunity to inspect, or have inspected, any and all property as defined above. The Buyer agrees to accept all property in its existing state. In witness, the parties execute on this Bill of Sale on _____, Signature of Buyer _____, Signature of Seller _____, Date _____ Additional Notes: 1. Make sure that this Bill of Sale document is completed, and signed by both parties. Once signed, it will go into effect on the effective date specified in the document. 1. The Buyer should be provided with the original document, and a copy should be made and provided to the seller. 1. This document cannot be used to legally buy or sell real estate. "
29        }
30      }
31    ]
32  }
33 }
```

Elasticsearch index

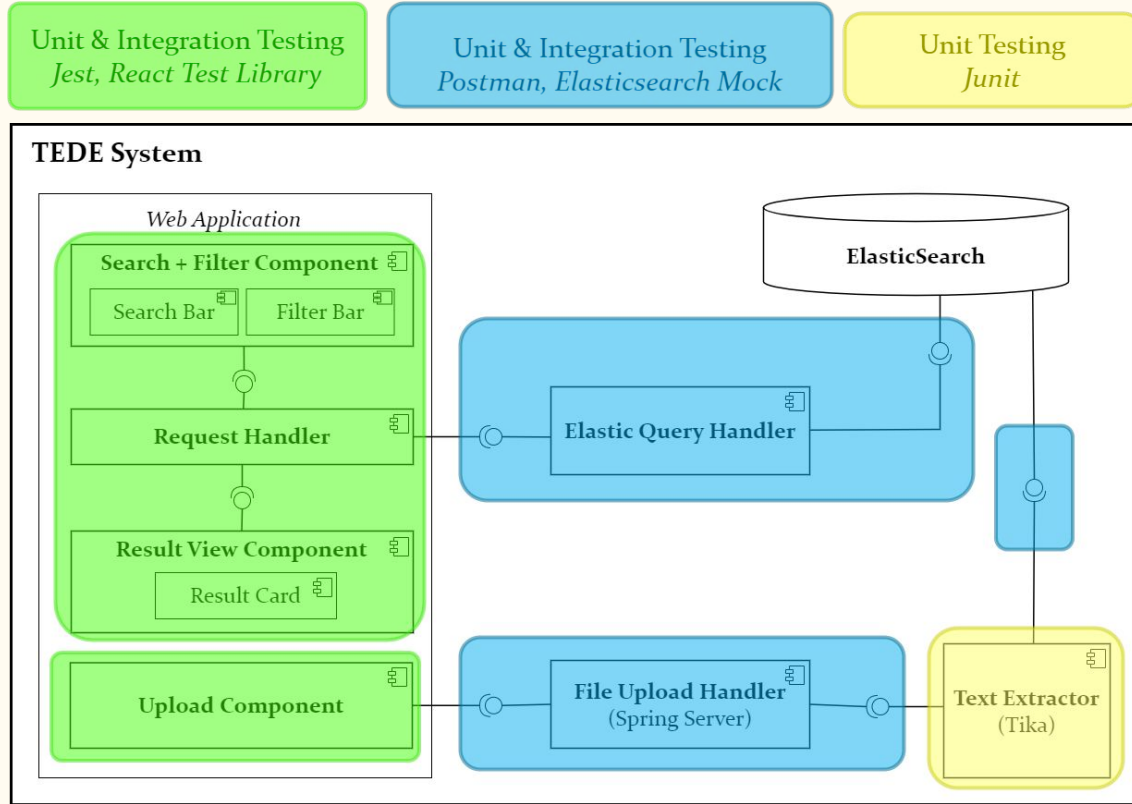
Implementation Challenges

- Elasticsearch authentication and communication
 - Certificate errors
 - Outdated documentation
- Architectural and requirement changes
 - Reallocating time to develop Tika server
 - Allocating more time for integration



Testing

Testing Plan



Acceptance Testing During BiWeekly Meetings with Client

Demo



Demo

Video

Questions ?



