# CE306- INFORMATION RETRIEVAL (REPORT)

**Student ID Number: 1403713/ 1402039**

CE306  Information Retrieval

# **Report Contents**

## Contents

## Instructions for running our system

The first step into the 'Indexing' aspect of our complete system was to identify a document collection of our choice to index. Therefore, we decided to go for a more sophisticated approach by indexing a website instead of the sample datasets from Elasticsearch. The URL we decided to index derived from a BBC website. This is because it is from a trusted source and have a range of content within the webpage. The queries devised for the test collection will be more effective for analysing if indexing from an efficient website.

To complete the indexing of a website, a crawler will need to be employed that gathers the data to pass over to Elasticsearch. We have attempted various crawlers and tools to work with Elasticsearch. For instance, we tested on Nutch, GNU wget tool and tapping into an RSS feed using Elasticsearch's on logstash. However, these did not give the results we want and were some issues in the process of scrapping. Therefore, we used another alternative by using 'HTTrack' to scrap our documents. In combination with Elasticsearch, the document was able to be indexed. This can be seen by the structed index format as he results.

The next stage of the process for the system will be to provide 'Searching' options for the user. From the successful indexing of the website and datasets, the user must be able to search it. We have decided to take an enhanced approach and not use the command line for this. Instead, we have used 'Kibana' to provide an interactive system for the user. The choice to use Kibana is because it provides visualisation of the data in Elasticsearch, suggesting better analysis of data.

The following stage is to 'Build a Test Collection', in which ten unique queries are devised with their expected results. This leads to the final stage of our system is 'Evaluation', this includes exploring different search engine settings to discover what effect they have on

the evaluation results. To accomplish this phase, we have identified a suitable metric. To provide the best approach for the evaluation results, we tried different retrieval models such as Boolean against TF.IDF. The numerical statistic of TFIDF provides us the reflection of how important a word is to a document in a collection or corpus. The value of tf-idf will increase proportionally to the number of times a word appears in the document. On the other hand, a Bool Query can be used that includes the use of AND/OR/NOT operators. This is used to fine tune the search queries to provide more relevant or specific results.

## Description of the document collection, Test Queries & Expected results:

The document collection we have decided to use from Elasticsearch are the complete works of William Shakespeare. This is because the Shakespeare collection is suitably parsed into fields, in which it is organised in the following schema:

```
{
   "line_id": INT,
   "play_name": "String",
   "speech_number": INT,
   "line_number": "String",
   "speaker": "String",
   "text_entry": "String",
}
```

Mappings for the fields were required before loading the Shakespeare and logs data sets. This allow the documents in the index to be sectioned into logical groups.

We have developed a variety of test queries to be evaluated from the 'Kibana User Interface'. By using a range of test queries, it could fully test the efficiency and effectiveness of varying information retrieval systems.

***For instance, from the 10-test queries we incorporated:***

- **Basic Match Query:**

```
GET /shakespeare/_search?
{
   "query": { "match": {"text_entry":"And borne her cleanly by the keepers nose?"} }
}
```

As the name suggests, a basic match query just searches the string that is provided in all the fields. In the test query below, if the text_entry field contains any of the values in 'And borne her cleanly by the keepers' then it will be displayed as the search results. The suggests the expected hits will be significantly high going to tens of thousands.

```
1 ▾  {
2      "took": 172,
3      "timed_out": false,
4 ▾    "_shards": {
5        "total": 5,
6        "successful": 5,
7        "failed": 0
8 ▲    },
9 ▾    "hits": {
10       "total": 46675,
11       "max_score": 37.00957,
12 ▾     "hits": [
13 ▾       {
14           "_index": "shakespeare",
15           "_type": "doc",
16           "_id": "97102",
17           "_score": 37.00957,
18 ▾         "_source": {
19             "type": "line",
20             "line_id": 97103,
21             "play_name": "Titus Andronicus",
22             "speech_number": 19,
23             "line_number": "2.1.99",
24             "speaker": "DEMETRIUS",
25             "text_entry": "And borne her cleanly by the keepers nose?"
26 ▲         }
27 ▲       },
28 ▾       {
29           "_index": "shakespeare",
30           "_type": "doc",
31           "_id": "58685",
32           "_score": 13.7950735,
33 ▾         "_source": {
34             "type": "line",
35             "line_id": 58686,
36             "play_name": "Measure for measure",
37             "speech_number": 5,
38             "line_number": "1.3.30",
39             "speaker": "DUKE VINCENTIO",
40             "text_entry": "And liberty plucks justice by the nose;"
41 ▲         }
42 ▲       },
43 ▾       {
44           "_index": "shakespeare",
45           "_type": "doc",
46           "_id": "90470",
47           "_score": 13.183935,
48 ▾         "_source": {
49             "type": "line",
50             "line_id": 90471,
51             "play_name": "Taming of the Shrew",
52             "speech_number": 66,
53             "line_number": "4.1.182",
54             "speaker": "PETRUCHIO",
55             "text_entry": "To make her come and know her keepers call,"
56 ▲         }
57 ▲       },
```

From the above test query, it shows there are 46675 hits which was what we expected to be tens of thousands. This test query is not a useful search as it allows far too much results.

