
KBaseSearchEngine Documentation

Release 0.01

Gavin Price, Arfath Pasha

Dec 22, 2017

CONTENTS:

| | | |
|-----------|---|-----------|
| 1 | Introduction | 1 |
| 2 | About This Release | 3 |
| 3 | Known Issues | 5 |
| 4 | System Requirements | 7 |
| 5 | Configuration Details | 9 |
| 6 | Installation Instructions | 11 |
| 6.1 | Local Deployment | 11 |
| 6.2 | Production Deployment | 14 |
| 7 | API | 15 |
| 8 | Functional Requirements | 17 |
| 9 | NonFunctional Requirements | 19 |
| 10 | Architecture - High Level Design | 21 |
| 10.1 | Data Flow | 21 |
| 11 | Architecture - Low Level Design | 23 |
| 12 | Formal Specs | 25 |
| 12.1 | Type Mapping | 25 |
| 13 | Performance Evaluation | 27 |
| 14 | Operations | 29 |
| 15 | TODOs (ideas for improvement) | 31 |
| 16 | Release Notes | 33 |
| 17 | Indices and tables | 35 |

INTRODUCTION

This project consists of,

1. an ETL application that extracts information from Workspace objects, transforms and loads a subset of this information into an ElasticSearch Index.
2. a search API that makes queries into the ElasticSearch index.

ABOUT THIS RELEASE

KNOWN ISSUES

SYSTEM REQUIREMENTS

1. [KBase Workspace](#) (See section on [Deploying the Workspace Service locally](#))
2. [ElasticSearch v5.5.2](#). Note: the current implementation is not compatible with ElasticSearch v6+.
3. [Kibana v5.5.2](#)

CONFIGURATION DETAILS

INSTALLATION INSTRUCTIONS

6.1 Local Deployment

Follow these instructions for a local deployment once the *System Requirements* have been satisfied. These instructions are known to work on Ubuntu 16.04 LTS. The rest of this playbook assumes that you have all dependency binaries in your system environment path variable. At a high level, the steps are -

1. Start ElasticSearch
2. Start Kibana
3. Configure Workspace listeners to write events to Search mongo db
4. Restart Workspace service
5. Create a Workspace data type
6. Configure KBaseSearchEngine
7. Start worker
8. Start coordinator
9. Verify ElasticSearch index

1. Open a new terminal and start ElasticSearch.

Note: Elastic Search can only be started up by a non-root user

```
$ elasticsearch
```

2. Open a new terminal and start Kibana.

```
$ kibana
```

3. Configure the Workspace listeners to write events to the Search mongodb.

```
$ gedit [PATH_TO_YOUR_WORKSPACE_DIR]/deploy.cfg
```

Add the following lines under the listener configuration section -

```
listeners = Search
listener-Search-class = us.kbase.workspace.modules.SearchPrototypeEventHandlerFactory
listener-Search-config-mongohost = localhost
listener-Search-config-mongodatabase = Search_test
listener-Search-config-mongouser = ""
listener-Search-config-mongopwd = ""
```

4. Restart the Workspace Service. (See section on [Deploying the Workspace Service locally](#))

5. Open a new terminal and save the following document as `Empty.spec`. Then load into `ipython`, register the spec and save an object of this type to the Workspace. Saving a new object will cause the Workspace listener to write a new event to the mongo instance. Note that the `ws.administer()` command below requires administration privileges on the workspace.

```
module Empty {  
  
    /* @optional foo */  
    typedef structure {  
        int foo;  
    } AType;  
};
```

```
$ ipython  
  
In [1]: spec = open("[PATH_TO_SPEC]/Empty.spec").read()  
In [2]: ws.request_module_ownership('Empty')  
In [3]: ws.administer({'command': 'listModRequests'})  
Out[4]:  
[{'moduleName': u'Empty', ...}]  
In [5]: ws.administer({'command': 'approveModRequest', 'module': 'Empty'})  
In [6]: ws.register_typespec({'spec': spec, 'new_types': ['AType'], 'dryrun': 0})  
Out[7]: {u'Empty.AType-0.1': ....}  
In [8]: ws.release_module('Empty')  
Out[9]: [u'Empty.AType-1.0']  
In [10]: ws.save_objects({'id': 1, 'objects': [{'type': 'Empty.AType', 'data': {'bar  
↪': 'baz'}, 'name': 'myobj'}]})  
Out[11]:  
[[1,  
u'myobj',  
...  
]]
```

Create a new terminal and start mongo to check to make sure the event has been written. Note that the status is UNPROC (unprocessed event).

```
$ mongo  
> show dbs  
Search_test  
admin  
local  
workspace  
ws_types  
> use Search_test  
switched to db Search_test  
> db.getCollectionNames()  
["searchEvents"]  
> db.searchEvents.findOne()  
{  
  "_id": ...,  
  "strcde": "WS",  
  "accgrp": 1,  
  ...  
  "status": "UNPROC"  
}
```