Extra Supporting Slides



Task Responsibility

FRONT-END

Tiffany Mayberry (SE)

Rushal Sohal (CPR E)

BACK-END

Jared Hayashi (SE)

Bruce Bitwayiki (CPR E)

Constraints and Considerations



Frontend should be a
web or desktop
application



Use **Elasticsearch** to
store the metadata



Application should **limit**
the returned results

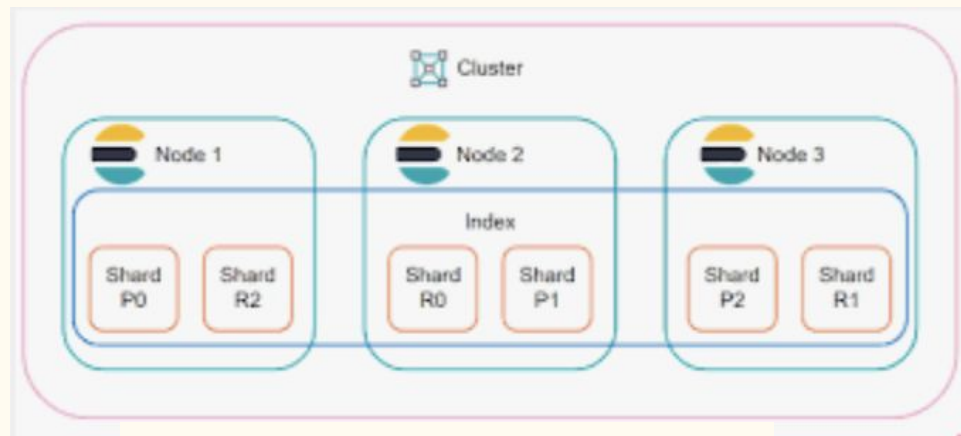
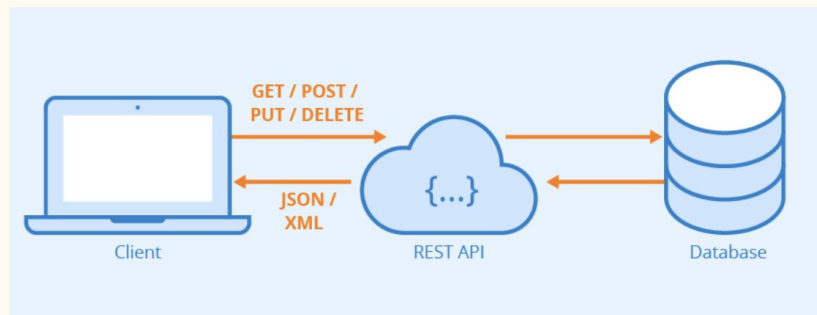
Elasticsearch

- **Free and open search and analytics engine**
- Based on the Lucene library
- Elastic Stack
- **At core - process JSON requests and return JSON data**
- **Index based searching**
- Structure based on documents
- Netflix, Ebay, Walmart, etc.



