

TOWARDS A CONCURRENT IMPLEMENTATION OF KEYWORD SEARCH OVER RELATIONAL DATABASES

by

Richard J.I. Drake

A Thesis Submitted in Partial Fulfillment
of the Requirements for the Degree of

Master of Science (MSc)

in

Faculty of Science

Computer Science

University of Ontario Institute of Technology

Supervisor: Dr. Ken Q. Pu

December 2013

Copyright © Richard J.I. Drake 2013

Abstract

Traditional relational database systems offer powerful data modelling and querying capabilities. Unfortunately a relational database does not permit users to perform natural (keyword) queries. We present a system for automatic mapping of a relational schema into a document schema in order to facilitate fast and powerful full-text (keyword) search. By providing a facility for users to conduct keyword searches, we improve accessibility and decrease search time complexity. In addition, we explore further improvements to time complexity by utilizing concurrent implementations of our search algorithms.

Keywords relational database; full-text search

Preface

Background and Motivation

The introduction of keyword search has revolutionized how we find information. Over the years, numerous techniques have been developed which make searching through large amounts of information for one or more keywords extremely fast. With the advent of faster computer hardware, we are able to not only search through small attributes of a document (eg. Title, synopsis, etc.) but rather the entirety of the document itself.

An example of an early system which utilized keyword search would be a library catalogue. Such a system would allow a user to search, by keyword, for the title of a book, manuscript, etc. The results would show item titles matching the keyword(s), as well as other information such as whether or not the item is in circulation, as well as where it is located within the library. This information would come from a relational database.

With the rise of the World Wide Web, much information was placed online. This information would be easy to access if one knew how to locate it. Unfortunately, over time, so much information existed on the World Wide Web that it became difficult to keep track of it all. There was a need to index all of this information and make it accessible. This need was filled by a Web search engine.

Initial Web search engines comprised of simple scripts that gathered listings of files on FTP servers; they were essentially link farms. A few short years later, the first full-text (keyword) search engine, WebCrawler, was released.

While full-text search engines provided an excellent means for locating information, as the

Web grew larger, the volume of noise also grew larger. In addition, every Web page could be structured in a different way; the Web was largely a collection of unstructured documents. That is, there were few obvious links between them.

Search engines such as Google attempted to solve this problem by introducing new algorithms, such as PageRank, to rank Web pages on both the relevance of their content as well as their reputation. The idea was if a page is linked to often, it is considered to be more authoritative on a subject than a page with fewer links. This allowed the relevance of a page to be computed based on not only its contents, but its artificial importance.

Molly attempts to avoid some of the issues plaguing search engines. It deals primarily with structured, filtered data. This allows us to provide results with less noise. In addition, the fact that it deals with structured data means links between documents are explicitly stated. Rather than inferring a link between documents based on hyperlinks, we know when two documents are linked together.

This thesis provides an overview of the Molly system.

Outline

Changed.

Contents

1	Background (2 days)	1
2	A Tale of Two Data Models	2
2.1	Relational model with star schema	2
2.1.1	Relational Model of Data	2
2.1.2	Star Join Schema to Form Entity Groups	3
2.1.3	Instances of an Entity Group	4
2.2	Pros and Cons of the Relational Model	5
2.2.1	Pros	5
2.2.2	Cons	6
2.3	Document model (4 days, week 2)	6
2.4	Pro and con of document model (1 day)	7
2.5	Best of both worlds (4 days, week 3)	7
3	Along came Clojure	8
3.1	Basic principles of functional programming (2 days)	8
3.2	Features of Clojure (2 days)	8
4	Search w/ Clojure	9
4.1	Thirdparty libraries (1 day, week 4)	9
4.2	Indexing of relational objects (5 days, week 5)	9
4.3	Keyword Search in document space (5 days, week 6)	9

4.4	Graph Search in document space (5 days, week 7)	10
5	Experimental evaluation (5 days, week 8)	11
5.1	Implementation	11
5.2	The data set	11
5.3	Runtime Evaluation	11
5.4	Lessons learned	12
6	Conclusion (0 days)	13
7	Definitions	14
7.1	Data Representation & Notation	14
8	RDBMS to Document Store	16
9	Concurrency	20
10	Theoretical	24
10.1	Data Representation & Notation	25
10.2	Relational Database Abstraction	26
10.3	Data Corpus	26
10.4	Graph Search Algorithms	28
11	Implementation	30
11.1	Functional vs. Procedural Programming of Algorithms	30
11.1.1	Functional Data Structures	30
11.2	Tuneable Parameters	30
12	System Implementation	31
12.1	Technology Selection	31
12.1.1	Programming Language	31

12.1.2	Relational Database	32
12.1.3	Full-Text Search Database	33
12.1.4	Web Stack	33
12.2	System Design	33
12.2.1	Configuration	34
12.3	Implementation Issues	35
13	Performance & Evaluation	36
13.1	Environment	36
13.2	Methodology	36
13.2.1	Pitfalls of the Java Virtual Machine	37
13.2.2	Mitigating JVM Pitfalls	38
14	Conclusion	48
A	Source Code	49
A.1	molly	49
A.1.1	molly.core	49
A.2	molly.conf	52
A.2.1	molly.conf.config	52
A.2.2	molly.conf.mycampus	53
A.3	molly.datatypes	57
A.3.1	molly.datatypes.database	57
A.3.2	molly.datatypes.entity	58
A.3.3	molly.datatypes.schema	61
A.4	molly.index	62
A.4.1	molly.index.build	62
A.5	molly.util	63
A.5.1	molly.util.nlp	63

A.6	molly.search	64
A.6.1	molly.search.lucene	64
A.6.2	molly.search.query_builder	66
A.7	molly.server	67
A.7.1	molly.server.serve	67
A.8	molly.algo	70
A.8.1	molly.algo.common	70
A.8.2	molly.algo.bfs	72
A.8.3	molly.algo.bfs_atom	73
A.8.4	molly.algo.bfs_ref	74
A.9	molly.bench	75
A.9.1	molly.bench.benchmark	75

List of Tables

10.1	Courses entity schema	26
10.2	Instructors entity schema	27
10.3	Sections entity schema	27
10.4	Schedules entity schema	28
10.5	Teaches entity schema	28
13.1	Reported resolutions of <code>currentTimeMillis()</code> by platform [1].	37

List of Figures

7.1	The structure of an entity	14
7.2	The structure of an entity group	15
8.1	Determining whether or not a field is classified as “special.”	16
8.2	Unique identifier generation.	16
8.3	Creation of a field.	17
8.4	Creation of a document.	17
8.5	Transformation of a row into the internal representation.	18
8.6	Transformation of a document into the internal representation.	19
8.7	Transformation of the internal representation into a document.	19
9.1	Retrieves the entity with the given ID.	20
9.2	Finds all groups containing an entity with the given ID.	20
9.3	Finds adjacent nodes in the graph (entities linked in groups).	21
9.4	Discovers new nodes along the frontier and updates the state accordingly.	21
9.5	BFS using recursion.	21
9.6	Discovers new nodes along the frontier and updates the state accordingly.	22
9.7	Implements a concurrent version of BFS using atoms.	22
9.8	Discovers new nodes along the frontier and updates the state accordingly.	23
9.9	Implements a concurrent version of BFS using references.	23
10.1	The structure of an entity	25

10.2	The structure of an entity group	26
13.1	core.clj	38
13.2	1 Hop	39
13.3	2 Hops	40
13.4	3 Hops	41
13.5	4 Hops	42
13.6	5 Hops	43
13.7	6 Hops	44
13.8	7 Hops	45
13.9	8 Hops	46
13.10	Comparison between single threaded and concurrent graph search	47

List of Algorithms

Chapter 1

Background (2 days)

Literature search on:

- DBExplore
- XRank
- BANKS
- ...