6. Create a new terminal and edit `search_tools.cfg`, create a test data type and build the executable script.


```
$ cd [PATH_TO_YOUR_KBaseSearchEngine_DIR]
$ git checkout master
$ git pull
$ cp search_tools.cfg.example search_tools.cfg
$ gedit search_tools.cfg
```

Make the following edits. Note: the user for the token used below must have workspace admin privileges.

```
search-mongo-host=localhost
search-mongo-db=Search_test
elastic-host=localhost
elastic-port=9200
scratch=[PATH_TO_DIR_WHERE_TEMP_FILES_CAN_BE_STORED_BY_APP]
workspace-url=http://localhost:7058
auth-service-url=https://ci.kbase.us/services/auth/api/legacy/KBase/Sessions/Login
indexer-token=[YOUR_CI_TOKEN]
types-dir=[PATH_TO_YOUR_KBaseSearchEngine_DIR]/KBaseSearchEngine/test_types
type-mappings-dir=[PATH_TO_YOUR_KBaseSearchEngine_DIR]/KBaseSearchEngine/test_type_
↪mappings
workspace-mongo-host=fake
workspace-mongo-db=fake
```

```
$ mkdir test_types
$ cd test_types
$ gedit Empty.json
```

```
{
  "global-object-type": "EmptyAType2",
  "ui-type-name": "A Type",
  "storage-type": "WS",
  "storage-object-type": "Empty.AType",
  "indexing-rules": [
    {
      "path": "whee",
      "keyword-type": "string"
    },
    {
      "path": "whee2",
      "keyword-type": "string"
    }
  ]
}
```

```
$ cd ..
$ mkdir test_type_mappings
$ make build-executable-script JARS_DIR=[ABSOLUTE_PATH_TO_KBASE_JARS_DIR] KB_
↪RUNTIME=[PATH_TO_YOUR_ANT_INSTALL_DIR (example /usr/share)]
```

8. Start a worker

```
$ bin/search_tools.sh -c search_tools.cfg -k myworker
Press return to shut down process
```

9. Start the coordinator. Note that the event is processed and data has been indexed.

```
$ bin/search_tools.sh -c search_tools.cfg -s
Press return to shut down process
```

```
Moved event xxx NEW_VERSION WS:1/1/1 from UNPROC to READY
Event xxx NEW_VERSION WS:1/1/1 completed processing with state INDX on myworker
```

10. Open Kibana in browser with url `localhost:5601/app/kibana#/dev_tools/console?_g=()`

On Kibana console, make the following query

```
GET _search
{
  "query": {
    "match_all": {}
  }
}

GET _cat/indices

GET kbase.1.emptytype2/data/_search
```

The results for the query should appear on the right panel.

6.2 Production Deployment

FUNCTIONAL REQUIREMENTS

NONFUNCTIONAL REQUIREMENTS

ARCHITECTURE - HIGH LEVEL DESIGN

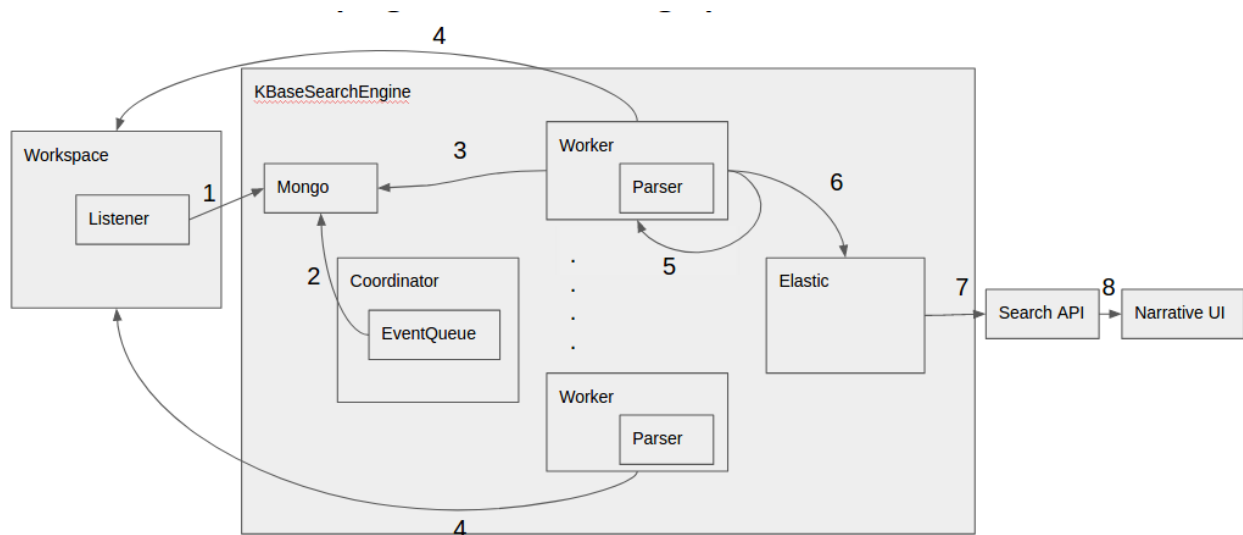


Fig. 10.1: KBaseSearchEngine component diagram.