- **Match Phrase Query:**

```
1   GET /_cat/indices?v
2
3   GET /shakespeare/_search?
4 ▾ {
5
6 ▾    "query": {
7
8 ▾      "match_phrase": {
9
10           "text_entry": "They are all welcome."
11
12 ▴      }
13
14 ▴    }
15
16 ▴  }
17
18 ▴ }
19
```

By using a 'Match Phrase' query, it will require that all the term in the query string be present in the document. As seen, in the above result the order specified in the query string needs to be exact and as close to each other. For instance, the match query we thought of was to search for text_entry of "They are all welcome". We expect for this query to return 1 hit that has the field 'text_entry' as 'They are all welcome'.

```
1 ▾ {
2    "took": 334,
3    "timed_out": false,
4 ▾  "_shards": {
5      "total": 5,
6      "successful": 5,
7      "failed": 0
8 ▴  },
9 ▾  "hits": {
10     "total": 1,
11     "max_score": 20.276226,
12 ▾   "hits": [
13 ▾     {
14         "_index": "shakespeare",
15         "_type": "doc",
16         "_id": "46970",
17         "_score": 20.276226,
18 ▾       "_source": {
19           "type": "line",
20           "line_id": 46971,
21           "play_name": "Julius Caesar",
22           "speech_number": 27,
23           "line_number": "2.1.102",
24           "speaker": "BRUTUS",
25           "text_entry": "They are all welcome."
26 ▴       }
27 ▴     }
28 ▴   ]
29 ▴ }
30 ▴ }
```

From the above test query, it shows there are 1 hit which was what we expected to be less than 5 hits. This test query is a useful search as it allows narrowing down of search results.

```
GET /shakespeare/_search?
{
  "query": { "match_phrase": {"text_entry":"To chase these pagans in those holy fields"} }
}
```

This is another type of test that makes use of 'Match Phrase' query. We expect the results in the text_entry field to contain 'To chase these pagans in those holy fields'. Since it requires to match exactly the full values of the field, the expected results will be limited to a few hits.

```
1  {
2    "took": 2,
3    "timed_out": false,
4    "_shards": {
5      "total": 5,
6      "successful": 5,
7      "failed": 0
8    },
9    "hits": {
10     "total": 1,
11     "max_score": 39.922577,
12     "hits": [
13       {
14         "_index": "shakespeare",
15         "_type": "doc",
16         "_id": "26",
17         "_score": 39.922577,
18         "_source": {
19           "type": "line",
20           "line_id": 27,
21           "play_name": "Henry IV",
22           "speech_number": 1,
23           "line_number": "1.1.24",
24           "speaker": "KING HENRY IV",
25           "text_entry": "To chase these pagans in those holy fields"
26         }
27       }
28     ]
29   }
30 }
```

As expected, this test query provides just 1 hit due to it requiring specific values from the 'test_entry' field.

- **Bool Query:**

```
1   POST /shakespeare/_search
2 ▾ {
3 ▾     "query": {
4 ▾         "bool": {
5 ▾             "must": {
6 ▾                 "bool" : { "should": [
7                         { "match": { "play_name": "Henry IV" }},
8                         { "match": { "play_name": "Julius Caesar" }} ] }
9 ▴             },
10            "must": { "match": { "speaker": "BRUTUS" }},
11            "must_not": { "match": {"speaker": "CASSIUS" }}
12 ▴         }
13 ▴     }
14 ▴ }
```