Chapter 2

A Tale of Two Data Models

2.1 Relational model with star schema

2.1.1 Relational Model of Data

$\langle DB \rangle ::= \langle name \rangle$

Definition 1. Database

Let d be a database instance. A database is comprised of three main components:

- $NAME[d] : \text{string}$
- $REL[d] : \text{list}(\text{REL})$
- $FK[d] : \text{list}(\text{FK})$

Definition 2. Relation

Let $r \in REL[d]$, where d is defined in Definition 1. A relation is comprised of three main components:

- $NAME[r] : \text{string}$
- $ATTR[r] : \text{list}(\text{ATTR})$
- $KEY[r] : \text{list}(\text{ATTR})$

The first is a **string** representation of the relation. The second is a list of attributes that make up entries, or tuples, in the relation. The third is a list of the relation's attributes that uniquely identify the tuple within the relation. That is, $\text{KEY}[r] \subseteq \text{ATTR}[r]$.

Definition 3. Attribute

Let $a \in \text{ATTR}[r]$, where r is defined in Definition 2. An attribute is comprised of two main components:

- $\text{NAME}[a]$
- $\text{TYPE}[a]$

Note: Lower case letters (e.g. a, b, c, \dots) are attributes.

Definition 4. Foreign Key

Let $\theta \in \text{FK}[d]$ be a FK instance, where d is defined in Definition 1. A foreign key is comprised of two main components:

- $\text{FROM}[fk] = (\text{REL}_s[\theta], \text{ATTR}_s[\theta])$
- $\text{TO}[fk] = (\text{REL}_t[\theta], \text{ATTR}_t[\theta])$

Note: θ, ϕ are the FK constraints

2.1.2 Star Join Schema to Form Entity Groups

A network (forest) of tuples, jointed via some existing $\theta \in \text{FK}[d]$.

The schema of entity group G is defined as vertices of G :

$$V(G) \subseteq \text{REL}(DB)$$

Relation r in the space of vertices of G , $V(G)$, may be a table or a computed view.

$$r \in V(G)$$

It is also defined as the edges of G , or $E(G)$, in the form of

$$r(a_1, a_2, \dots, a_k) \rightarrow s(b_1, b_2, \dots, b_k)$$

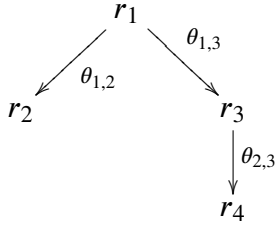
where $a_i \in \text{ATTR}[r]$, $b_i \in \text{ATTR}[s]$, with the additional constraint of $r, s \in V(G)$.

Example 1.

$$\text{Instructor}(\text{name}) \rightarrow \text{Schedule}(\text{instructor})$$

$$\text{Schedule}(\text{code}) \rightarrow \text{Course}(\text{id})$$

2.1.3 Instances of an Entity Group



Instances are obtained by the following process.

For $r_i(a_{i,1}, a_{i,2}, \dots, a_{i,k}) \rightarrow r_j(b_{j,1}, b_{j,2}, \dots, b_{j,k})$

$$c_{ii,j} = \bigwedge_{n=1}^k (a_{i,n} = b_{j,n})$$

$$\begin{aligned} \text{VIEW}[G] &= \bowtie_{\theta_{ij}} (r_i, r_j) \\ &= r_1 \bowtie_{\theta_{1,2}} r_2 \bowtie_{\theta_{2,3}} r_3 \dots \bowtie_{\theta_{n,n+1}} r_n \end{aligned}$$

where $r_1, r_2, r_3, \dots, r_n$ are relations discovered by a depth-first search traversal of G .

Each tuple in $\text{VIEW}[G]$ is an instance of entity group G .

Motivation:

To do (1)

Database schema

Vertices: $REL[d]$

Edges: $FK[d]$

The entity graphs are overlapping subgraphs at the schema level.

Question:

How to determine connectivity at the instance level?

- ER style relational schema
- Star join schema to form entity groups
- Expressing relational objects using universal design pattern (describing data using scalar, lists and dictionaries).
- Relational object graph

2.2 Pros and Cons of the Relational Model

In order to better understand the motivation behind this work, it is important to examine the strong as well as weak points of the relational model.

2.2.1 Pros

- Well supported by relational algebra and relational databases (RDBMS)
- Clean and consistent database instances (ACID?)
- Can use queries to resolve instance-level connectivity
 - How is "Ken" connected to "CSCI 3030U"?

To do (2)

2.2.2 Cons

- Must know the relational schema
 1. Know table/attribute names
 2. Know join paths (schema)
- Inflexible string matching options (basically just have `LIKE`), substring matching
- Must know SQL
- All queries must be re-written upon schema changes (rename, change in join path, etc.)
- Not adaptive to new join path (e.g. newly created entity group, deleted E.G. etc.)
- Good for analytics (aggregation, selection) if user has domain knowledge of the schema.
- Bad for exploratory queries.
- Bad if user doesn't know SQL
- Bad for flexibility

2.3 Document model (4 days, week 2)

- Definition: documents, terms, and the bag of terms model for documents and queries
- Definition of keyword search queries: vectorization of documents (tf-idf) and queries. Models of distance between documents and queries (cosine-distance, jaccard distance, BM25).
- Extended document model with attributes and fields
- Expressing documents in the universal design pattern (aka list+dict)
- Document graph

2.4 Pro and con of document model (1 day)

- Good: exploratory queries using keywords (google)
- Good: easy (or no) syntax
- Good: fuzzy matching (using n-gram)
- Bad: No analytics

2.5 Best of both worlds (4 days, week 3)

- Hybrid database defined by both the relational model and the document model
- Translation between relational objects (entities and entity) groups to documents.
- Translation of documents back to relational objects.
- Proof of lossless translation between relational space and document space

Chapter 3

Along came Clojure

3.1 Basic principles of functional programming (2 days)

- immutable data structures
- persistent data structures using multi-versioning
- functions (and higher order functions) as values