10.1 Data Flow

1. The workspace pushes workspace level, object level and version level events into the KBaseSearchEngine MongoDB instance. The initial state of the events is UNPROC (or unprocessed).
2. The EventQueue periodically fetches events from the database and sets those that can be processed into a READY state. The EventQueue is a three-level blocking queue that blocks events that may cause an out-of-order update on the index. For example, an object level event like “rename” must block another object level event like “delete”. i.e. these two events cannot be executed in parallel by the workers. Also, importantly, the queue prevents simultaneous updates on the same document in Elasticsearch, which can cause update conflicts.
3. The workers pull events that are READY for processing, set their state to PROCESSING in the database instance and begin processing the event.
4. When necessary (like for new version events for example) the workers make requests to the workspace to get object data for processing.
5. If the processing of an event fails due to network connectivity or other such recoverable issues, the event is reprocessed using a Retrier. The Retrier retries an event a finite number of times before setting the event processing state as FAILED in the mongo instance. A log is written out when processing of an event fails.

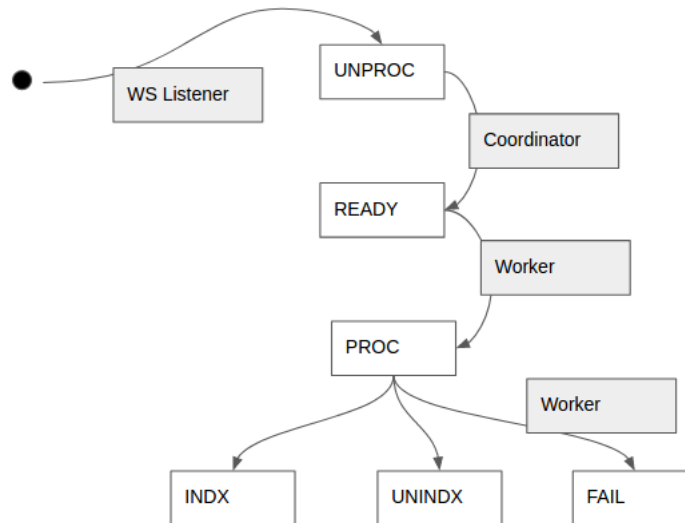


Fig. 10.2: Event state transition diagram.

6. Once the event has been processed successfully, the corresponding object is (re-)indexed into the ElasticSearch index and the event state in the mongo instance is set to INDX (or indexed).

7&8. Queries from the narrative UI are serviced by the search API which in turn makes queries to the ElasticSearch index.

ARCHITECTURE - LOW LEVEL DESIGN

FORMAL SPECS

12.1 Type Mapping

Type transformation specifications

Example type transformation *json* files are in *resources/types*.

TODO documentation of the transformation spec.

Type mapping specifications

Type mappings are optional *yaml* files that specify how to map data source types to search types. If provided for a source type, they override the mapping provided in the type transformation file(s) (in the *source-type* and *source-object-type* fields). In particular, the mapping files are aware of the source type version while the transformation files are not.

Mapping files also allow using the same set of transformation files for multiple environments where the source types may not have equivalent identifiers by providing environment-specific mapping files.

There is an example mapping file in *resources/typemappings* that explains the structure and how the mappings work.

PERFORMANCE EVALUATION

OPERATIONS

Backup and Restoration of index Deploying multiple indexes

TODOS (IDEAS FOR IMPROVEMENT)

1. After the first release, check if it makes sense to make `global-object-type` the same as the `storage-object-type`. This can help simplify the type spec by getting rid of one of the key-value pairs. It is also nice to keep types consistent across systems. The benefit that should be sought is if it will allow users to search specifically within the `kbaseGenomes` set, `someOtherGenomes` sets, and all genomes in general. Currently, it is not clear what the final indexed genome object looks like and what its searchability is like.
2. It seems like it will be better to completely remove “full-text” as a key in the type specs and instead use “keyword-type: text”. It simplifies the spec. If “keyword-type: text” is currently substitutable for “full-text: true” then I think we should make this change and remove the “full-text” as a valid key. It seems like the option “full-text: false” and “keyword-type: string” (for structured content) is supported in ESV5 with just “keyword-type: keyword”.
3. Make all type specs YAML and escape all json files on loading type specs.
4. Must write a formal spec for type spec

RELEASE NOTES

INDICES AND TABLES

- `genindex`
- `modindex`
- `search`