The use of Bool Query includes AND/OR/NOT operators and can be used to enhance the efficiency and effectiveness of our search queries. The use of Bool Query provides more relevant or specific results. For example, in the above test query we expected result to contain either Henry IV or Julius Caesar for the 'play_name' field. Furthermore, the speaker field should be 'BRUTUS' and cannot be 'CASSIUS'. The amount of hits will be extensive since many fields contain Henry IV and Julius Caesar as the values.

```
1 ▾ {
2     "took": 6,
3     "timed_out": false,
4 ▾   "_shards": {
5       "total": 5,
6       "successful": 5,
7       "failed": 0
8 ▴   },
9 ▾   "hits": {
10      "total": 833,
11      "max_score": 12.66049,
12 ▾    "hits": [
13 ▾      {
14          "_index": "shakespeare",
15          "_type": "doc",
16          "_id": "46370",
17          "_score": 12.66049,
18 ▾        "_source": {
19            "type": "line",
20            "line_id": 46371,
21            "play_name": "Julius Caesar",
22            "speech_number": 25,
23            "line_number": "1.2.32",
24            "speaker": "BRUTUS",
25            "text_entry": "I am not gamesome: I do lack some part"
26 ▴        }
27 ▴      },
28 ▾      {
29          "_index": "shakespeare",
30          "_type": "doc",
31          "_id": "46381",
32          "_score": 12.66049,
33 ▾        "_source": {
34            "type": "line",
35            "line_id": 46382,
36            "play_name": "Julius Caesar",
37            "speech_number": 27,
38            "line_number": "1.2.43",
39            "speaker": "BRUTUS",
40            "text_entry": "I turn the trouble of my countenance"
41 ▴        }
42 ▴      },
```

As expected, this test query provides many hits at 833 due to the broad range of Boolean values accepted.

```
POST /shakespeare/_search
{
   "query": {
      "bool": {
         "must": {
            "bool" : { "should": [
                   { "match": { "speaker": "FALSTAFF" }},
                   { "match": { "speaker": "WESTMORELAND" }} ] }
         },
         "must": { "match": { "text_entry": "lord of the council rated me the other day in the" }},
         "must_not": { "match": {"text_entry": "with vanity. I would to God thou and I knew where a" }}
      }
   }
}
```

This is another form of a 'Bool Query', which makes use of the 'OR' operator. From this query we expect the results to show either 'FALSTAFF' or 'WESTMORELAND' from the speaker field. Since, it also requires the 'text_entry' field to match a specific value the total hits would be low at around 200 hits.

```
1 ▾  {
2      "took": 56,
3      "timed_out": false,
4 ▾    "_shards": {
5        "total": 5,
6        "successful": 5,
7        "failed": 0
8 ▴    },
9 ▾    "hits": {
10       "total": 184,
11       "max_score": 36.11731,
12 ▾     "hits": [
13 ▾        {
14            "_index": "shakespeare",
15            "_type": "doc",
16            "_id": "192",
17            "_score": 36.11731,
18 ▾          "_source": {
19              "type": "line",
20              "line_id": 193,
21              "play_name": "Henry IV",
22              "speech_number": 27,
23              "line_number": "1.2.79",
24              "speaker": "FALSTAFF",
25              "text_entry": "lord of the council rated me the other day in the'
26 ▴          }
27 ▴        },
28 ▾        {
29            "_index": "shakespeare",
30            "_type": "doc",
31            "_id": "36995",
32            "_score": 17.664215,
33 ▾          "_source": {
34              "type": "line",
35              "line_id": 36996,
36              "play_name": "Henry V",
37              "speech_number": 19,
38              "line_number": "1.2.174",
39              "speaker": "WESTMORELAND",
40              "text_entry": "Playing the mouse in absence of the cat,"
41 ▴          }
42 ▴        },
43 ▾        {
44            "_index": "shakespeare",
45            "_type": "doc",
46            "_id": "54",
47            "_score": 16.046425,
48 ▾          "_source": {
49              "type": "line",
50              "line_id": 55,
51              "play_name": "Henry IV",
52              "speech_number": 4,
53              "line_number": "1.1.52",
54              "speaker": "WESTMORELAND",
55              "text_entry": "On Holy-rood day, the gallant Hotspur there,"
56 ▴          }
57 ▴        },
```

As expected, this test query provides not as much hits at 184 due to some specific Boolean values accepted.

```
POST /shakespeare/_search
{
    "query": {
        "bool": {
            "must": {
                "bool" : { "should": [
                        { "match": { "speaker": "FALSTAFF" }},
                        { "match": { "speaker": "PRINCE HENRY" }} ] }
            },
            "must": { "match": { "text_entry": "lord of the council rated me the other day in the" }},
            "must_not": { "match": {"text_entry": "with vanity. I would to God thou and I knew where a" }}
        }
    }
}
```