3.2 Features of Clojure (2 days)

- Data structures supporting the universal design pattern
- Concurrency + STM
- Interoperability with JVM (including Lucene)

Chapter 4

Search w/ Clojure

4.1 Thirdparty libraries (1 day, week 4)

- Lucene

4.2 Indexing of relational objects (5 days, week 5)

- Schema definition
- Crawling using SQL
- Indexing using relational objects
- Fuzzy indexing of values (typed by classes)

4.3 Keyword Search in document space (5 days, week 6)

- Disambiguate keywords using fuzzy search (suggestion, overloaded terms)
- Flexibility keyword search for documents
- Translate search result back to relational space

4.4 Graph Search in document space (5 days, week 7)

- Why we need graph search
- Search in document graph using graph search algorithms with functional implementations:
(Ford Fulkerson, BFS)
- Speed up using concurrency
- Clojure specific optimization: ref + atom

Chapter 5

Experimental evaluation (5 days, week 8)

5.1 Implementation

- Choice of language
- Statistics about the code base: LOC, classes, ?
- Github hosted

5.2 The data set

- Description of the data set
- Statistics of the data set

5.3 Runtime Evaluation

- Index speed
- Keyword search speed
- Graph search speed:
 - Ford Fulkerson

- BFS
- Concurrent BFS using refs
- Concurrent BFS using atoms

5.4 Lessons learned

- Simple algorithms are easier to parallelize
- STM is effective: transactions do not rollback (that much), so we observe impressive speed-up in concurrent versions.
- Fine tuning is beneficial: atom is better than ref.
- The clojure way: correctness first, runtime optimization latter (ref to atom is natural).

Chapter 6

Conclusion (0 days)

Survived Clojure.

Chapter 7

Definitions

7.1 Data Representation & Notation

The data from the database is represented in various data structures. There are separate representations for each type of data: values, entities, and entity groups.

Value

Definition 5. A **Value** represents a single piece of information. To avoid repetition, each value is unique. That is, $\exists! v \in V$, where v is a value in the set V of all values.

Entity

Definition 6. An **Entity** is a collection of attributes, a_n , each mapped to a single value, v_n . An entity also includes additional information such as a unique identifier.

id	$T_n v_{id}$
a_1	v_1
a_2	v_2
\vdots	\vdots
a_n	v_n

Figure 7.1: The structure of an entity

Entities are analogous to rows in a database table. Thus, the unique identifier is generated based on the table name, T_n , as well as unique key in the table, v_{id} . The unique key identifies the row, and the table name identifies the table. Together they uniquely identify the entity within the entire database. Attributes are analogous to columns in a database table.

$\exists! e_{id} \in E$, where E is the set of all entities.

Entity Group

Definition 7. An **Entity Group** joins together two or more entities. These entity groups can also have attributes, a_n , and values, v_n , associated with them much like entities.

$$\begin{array}{ll} e_L & [e_1, e_2, \dots, e_n] \\ a_1 & v_1 \\ a_2 & v_2 \\ \vdots & \vdots \\ a_n & v_n \end{array}$$

Figure 7.2: The structure of an entity group

Chapter 8

RDBMS to Document Store

```
9 (defn special?
10   [field-name]
11   (and (.startsWith field-name "__") (.endsWith field-name "__"))))
```

Figure 8.1: Determining whether or not a field is classified as “special.”

```
13 (defn uid
14   "Possible inputs include:
15   row :T :ID
16   row [[:T :ID] [:T :ID]]
17   row [[:T :ID :desc] [:T :ID :desc]]"
18   ([row C id]
19    (if (nil? (row id))
20        (throw
21          (Exception.
22            (str "ID column " id " does not exist in row " row ".")))
23        (str (name C)
24              "|"
25              (clojure.string/replace (row id) #"\s+" "_"))))
26   ([row Tids]
27    (clojure.string/join " " (for [[C id] Tids]
28                                (uid row C id)))))
```

Figure 8.2: Unique identifier generation.

```
30 (defn field
31   [field-name field-value]
32   (Field. field-name
33           field-value
34           Field$Store/YES
35           Field$Index/ANALYZED))
```

Figure 8.3: Creation of a field.

```
37 (defn document
38   [fields]
39   (let [doc (Document.)]
40     (do
41       (doseq [[field-name field-value] fields]
42         (.add doc (field (name field-name) (str field-value))))
43       doc)))
```

Figure 8.4: Creation of a document.

```

45 (defn row->data
46   ^{:doc "Transforms a row into the internal representation."}
47   [this schema]
48   (let [T      (schema :T)
49         C      (schema :C)
50         attr-cols (schema :attrs)
51         attrs    (if (nil? attr-cols)
52                     this
53                     (select-keys this attr-cols))
54         meta-data {:type T :class C}
55         id-col    (schema :ID)]
56     (with-meta (if (= T :group)
57                 (conj attrs {:entities (uid this id-col)})
58                 attrs)
59              (condp = T
60                  :value    (assoc meta-data
61                                   :class
62                                   (clojure.string/join "|"
63                                   (map name
64                                        [C (first attr-cols)])))
65                  :entity  (assoc meta-data :id
66                                   (if (coll? id-col)
67                                       (uid this id-col)
68                                       (uid this C id-col)))
69                  :group   (assoc meta-data
70                                   :entities
71                                   (uid this id-col))
72                  (throw
73                   (IllegalArgumentException.
74                    "I only know how to deal with types :value,
75                    :entity, and :group"))))))))

```

Figure 8.5: Transformation of a row into the internal representation.

```

77 (defn doc->data
78   ^{:doc "Transforms a Document into the internal representation."}
79   [this]
80   (let [fields      (.getFields this)
81         extract     (fn [x] [(keyword (clojure.string/replace
82                                     (.name x) "_" ""))
83                               (.stringValue x)])
84         check-special (fn [x] (special? (.name x)))
85         filter-fn     (fn [f] (apply hash-map
86                                   (flatten
87                                    (map extract
84                                     (filter f fields))))))]
88     (with-meta (filter-fn (fn [x] (not (check-special x))))
89               (filter-fn check-special))))
90

```

Figure 8.6: Transformation of a document into the internal representation.

```

92 (defn data->doc
93   ^{:doc "Transforms the internal representation into a Document."}
94   [this]
95   (let [int-meta  (meta this)
96         T         (int-meta :type)
97         all       (clojure.string/lower-case
98                   (clojure.string/join " "
99                                         (if (= T :entity)
100                                             (conj (vals this)
101                                                    (name
102                                                     (int-meta :class)))
103                                             (vals this))))
104         luc-meta  [[:__type__ (name T)]
105                   [[:__class__ (name (int-meta :class))]
106                    [[:__all__ (if (= T :value)
107                                   (q-gram all)
108                                   all))]]
109         raw-doc   (concat luc-meta
110                           this
111                           (condp = (int-meta :type)
112                               :value  [[:value all]]
113                               :entity [[:__id__ (int-meta :id)]]
114                               :group  [[]]))

```

Figure 8.7: Transformation of the internal representation into a document.