Risk Management & Mitigation

- Project Planning
- UI Integration
- REST API
- Elasticsearch Node reliability
- Text Extraction handler
- Testing



(⤴ Elasticsearch Structure)

Resource and Cost Estimate

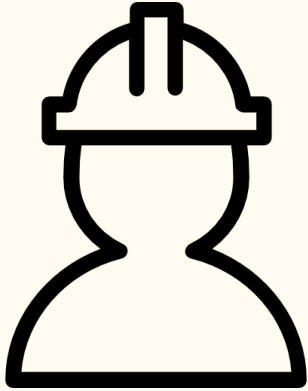


ETG Server



Open Source Packages

Users & Needs



Builders / Contractors
Access to
measurements, designs,
orders

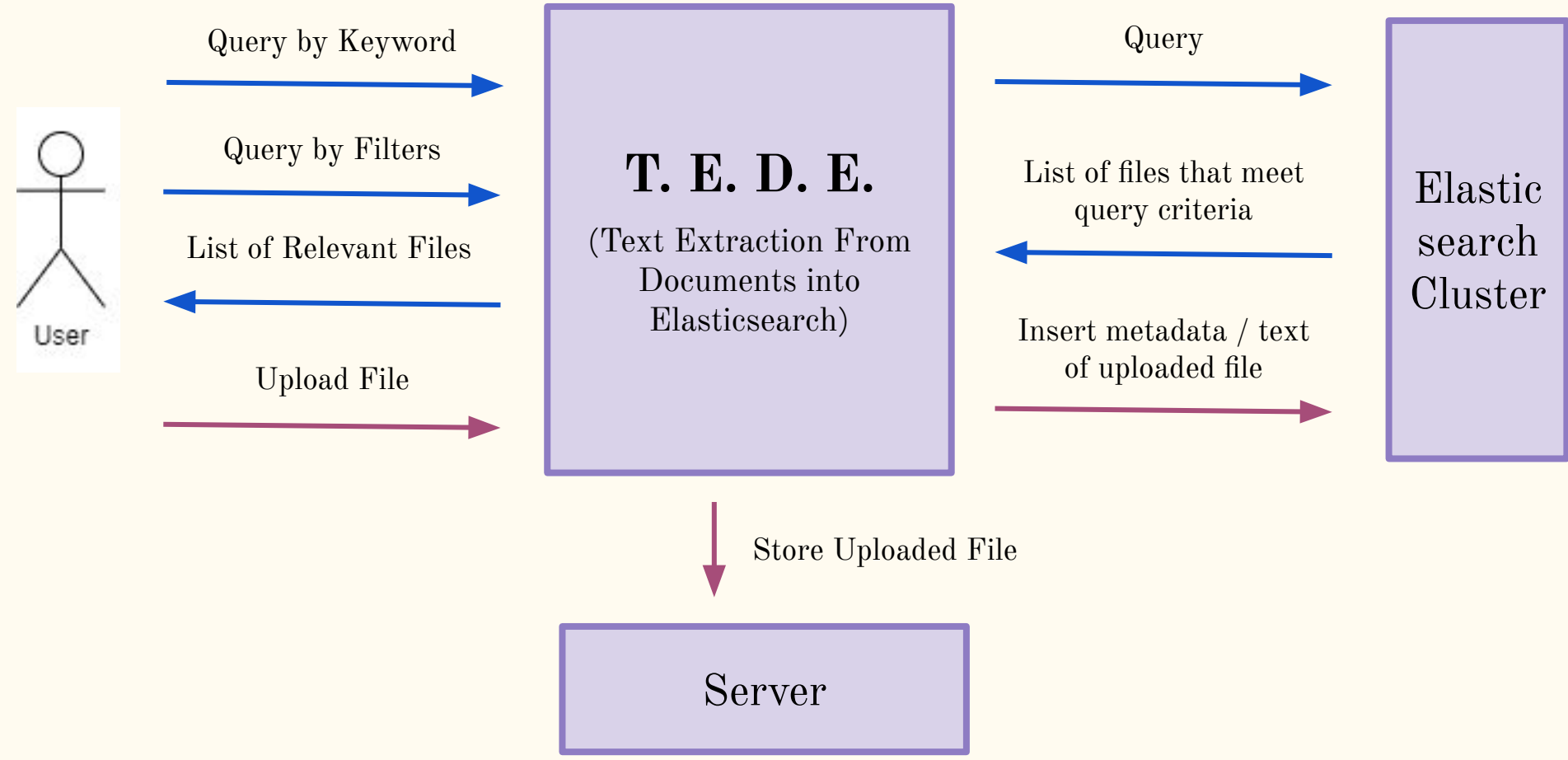


Homeowners
Contracts,
estimates/quotes,
specifications, notes



Buildertrend Staff
Support

Context Diagram



Original Testing Plan

Interface Testing
Postman

Unit Testing
Junit, Jest, React Test Library

Integration & System Testing