The above test query is like the previous one, in which the results will display the field of
speaker the values of either 'FALSTAFF' or 'PRINCE HENRY'. The expected hits should be
around 1000 as the values for the speaker field is quite common. Plus, the values for
text_entry to match can be any characters which narrows down the search by just a little.

```
1 ▾ {
2      "took": 9,
3      "timed_out": false,
4 ▾    "_shards": {
5          "total": 5,
6          "successful": 5,
7          "failed": 0
8 ▴    },
9 ▾    "hits": {
10         "total": 950,
11         "max_score": 36.11731,
12 ▾       "hits": [
13 ▾           {
14                 "_index": "shakespeare",
15                 "_type": "doc",
16                 "_id": "192",
17                 "_score": 36.11731,
18 ▾               "_source": {
19                     "type": "line",
20                     "line_id": 193,
21                     "play_name": "Henry IV",
22                     "speech_number": 27,
23                     "line_number": "1.2.79",
24                     "speaker": "FALSTAFF",
25                     "text_entry": "lord of the council rated me the other day in the"
26 ▴               }
27 ▴           },
28 ▾           {
29                 "_index": "shakespeare",
30                 "_type": "doc",
31                 "_id": "125",
32                 "_score": 20.611582,
33 ▾               "_source": {
34                     "type": "line",
35                     "line_id": 126,
36                     "play_name": "Henry IV",
37                     "speech_number": 2,
38                     "line_number": "1.2.12",
39                     "speaker": "PRINCE HENRY",
40                     "text_entry": "the time of the day."
41 ▴               }
42 ▴           },
43 ▾           {
44                 "_index": "shakespeare",
45                 "_type": "doc",
46                 "_id": "3176",
47                 "_score": 17.34315,
48 ▾               "_source": {
49                     "type": "line",
50                     "line_id": 3177,
51                     "play_name": "Henry IV",
52                     "speech_number": 4,
53                     "line_number": "5.5.18",
54                     "speaker": "PRINCE HENRY",
55                     "text_entry": "The fortune of the day quite turnd from him,"
56 ▴               }
57 ▴           },
```

As expected, this test query provides around 1000 hits at 950 due to some common values
used for searching.

```
POST /shakespeare/_search
{
    "query": {
        "bool": {
            "must":{"match":{"speaker":"DEMETRIUS"}},
            "must":{"match":{"line_number":"2.1.99"}},
            "must":{"match":{"play_name":"Titus Andronicus"}}
        }
    }
}
```

From this test query, it is very specific in what values it must match. Therefore, it is expected the results to less than 5 hits or just 1. This query requires the fields of 'speaker', 'line_number' and 'play_name' to be exact values.

```
1 ▾ {
2      "took": 2,
3      "timed_out": false,
4 ▾    "_shards": {
5        "total": 5,
6        "successful": 5,
7        "failed": 0
8 ▴    },
9 ▾    "hits": {
10       "total": 1,
11       "max_score": 22.289654,
12 ▾     "hits": [
13 ▾       {
14           "_index": "shakespeare",
15           "_type": "doc",
16           "_id": "97102",
17           "_score": 22.289654,
18 ▾         "_source": {
19             "type": "line",
20             "line_id": 97103,
21             "play_name": "Titus Andronicus",
22             "speech_number": 19,
23             "line_number": "2.1.99",
24             "speaker": "DEMETRIUS",
25             "text_entry": "And borne her cleanly by the keepers nose?"
26 ▴         }
27 ▴       }
28 ▴     ]
29 ▴   }
30 ▴ }
```

As expected, this test query provides very few hits with 1 only. This is due to it requiring very specific results as it must match certain values for the fields.

- **Regexp Query:**

```
POST /shakespeare/_search
{
    "query": {
        "regexp" : {
            "speaker" : "b[a-z]*s"
        }
    },
    "_source": ["line_number", "speaker"],
    "highlight": {
        "fields" : {
            "speaker" : {}
        }
    }
}
```

The use of Regexp queries ensures that more complex patterns are specified in comparison to wildcard queries. Instead of using '?' to specify the remaining letters of the query, it uses '[a-z]' to specify any characters that can occur between the starting letter of 'b' and the ending letter of 's'. Therefore, in this test query the expected result would be anything that has in the speaker field beginning with 'b' and ending in 's'. In this case the output should have in the speaker field 'BRUTUS'. This suggests hits will be around 1000 to 2000.

```
1 ▾ {
2     "took": 30,
3     "timed_out": false,
4 ▾   "_shards": {
5         "total": 5,
6         "successful": 5,
7         "failed": 0
8 ▴   },
9 ▾   "hits": {
10        "total": 1586,
11        "max_score": 1,
12 ▾      "hits": [
13 ▾          {
14                 "_index": "shakespeare",
15                 "_type": "doc",
16                 "_id": "8483",
17                 "_score": 1,
18 ▾              "_source": {
19                     "line_number": "4.2.5",
20                     "speaker": "BEVIS"
21 ▴              },
22 ▾              "highlight": {
23 ▾                  "speaker": [
24                         "<em>BEVIS</em>"
25 ▴                  ]
26 ▴              }
27 ▴          },
28 ▾          {
29                 "_index": "shakespeare",
30                 "_type": "doc",
31                 "_id": "8494",
32                 "_score": 1,
33 ▾              "_source": {
34                     "line_number": "4.2.16",
35                     "speaker": "BEVIS"
36 ▴              },
37 ▾              "highlight": {
38 ▾                  "speaker": [
39                         "<em>BEVIS</em>"
40 ▴                  ]
41 ▴              }
42 ▴          },
43 ▾          {
44                 "_index": "shakespeare",
45                 "_type": "doc",
46                 "_id": "8500",
47                 "_score": 1,
48 ▾              "_source": {
49                     "line_number": "4.2.22",
50                     "speaker": "BEVIS"
51 ▴              },
52 ▾              "highlight": {
53 ▾                  "speaker": [
54                         "<em>BEVIS</em>"
55 ▴                  ]
56 ▴              }
57 ▴          },
```