Chapter 9

Concurrency

```
6 (defn find-entity-by-id
7   [G id]
8   (let [query (boolean-query [[(query :type :entity) :and]
9                               [(query :id id) :and]])]
10     (map doc->data (idx-search G query 10))))
```

Figure 9.1: Retrieves the entity with the given ID.

```
12 (defn find-group-for-id
13   [G id]
14   (let [query (boolean-query [[(query :type :group) :and]
15                               [(query :entities id) :and]])
16         results (map doc->data (idx-search G query 10))
17         big-str (clojure.string/join " "
18                                       (map #(% :entities) results))]
19     (distinct (clojure.string/split big-str #"\\s{1}"))))
```

Figure 9.2: Finds all groups containing an entity with the given ID.

```

21 (defn find-adj
22   [G v]
23   (remove #{v} (find-group-for-id G v)))

```

Figure 9.3: Finds adjacent nodes in the graph (entities linked in groups).

```

4  (defn update-adj
5    [G marked dist prev u max-hops]
6    (loop [adj      (find-adj G u)
7            marked   marked
8            dist     dist
9            prev     prev
10           frontier []]
11      (if (or (empty? adj) (>= (dist u) max-hops))
12          [(conj marked u) dist prev frontier]
13          (let [v      (first adj)
14                adj'   (rest adj)]
15              (if (marked v)
16                  (recur adj' marked dist prev frontier)
17                  (let [dist'   (assoc dist v (inc (dist u)))
18                        prev'   (assoc prev v u)
19                        frontier' (conj frontier v)]

```

Figure 9.4: Discovers new nodes along the frontier and updates the state accordingly.

```

21
22 (defn bfs
23   [G s t max-hops]
24   (loop [Q      (-> (clojure.lang.PersistentQueue/EMPTY) (conj s))
25          marked #{s}
26          dist   {s 0}
27          prev   {s nil}]
28       (if (or (empty? Q)
29               (some (fn [node] (= node t)) marked))
30           [marked dist prev]
31           (let [u      (first Q)
32                 Q'     (rest Q)
33                 [marked' dist' prev' frontier]
34                 (update-adj G marked dist prev u max-hops)]

```

Figure 9.5: BFS using recursion.

```

4 (defn update-adj
5   [state-ref G u max-hops]
6   (let [marked? (@state-ref :marked)
7         deferred (if (>= ((@state-ref :dist) u) max-hops)
8                        []
9                        (doall
10                         (for [v (find-adj G u)]
11                           (if (marked? v)
12                               nil
13                               (future

```

Figure 9.6: Discovers new nodes along the frontier and updates the state accordingly.

```

15                                     state-ref
16                                     update-state
17                                     u
18                                     v
19                                     max-hops))))))]]
20   (doall (map deref-future deferred)))
21
22 (defn bfs-atom
23   [G s t max-hops]
24   (let [state-ref (atom (initial-state s))]
25     (while (and (not (empty? (@state-ref :Q)))
26                 (not (@state-ref :done)))

```

Figure 9.7: Implements a concurrent version of BFS using atoms.

```

4 (defn update-adj
5   [state-ref G u max-hops]
6   (let [marked? (@state-ref :marked)
7         deferred (if (>= ((@state-ref :dist) u) max-hops)
8                        []
9                        (doall
10                         (for [v (find-adj G u)]
11                           (if (marked? v)
12                               nil
13                               (future (dosync (alter
14                                           state-ref
15                                           update-state
16                                           u

```

Figure 9.8: Discovers new nodes along the frontier and updates the state accordingly.

```

18                                     max-hops)))))))]
19   (doall (map deref-future deferred))))
20
21 (defn bfs-ref
22   [G s t max-hops]
23   (let [state-ref (ref (initial-state s))]
24     (while (and (not (empty? (@state-ref :Q)))
25                (not (@state-ref :done)))
26       (let [u (first (@state-ref :Q))
27             Q' (pop (@state-ref :Q))]
28         (dosync (alter state-ref assoc :Q Q'))
29         (if (some (fn [node] (= node t)) (@state-ref :marked))

```

Figure 9.9: Implements a concurrent version of BFS using references.

Chapter 10

Theoretical

The problem of efficiently searching through a graph is the subject of countless articles and journal papers. Numerous algorithms have been proposed to perform graph search in an efficient manner. Many of these algorithms build upon their predecessors, making assumptions and changing aspects to better suit a particular problem area.

The applications of efficient graph search algorithms are endless. Algorithms such as Dijkstra's are used heavily in path finding, for example in a portable GPS unit. A* Search is used in the area of computer vision to approximate the best path to take. Companies such as Amazon make use of graph search algorithms in order to find related products for consumers to purchase.

This thesis concentrates on building a system which can be utilized in order to discover related information. The system that was built is capable of finding not only related information from neighbouring nodes, but also the best path between two arbitrary nodes.

Section 10.1 provides an overview of how the data is represented in the system. It also defines notation to represent this data. Section 10.2 outlines how the relational database is related to the data in the system. Section 10.3 provides an example data corpus which is utilized throughout this thesis. It details the "mycampus" dataset. Finally, Section 10.4 provides justification for the algorithm chosen to perform the graph search. It discusses multiple different graph search algorithms, as well as why the one used in this thesis was chosen.

10.1 Data Representation & Notation

The data from the database is represented in various data structures. There are separate representations for each type of data: values, entities, and entity groups.

Value

Definition 8. A **Value** represents a single piece of information. To avoid repetition, each value is unique. That is, $\exists! v \in V$, where v is a value in the set V of all values.

Entity

Definition 9. An **Entity** is a collection of attributes, a_n , each mapped to a single value, v_n . An entity also includes additional information such as a unique identifier.

$$\begin{array}{ll} \text{id} & T_n | v_{id} \\ a_1 & v_1 \\ a_2 & v_2 \\ \vdots & \vdots \\ a_n & v_n \end{array}$$

Figure 10.1: The structure of an entity

Entities are analogous to rows in a database table. Thus, the unique identifier is generated based on the table name, T_n , as well as unique key in the table, v_{id} . The unique key identifies the row, and the table name identifies the table. Together they uniquely identify the entity within the entire database. Attributes are analogous to columns in a database table.

$\exists! e_{id} \in E$, where E is the set of all entities.

Entity Group

Definition 10. An **Entity Group** joins together two or more entities. These entity groups can also have attributes, a_n , and values, v_n , associated with them much like entities.

$$\begin{array}{ll}
e_L & [e_1, e_2, \dots, e_n] \\
a_1 & v_1 \\
a_2 & v_2 \\
\vdots & \vdots \\
a_n & v_n
\end{array}$$

Figure 10.2: The structure of an entity group

10.2 Relational Database Abstraction

10.3 Data Corpus

For illustrative purposes, the mycampus dataset will be used. This data comes from UOIT's course registration system.

There are several different entities which comprise the mycampus dataset:

- Courses
- Instructors
- Schedules
- Sections
- Teaches

The **Courses** entity represents a course. A course has an individual code, along with a title and a description (See Table 10.1). The code uniquely identifies the course.

Column	Type	Description
code	VARCHAR	Unique course code
title	VARCHAR	Title of the course
description	TEXT	A brief description

Table 10.1: Courses entity schema

The **Instructors** entity represents an individual instructor. Each instructor has a unique identifier, as well as a name (See Table 10.2). An instructor can be a Professor, Lecturer, Sessional Instructor, or Teaching Assistant.

Column	Type	Description
id	INT	Unique identifier
name	VARCHAR	The instructor's name

Table 10.2: Instructors entity schema

The **Sections** entity represents a section of a course. Courses may have many sections. For example, one section could be a lecture, while another is a lab. Some courses may have over a dozen sections, depending on the associated term.

Each section contains a unique identifier, capacity information (students enrolled, spots open, etc.), when registration is open for the section, how many credits it is worth, what level (eg. undergraduate or graduate), and what year the section is offered in (See Table 10.3).

Column	Type	Description
id	INT	Unique identifier
actual	INT	Number of people enrolled in the course
campus	VARCHAR	String uniquely identifying the campus
capacity	INT	Maximum number of people that may be enrolled in the section
credits	FLOAT	Number of credits awarded upon successful completion
levels	VARCHAR	The level of the course (eg. undergraduate, graduate, etc.)
registration_start	DATE	Date registration for the section opens
registration_end	DATE	Date registration for the section ends
semester	VARCHAR	String that uniquely identifies the semester
sec_code	INT	Unique section code (called a CRN)
sec_number	INT	Sequential number identifying the number of the section
year	INT	The year the section is offered in

Table 10.3: Sections entity schema

The **Schedules** entity represents a scheduled meeting of a section. A section may have many schedules. For example, a lecture section may meet twice a week.

Each schedule has a unique identifier, a date range in which the schedule is active, the time of the schedule, the type, location, and the day which the class takes place (See Table 10.4).

Column	Type	Description
id	INT	Unique identifier
date_start	DATE	First day of class
date_end	DATE	Last day of class
day	VARCHAR	Single character representing the day of the week
schedtype	VARCHAR	Lecture, tutorial, lab, etc.
hour_start	INT	Hour the class starts at
hour_end	INT	Hour the class ends at
min_start	INT	Minute the class starts at
min_end	INT	Minute the class ends at
classtype	VARCHAR	The type of class
location	VARCHAR	Unique location name where class is held

Table 10.4: Schedules entity schema

A **Teaches** entity is used to link together an instructor and a schedule. Each link has a unique identifier along with the instructor's position (eg. Teaching Assistant, Lecturer, etc.) (See Table 10.5).

Column	Type	Description
id	INT	Unique identifier
position	VARCHAR	Position of the Instructor with regard to a Schedule

Table 10.5: Teaches entity schema

10.4 Graph Search Algorithms

Careful consideration was given to which graph search algorithm was to be used. Among the choices were:

- Breadth-First
- Bellman-Ford
- Dijkstra
- A*

The first choice, Breadth-First Search, is among the simplest of the algorithms. It has the advantage of being simple to implement. It can only handle fixed costs for travelling between nodes, which can be a disadvantage.

Bellman-Ford is similar to BFS. It has the ability to deal with variable cost. This comes at the cost of increased difficulty in implementation. When the cost between nodes is fixed, Bellman-Ford essentially becomes BFS.

A greedy version of BFS is Dijkstra's Algorithm. It utilizes a priority queue rather than a regular queue, allowing it to be faster. As it is a greedy algorithm, Dijkstra's Algorithm may not return the optimal result.

An extension to Dijkstra's is A* search. Graph search spaces can be rather large. A* attempts to prune the search space based on a heuristic. A* is a natural choice to perform graph search in certain areas such as computer vision where an obvious heuristic exists. It has a disadvantage of consuming large amounts of memory (though IDA* attempts to limit memory consumption).

There are many more graph search algorithms. The above were primarily considered as many of the other algorithms are simple extensions with different data structures.

For this application, there is no obvious heuristic. This eliminates A*. There is also no obvious cost function. This eliminates Dijkstra's Algorithm. Bellman-Ford is very similar to BFS with the added ability to deal with variable cost. As the cost function is not obvious and thus constant, Bellman-Ford essentially reverts back to BFS.

For these reasons, BFS was chosen as the graph search algorithm. While the other candidates and others provide numerous advantages over BFS in many situations, this is not one of them.

Chapter 11

Implementation

11.1 Functional vs. Procedural Programming of Algorithms

Functional and procedural programming languages are entirely different paradigms. In functional programming, data is immutable. Functions must be pure and predictable (free of “side effects”).

Loops and control flow are accomplished with recursion and pattern matching.

Procedural languages allow for more flexibility at a cost of unpredictability.

Key differences between FP and procedural (for algorithms)

11.1.1 Functional Data Structures

Closure immutable data structures, versioning (persistent data structures), STM vs. locking

11.2 Tuneable Parameters

Tuneable parameters

Chapter 12

System Implementation

12.1 Technology Selection

12.1.1 Programming Language

Numerous programming languages were considered for the implementation. Among them were: Ruby, Python, and Clojure. Each of the languages presented both pros and cons. Ruby and Python are rather similar object oriented languages. Clojure is a functional language.

Ruby features a clean syntax and is very object oriented. Its object model was inspired by that of Smalltalk. It has a very active web development community and numerous web frameworks (eg. Rails, Sinatra, Padrino, etc.). In contrast to Python, it features a lean core, instead choosing to depend on third party libraries.

Python also features clean syntax and is object oriented. It embraces a “batteries included” philosophy to its standard library, featuring libraries for everything from serial communications to importing/exporting CSV files. It too has a strong web development community and numerous web frameworks (eg. Django, Flask, web2py, etc.).

Python is used extensively by scientists. As such, it has numerous plotting, computation, and simulation libraries. Examples include Matplotlib, SciPy, NumPy, and NetworkX.

Both of the above languages embrace functional programming elements. Python has filter, reduce, and map. Ruby’s collections are also capable of performing filter (select), reduce (inject),

and map. They both lack an important aspect of functional programming: immutable objects.

Clojure is a purely functional language built on top of the JVM. Like Python and Ruby, it is dynamically typed. Unlike Python and Ruby, it is a Lisp dialect and as such features macros and code-as-data.

As it runs on the JVM, Clojure allows us to utilize existing Java languages. Using Java libraries through Python or Ruby requires using their respective native code interfaces to Java's JNI. This process is typically automated through a tool such as SWIG.

Clojure was chosen as the implementation language for a number of reasons. First of all, as it is built on top of the JVM, it can utilize JDBC for database access. It can also make use of several other useful libraries. Secondly, its dynamic Lisp nature is a natural way of processing data. Clojure allows us to model our data structures directly after the data. In a language such as Python, one would need to make use of classes.