As expected, this test query provides hits into thousands category with 1586. This is due to the flexibility in the values that are allowed in the query.

- **Match Phrase Prefix:**

```
POST /shakespeare/_search
{
    "query": {
        "match_phrase_prefix" : {
            "text_entry": {
                "query": "so prou",
                "slop": 3,
                "max_expansions": 10
            }
        }
    },
    "_source": [  "line_id", "play_name", "text_entry" ]
}
```

By using 'Match Phrase Prefix', the queries functions like an autocomplete system, in which the user can complete the search without inputting the whole query. In the test query above, it autocompletes 'prou' that is most likely to be 'proud'. It also makes use of 'max_expansions' that is set to 10. The idea of this is that the auto-complete field is limited 10 characters.

```
1 ▾ {
2     "took": 4,
3     "timed_out": false,
4 ▾   "_shards": {
5         "total": 5,
6         "successful": 5,
7         "failed": 0
8 ▴   },
9 ▾   "hits": {
10        "total": 7,
11        "max_score": 47.688564,
12 ▾      "hits": [
13 ▾        {
14              "_index": "shakespeare",
15              "_type": "doc",
16              "_id": "67460",
17              "_score": 47.688564,
18 ▾            "_source": {
19                  "play_name": "A Midsummer nights dream",
20                  "text_entry": "Have every pelting river made so proud",
21                  "line_id": 67461
22 ▴            }
23 ▴        },
24 ▾        {
25              "_index": "shakespeare",
26              "_type": "doc",
27              "_id": "24755",
28              "_score": 40.532585,
29 ▾            "_source": {
30                  "play_name": "Coriolanus",
31                  "text_entry": "Was ever man so proud as is this Marcius?",
32                  "line_id": 24756
33 ▴            }
34 ▴        },
35 ▾        {
36              "_index": "shakespeare",
37              "_type": "doc",
38              "_id": "11612",
39              "_score": 40.033813,
40 ▾            "_source": {
41                  "play_name": "Henry VI Part 3",
42                  "text_entry": "Ha! durst the traitor breathe out so proud words?",
43                  "line_id": 11613
44 ▴            }
45 ▴        },
46 ▾        {
47              "_index": "shakespeare",
48              "_type": "doc",
49              "_id": "12223",
50              "_score": 40.033813,
51 ▾            "_source": {
52                  "play_name": "Henry VI Part 3",
53                  "text_entry": "And so, proud-hearted Warwick, I defy thee,",
54                  "line_id": 12224
55 ▴            }
56 ▴        },
```

As expected, this test query provides few hits with only 7. This is due it autocompletes specific values of 'proud'.

- **Match Phrase Query:**

```
POST /shakespeare/_search
{
    "query": {
        "multi_match" : {
            "query": "gentleman",
            "fields": ["text_entry", "speaker"],
            "type": "phrase",
            "slop": 2
        }
    },
    "_source": [ "text_entry", "speaker" ]
}
```