12.1.2 Relational Database

Several relational database management engines (RDBMS) were evaluated for this project. Key requirements were SQL support, and JDBC driver availability.

SQL, or Structured Query Language, is a powerful way to make complex queries. It is loosely based on Relational Algebra (RA)[3, p. 243]. In addition to being a data-manipulation language (similar to Relational Algebra), it is also a data-definition language.

JDBC, or Java Database Connectivity, was required in order for our application to access the database. JDBC is an API which allows us to access a database using a standard interface. It is the responsibility of the JDBC Driver to translate the API calls into the correct syntax for the database engine.

Any RDBMS that supports SQL-92 or higher and has a JDBC driver for it should be usable in this system. For testing purposes, we chose to use SQLite. SQLite provided a SQL interface, had a stable JDBC driver, and is ubiquitous. SQLite is an embedded database, allowing for local development without running a separate daemon process.

12.1.3 Full-Text Search Database

Two full-featured full-text search databases (FTSDB) were evaluated. The first, Xapian, is written as a C++ library. It features a rich query language, high performance, and a scalable design[5]. Bindings for Java are made available via a simple JNI library.

The second, Lucene, is a pure-Java library. It too features a rich query language and reasonable performance. Lucene also forms the basis of the Apache Solr search platform. As it is written in Java, Lucene runs natively on the JVM.

Both choices provided plenty of features, Java interfaces, and excellent performance. Ultimately Lucene was chosen. As it ran on the JVM, interfacing with the Lucene library from our Clojure code was simple. The fact that it was written in pure-Java and required no native libraries meant the entire project could be bundled as a Java Archive (JAR).

12.1.4 Web Stack

Clojure is a young language. As such, few web frameworks exist specifically for it (it can, however, interface with any Java web framework).

Clojure web frameworks are built on top of Ring, itself an HTTP abstraction layer with adapters for Java servlets as well as the Jetty web server. Ring is similar to WSGI in Python, or Rack in Ruby.

The web framework chosen for this project was Noir. Noir provided a simple Clojure API for building the JSON API to allow for client-side querying of the project's API.

12.2 System Design

The system is split into numerous components.

A configuration file is used to describe the entities and entity groups. It defines the relationship between entities, as well as their attributes and values. This configuration file is used by the crawler and indexer in order to retrieve data from the relational database and index it in a full-text search database.

Once the index is built, another component provides an internal API for querying the full-text search database. It is a thin Clojure wrapper over the Lucene API which provides the basis for search and retrieval of results. This wrapper is utilized by the system to discover relationships among entities.

Finally, the system can be interfaced with via a JSON API. The JSON API provides a simple HTTP wrapper around the search and discovery features of the system. The use of JSON allows client-side JavaScript to send and receive queries with XMLHttpRequest.

To do (3)

12.2.1 Configuration

Configuration files are simply Clojure code that follow a well-defined protocol. The protocol defines three main components of the system:

connection The relational database connection. It must be an instance of the **Database** protocol.

index The path to the full-text search database index file.

schema A definition of the database schema. The definition includes both the definition of entities, as well as the definition of entity groups. The schema is a list of instances of the **Schema** protocol.

Database Protocol

The **Database** protocol is defined as follows.

```
(defprotocol Database
  (execute-query [this query f]))
```

Each database instance must provide a function which accepts a query and function as arguments. It must execute the query provided and call function **f** on each row returned.

Schema Protocol

The Schema protocol is defined as follows.

```
(defprotocol Schema
  (crawl [this db-conn idx-w])
  (klass [this])
  (schema-map [this]))
```

A schema consists of the following information:

- A symbol (typically the name of the table).
- A hash map containing the structure of the schema.

Schemas are also responsible for populating the full-text search index with data from the relational database.

A concrete implementation for this protocol is provided by the `EntitySchema` type.

12.3 Implementation Issues

Chapter 13

Performance & Evaluation

13.1 Environment

All benchmarks were run on the same machine using the same software versions. The machine has the following specs.

Item	Description
Model	Dell PowerEdge 2950
Processor	2 x Quad-Core Intel® Xeon™ 5300 3.0GHz
Memory	8 GB DDR2 ECC RAM 667 MHz

Along with the following software versions.

Item	Description
Operating System	GNU/Linux Ubuntu 10.04.4 LTS
Kernel	2.6.32-28-generic
Java™ SE Runtime Environment	1.6.0_26-b03
Java HotSpot™ 64-Bit Server VM	20.1b02

13.2 Methodology

Evaluating the performance of code is a difficult problem at best. It is difficult to determine the impact on performance of various uncontrollable factors. Virtual Machines add another layer of abstraction which introduces even more factors.

As this project runs a top of the Java Virtual Machine (JVM), there are many factors to consider with regard to performance. In addition to uncontrollable events such as garbage collection, the JVM's behaviour can differ based on the operating system (OS) it is running on.

13.2.1 Pitfalls of the Java Virtual Machine

Reliable benchmarking on the JVM is a difficult problem. The standard HotSpot VM utilizes a non-deterministic garbage collector (GC). In addition, the JIT is also non-deterministic. These two problems combine into a larger issue.

Timing Execution

There are two methods provided by the JVM for getting a precise time from the OS; `System.currentTimeMillis()` and `System.nanoTime()`. The former returns the “wall” time. It is possible for time to leap forward or backward using this method. If daylight savings time occurs between calls to `System.currentTimeMillis()`, it can result in a negative time.

The latter uses a variety of methods in order to obtain the precise time. The method it uses depends on both the OS and the hardware itself. Under Windows, `System.nanoTime()` makes use of the `QueryPerformanceCounter(QPC)` call. This call may make use of the “programmable-interval-timer (PIT), or the ACPI power management timer (PMT), or the CPU-level timestamp-counter (TSC).” [4] Accessing the PIT and PMT is a slow operation. Accessing the TSC is a very fast operation, but may result in varying numbers [2].

Linux attempts to use the TSC when possible. If it finds the values to be unreliable (eg. different cores vary too much), it makes use of the High Precision Event Timer (HPET) [2].

Platform	Resolution (ms)
Windows 95/98	55
Windows NT, 2000, XP (single processor)	10
Windows XP (multi processor)	15.625
Linux 2.4 Kernel	10
Linux 2.6 Kernel	1

Table 13.1: Reported resolutions of `currentTimeMillis()` by platform [1].

The JVM's timing functions are one area where the JVM's behaviour differs based on the OS. Different operating systems provide varying degrees of time resolution. The speed at which the OS-level calls to these timing functions return can even differ.

Garbage Collection

Just-in-Time Compilation

A typical JVM only loads classes when they're first used [1]. This task typically involves disk I/O, along with extra processing. As a result, the initial run of a task may only...