For this test query, it makes use of match phrase query with the requirement of having 'gentleman' in either the fields of 'text_entry' or 'speaker'. It considers the order specified in the query string and needs to be close to each other. The use of the field of 'slop' indicates how far apart terms can be while still considering the match. The expected hits for this test query should be around 1000 hits as it's not a specific search.

```
1   {
2     "took": 1,
3     "timed_out": false,
4     "_shards": {
5       "total": 5,
6       "successful": 5,
7       "failed": 0
8     },
9     "hits": {
10      "total": 1097,
11      "max_score": 9.256456,
12      "hits": [
13        {
14          "_index": "shakespeare",
15          "_type": "doc",
16          "_id": "103318",
17          "_score": 9.256456,
18          "_source": {
19            "text_entry": "A gentleman! what gentleman?",
20            "speaker": "OLIVIA"
21          }
22        },
23        {
24          "_index": "shakespeare",
25          "_type": "doc",
26          "_id": "16144",
27          "_score": 8.95748,
28          "_source": {
29            "text_entry": "Gentleman,",
30            "speaker": "ROSALIND"
31          }
32        },
33        {
34          "_index": "shakespeare",
35          "_type": "doc",
36          "_id": "27250",
37          "_score": 8.165185,
38          "_source": {
39            "text_entry": "A gentleman.",
40            "speaker": "CORIOLANUS"
41          }
42        },
43        {
44          "_index": "shakespeare",
45          "_type": "doc",
46          "_id": "52667",
47          "_score": 8.165185,
48          "_source": {
49            "text_entry": "Exit Gentleman",
50            "speaker": "ALBANY"
51          }
52        },
53        {
54          "_index": "shakespeare",
55          "_type": "doc",
56          "_id": "29550",
57          "_score": 8.003025,
58          "_source": {
59            "text_entry": "A gentleman.",
60            "speaker": "CLOTEN"
61          }
62        },
```

As expected, this test query provides around a thousand hits and it was 1097. This is due to the effectiveness of using multi match for the search.

```
POST /shakespeare/_search
{
    "query": {
        "multi_match" : {
            "query": "demetrius",
            "fields": ["text_entry", "speaker"],
            "type": "phrase",
            "slop": 2
        }
    },
    "_source": [ "text_entry", "speaker" ]
}
```

For this test query, it is like the previous one but this time it allows 'demetrius' to appear either in the 'text_entry' or 'speaker' fields. The type of this query is also set to 'phrase'.

```
1 ▼ {
2     "took": 2,
3     "timed_out": false,
4 ▼   "_shards": {
5         "total": 5,
6         "successful": 5,
7         "failed": 0
8 ▲   },
9 ▼   "hits": {
10        "total": 15,
11        "max_score": 16.21323,
12 ▼      "hits": [
13 ▼          {
14                "_index": "shakespeare",
15                "_type": "doc",
16                "_id": "67618",
17                "_score": 16.21323,
18 ▼              "_source": {
19                    "type": "line",
20                    "line_id": 67619,
21                    "play_name": "A Midsummer nights dream",
22                    "speech_number": 35,
23                    "line_number": "",
24                    "speaker": "HELENA",
25                    "text_entry": "Exit DEMETRIUS"
26 ▲              }
27 ▲          },
28 ▼          {
29                "_index": "shakespeare",
30                "_type": "doc",
31                "_id": "68374",
32                "_score": 15.26428,
33 ▼              "_source": {
34                    "type": "line",
35                    "line_id": 68375,
36                    "play_name": "A Midsummer nights dream",
37                    "speech_number": 78,
38                    "line_number": "3.2.334",
39                    "speaker": "HELENA",
40                    "text_entry": "With Demetrius."
41 ▲              }
42 ▲          },
43 ▼          {
44                "_index": "shakespeare",
45                "_type": "doc",
46                "_id": "67756",
47                "_score": 13.69783,
48 ▼              "_source": {
49                    "type": "line",
50                    "line_id": 67757,
51                    "play_name": "A Midsummer nights dream",
52                    "speech_number": 17,
53                    "line_number": "2.2.96",
54                    "speaker": "HELENA",
55                    "text_entry": "Therefore no marvel though Demetrius"
56 ▲              }
57 ▲          },
```

As expected, this test query provides few hits at 15. This is due to it requiring the specific value of 'demetrius' in the 'text_entry' field.

# Design decisions of overall architecture

We used Kibana as our user interface because it provides a user-friendly interface and has great compatibility. Kibana provides advantages such as allowing visualisation of our Elasticsearch data and is easy efficient to use. It also has a dashboard that includes essential data such as the different requests running.

We also made the decision of using Elasticsearch because it is one of the most popular and available open source search engine.

# How we attempted to Scrape and Crawl a website

Scrapper: HTTrack
http://www.httrack.com/
Install current machine version from download page.
Visit: http://www.httrack.com/html/shelldoc.html for Windows or Linux system set up guide.
In our version we used Linux ubuntu distribution and were running from terminal window.
$ webttrack
Once that done open localhost:8080 in your web browser to see actual HTTrack web plugin and follow previously mentioned guide for set up.
After you extracted your website locate HTML file.
Visit GitHub at https://github.com/scraperwiki/elasticsearch-tool and download tar.gz
Extract it where you like, elasticsearch-tool has more functions, but as right now we only going to use httrack_html_to_json.py file.
Once you have HTML file from HTTrack and elasticsearch-tool folder from github, move or copy HTML file to elasticsearch-tool folder and run httrack_html_to_json.py in terminal window.
$ python httrack_html_to_json.py website.html
This will give JSON file format with Text only from website
JSON has few attributes like "title", "url", "scrape_date", and "body". "Body" is where all the text is held.
Having only these few attributes we decided that it would be insufficient to have such low number of attributes to have complex queries.

# Evaluation results

```
POST /shakespeare/_search
{
    "query": {
            "bool" : { "should": [
                    { "match": { "text_entry": "Demetrius" }},
                    { "match": { "text_entry": "Lysander" }} ] }
        }}
```

This is an example of a Boolean Query that searches 'Demetrius' and 'Lysander' as the values for the field of 'text_entry'. Both the values could appear as single text_entry or together. This suggests if they both appear in the same query then the score would be higher.

```
1▼  {
2      "took": 1,
3      "timed_out": false,
4▼     "_shards": {
5          "total": 5,
6          "successful": 5,
7          "failed": 0
8▲     },
9▼     "hits": {
10         "total": 6,
11         "max_score": 14.988718,
12▼        "hits": [
13▼            {
14                 "_index": "shakespeare",
15                 "_type": "doc",
16                 "_id": "68864",
17                 "_score": 14.988718,
18▼                "_source": {
19                     "text_entry": "Enter LYSANDER, DEMETRIUS, HERMIA, and HELENA",
20                     "speaker": "THESEUS"
21▲                }
22▲            },
23▼            {
24                 "_index": "shakespeare",
25                 "_type": "doc",
26                 "_id": "68396",
27                 "_score": 13.173532,
28▼                "_source": {
29                     "text_entry": "Exeunt LYSANDER and DEMETRIUS",
30                     "speaker": "DEMETRIUS"
31▲                }
32▲            },
33▼            {
34                 "_index": "shakespeare",
35                 "_type": "doc",
36                 "_id": "68547",
37                 "_score": 12.744417,
38▼                "_source": {
39                     "text_entry": "SCENE I. The same. LYSANDER, DEMETRIUS, HELENA, and HERMIA",
40                     "speaker": "PUCK"
41▲                }
42▲            },
43▼            {
44                 "_index": "shakespeare",
45                 "_type": "doc",
46                 "_id": "68698",
47                 "_score": 10.422862,
48▼                "_source": {
49                     "text_entry": "Horns and shout within. LYSANDER, DEMETRIUS, HELENA, and HERMIA wake and start
   up",
50                     "speaker": "THESEUS"
51▲                }
52▲            },
```

The above shows the results of Boolean query that have 6 hits. Based on this query there were relevant documents at position 1, 3 and 6 with our average precision at 0.722, with the other documents being irrelevant.

TF.IDF search query

```
POST /shakespeare/_search
{
    "query": {
        "multi_match" : {
            "query": "lysander demetrius",
            "fields": ["text_entry", "speaker"],
            "type": "phrase",
            "slop": 2
        }
    },
    "_source": [ "text_entry", "speaker" ]
}
```

This query searches if 'lysander demetrius' appears in the fields of 'text_entry' and 'speaker'.

The results for this specific query are as follows:

```
1 ▼ {
2     "took": 1,
3     "timed_out": false,
4 ▼   "_shards": {
5         "total": 5,
6         "successful": 5,
7         "failed": 0
8 ▲   },
9 ▼   "hits": {
10        "total": 115,
11        "max_score": 18.99001,
12 ▼      "hits": [
13 ▼          {
14                "_index": "shakespeare",
15                "_type": "doc",
16                "_id": "68396",
17                "_score": 18.99001,
18 ▼              "_source": {
19                    "type": "line",
20                    "line_id": 68397,
21                    "play_name": "A Midsummer nights dream",
22                    "speech_number": 86,
23                    "line_number": "",
24                    "speaker": "DEMETRIUS",
25                    "text_entry": "Exeunt LYSANDER and DEMETRIUS"
26 ▲              }
27 ▲          },
28 ▼          {
29                "_index": "shakespeare",
30                "_type": "doc",
31                "_id": "67020",
32                "_score": 15.270048,
33 ▼              "_source": {
34                    "type": "line",
35                    "line_id": 67021,
36                    "play_name": "A Midsummer nights dream",
37                    "speech_number": 3,
38                    "line_number": "",
39                    "speaker": "THESEUS",
40                    "text_entry": "Enter EGEUS, HERMIA, LYSANDER, and DEMETRIUS"
41 ▲              }
42 ▲          },
43 ▼          {
44                "_index": "shakespeare",
45                "_type": "doc",
46                "_id": "68688",
47                "_score": 14.988719,
48 ▼              "_source": {
49                    "type": "line",
50                    "line_id": 68689,
51                    "play_name": "A Midsummer nights dream",
52                    "speech_number": 32,
53                    "line_number": "4.1.132",
54                    "speaker": "EGEUS",
55                    "text_entry": "And this, Lysander; this Demetrius is;"
56 ▲              }
57 ▲          },
```

The above shows the results of TF.IDF query that have 115 hits. Based on this query there were relevant documents at position 2 and 6 with our average precision at 0.416, with the other documents being irrelevant.

We decided to use Mean Average Precision (MAP), because it will allow us to compare between different information retrieval models. The chosen information retrieval models were TF.IDF and Boolean. Above search queries resulted in MAP of 0.569, which suggests it is quite accurate. The ideal result would be 1, but there will be no information retrieval model that are perfect.

## Engineering a Complete System

Complete system is created in tmp folder on Ubuntu operating system for simplicity.

To run script open terminal window inside file existing folder and type $ ./setupElastic

- ➢ Creating directory in tmp folder
- ➢ Copies all needed script and json dataset files
- ➢ Moves terminal to created folder at step 1
- ➢ Starts downloading all necessary files like, elasticsearch, kibana and java environment
- ➢ Next is unpacks everything in folder
- ➢ Java home is set for elastic search to use
- ➢ Posting indexed Json file into elastic search indices
- ➢ Checks if dataset is uploaded
- ➢ New terminal window is opened for kibana launch

```
1   #!/bin/dash
2   echo "Creating Elastic Search ENV"
3
4   mkdir /tmp/elastic_kibana
5   echo "/tmp/elastic_kibana CREATED"
6   cp ./setupElastic /tmp/elastic_kibana
7   cp ./runKibana /tmp/elastic_kibana
8   cp ./shakespeare_6.0.json /tmp/elastic_kibana
9   echo "Files copied to TMP Folder"
10  cd /tmp/elastic_kibana
11  echo "Moving to TMP folder to start download"
12  curl -L -O https://artifacts.elastic.co/downloads/elasticsearch/elasticsearch-5.1.2.tar.gz
13  curl -L -o jre8.tgz http://javadl.oracle.com/webapps/download/AutoDL?BundleId=216422
14  curl -L -O https://artifacts.elastic.co/downloads/kibana/kibana-5.0.0-linux-x86_64.tar.gz
15  echo "Unpacking downloaded files"
16  tar xf elasticsearch-5.1.2.tar.gz
17  tar xf jre8.tgz
18  tar xf kibana-5.0.0-linux-x86_64.tar.gz
19  echo "Setting up Java HOME for Elastic Search database"
20  JAVA_HOME=/tmp/elastic_kibana/jre1.8.0_111 PATH=$JAVA_HOME/bin:$PATH elasticsearch-5.1.2/bin/elasticsearch -d
21  echo "First Error is normal to set up local host....AFTER ERROR wait 15 seconds"
22  curl -XPOST 'localhost:9200/shakespeare/doc/_bulk?pretty' --data-binary @shakespeare_6.0.json
23  sleep 15s
24  curl -XPOST 'localhost:9200/shakespeare/doc/_bulk?pretty' --data-binary @shakespeare_6.0.json
25  echo "CHECKING indexed documents"
26  curl 'localhost:9200/_cat/indices?v'
27
28  echo "All OK, now Kibana Launching"
29
30  gnome-terminal -x ./runKibana
31
```

The screen capture below shows how Kibana is set up.

- ➢ Java Home environment is set for Kibana
- ➢ Wait 20 sec for default browser to open local host
- ➢ Start writing queries

```
1   #!/bin/dash
2   cd /tmp/elastic_kibana
3   echo "JAVA HOME setting up for kibana"
4   echo "wait up to 20 sec to OPEN BROWSER window with kibana running"
5   JAVA_HOME=/tmp/elasticAssic/jre1.8.0_111 PATH=$JAVA_HOME/bin:$PATH kibana-5.0.0-linux-x86_64/bin/kibana &
6   sleep 20s
7   sensible-browser http://localhost:5601/app/kibana#/dev_tools
8   echo "Just wait a second pleaaase :) opening"
9   bash
10
```

## Possible systems improvements & Extensions

The system can be improved with Logstash extensions for easy web scrapping integration into Elasticsearch. Enhancements can be made for ranked evaluation using the latest version of Elasticsearch 6.2. in addition, we could incorporate more search queries with different parameters to produce more accurate results score.

Possible improvements for the search queries could be having more alternative metrics to produce average precision that would be used to gather the MAP for overall system performance.