13.2.2 Mitigating JVM Pitfalls

JIT Compilation

```
33
34 (defn -main
35   [& args]
36   (let [[opts arguments banner] (parse-args (flatten args))]
37     (when (or (opts :help) (not (opts :config)))
38       (println banner)
39       (System/exit 0))
40
41     (let [properties (load-props (opts :config))
```

Figure 13.1: core.clj

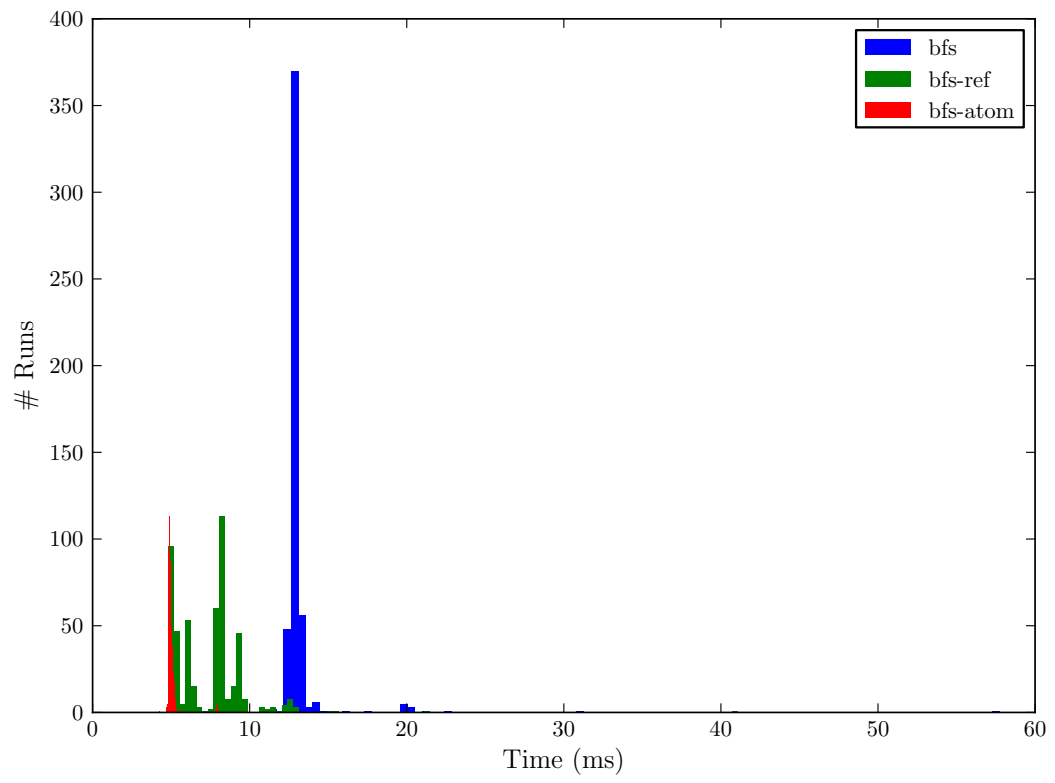


Figure 13.2: 1 Hop

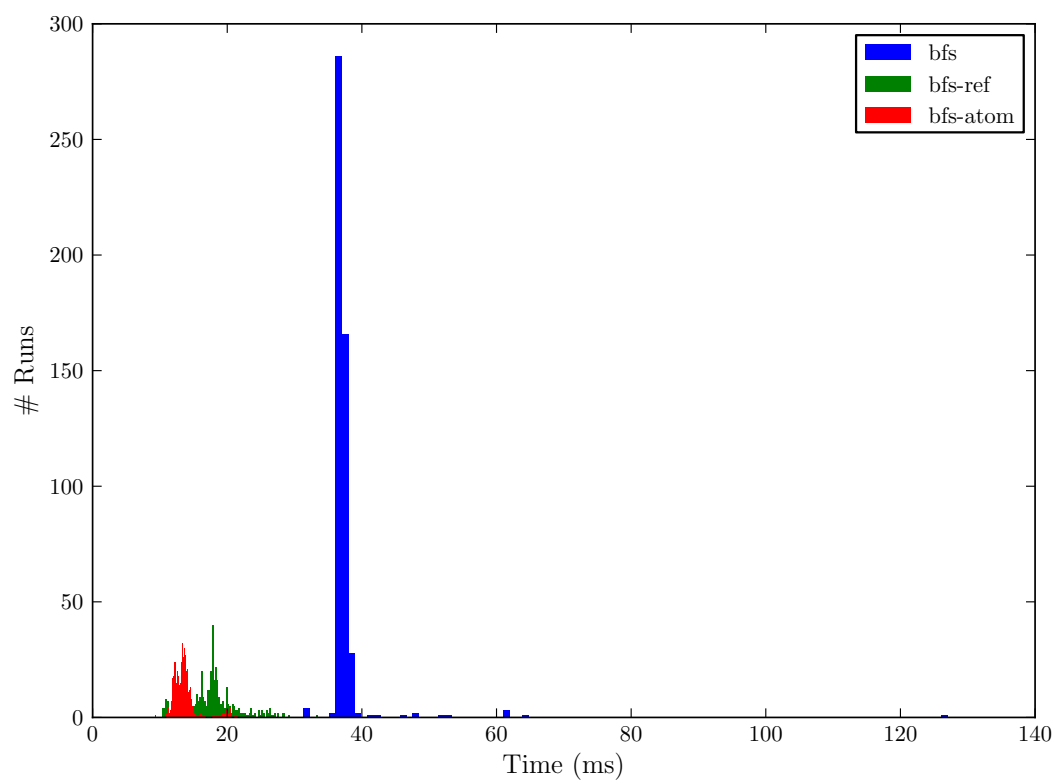


Figure 13.3: 2 Hops

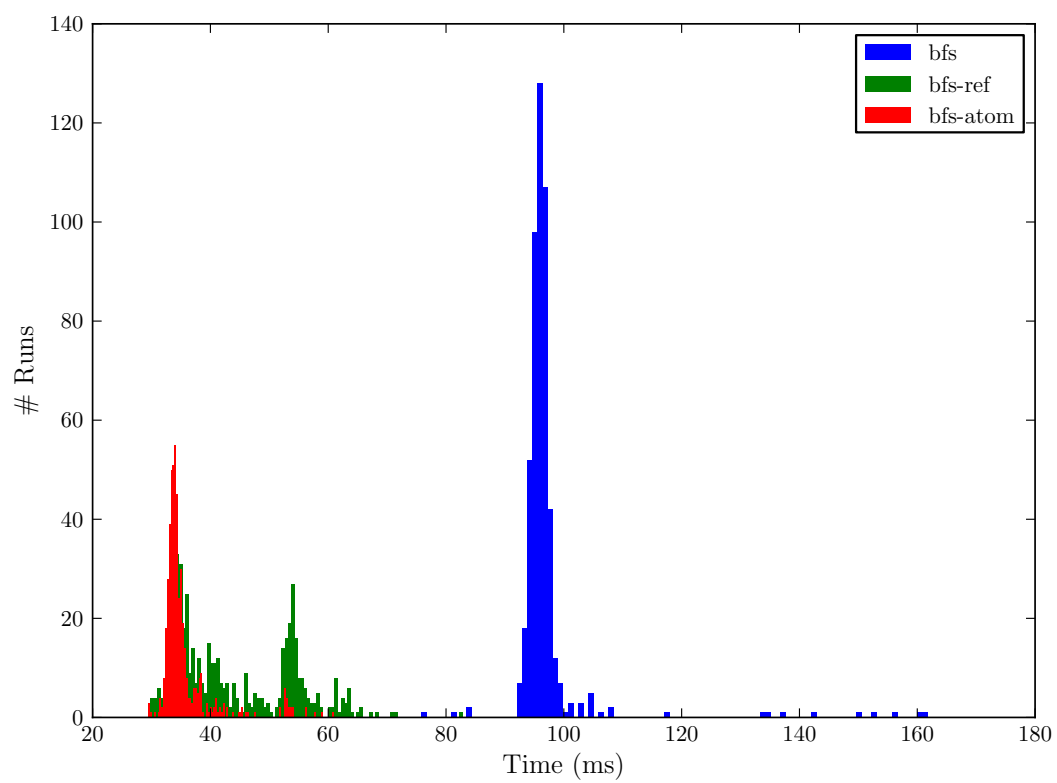


Figure 13.4: 3 Hops

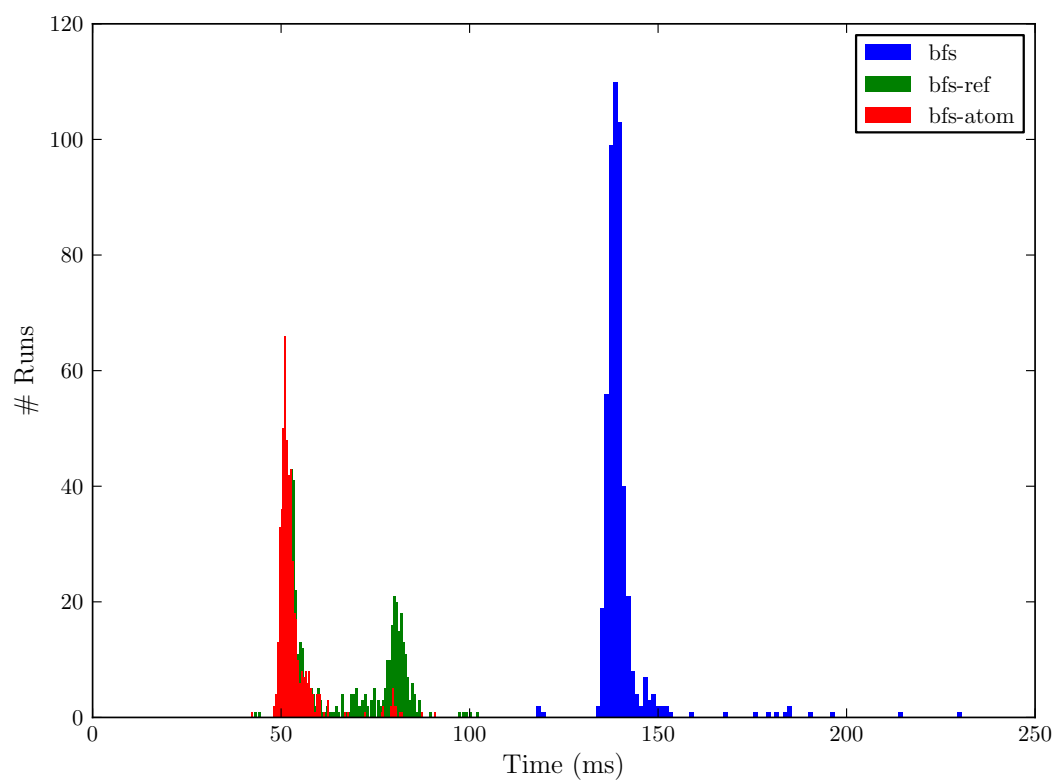


Figure 13.5: 4 Hops

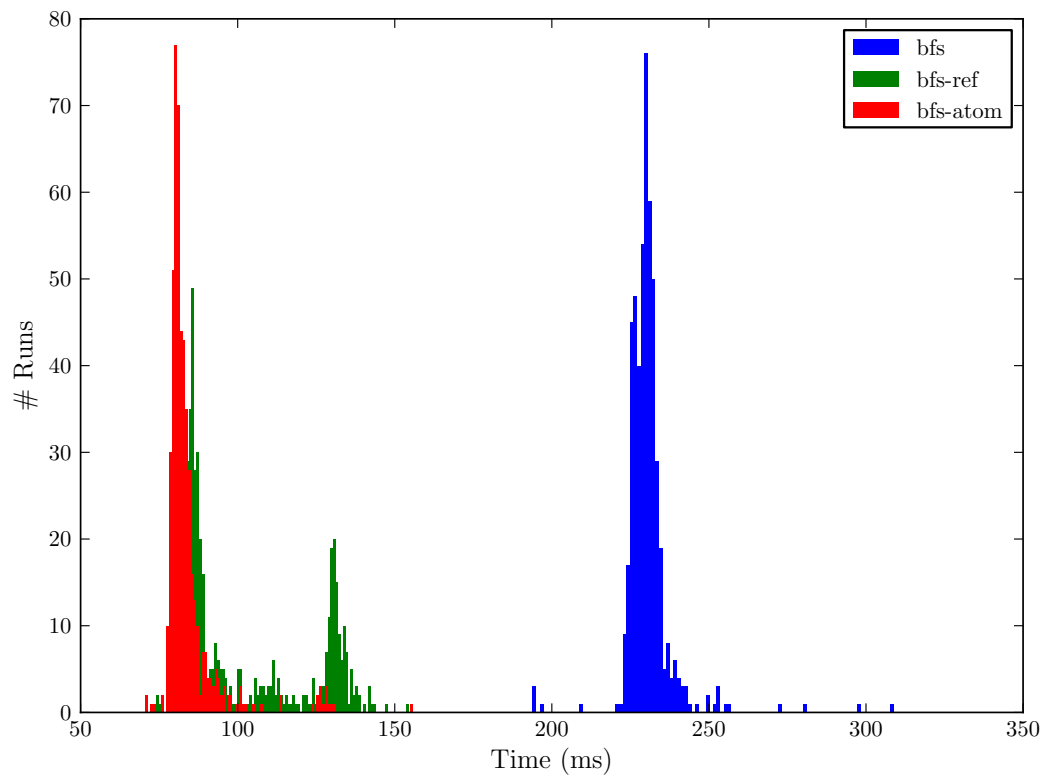


Figure 13.6: 5 Hops

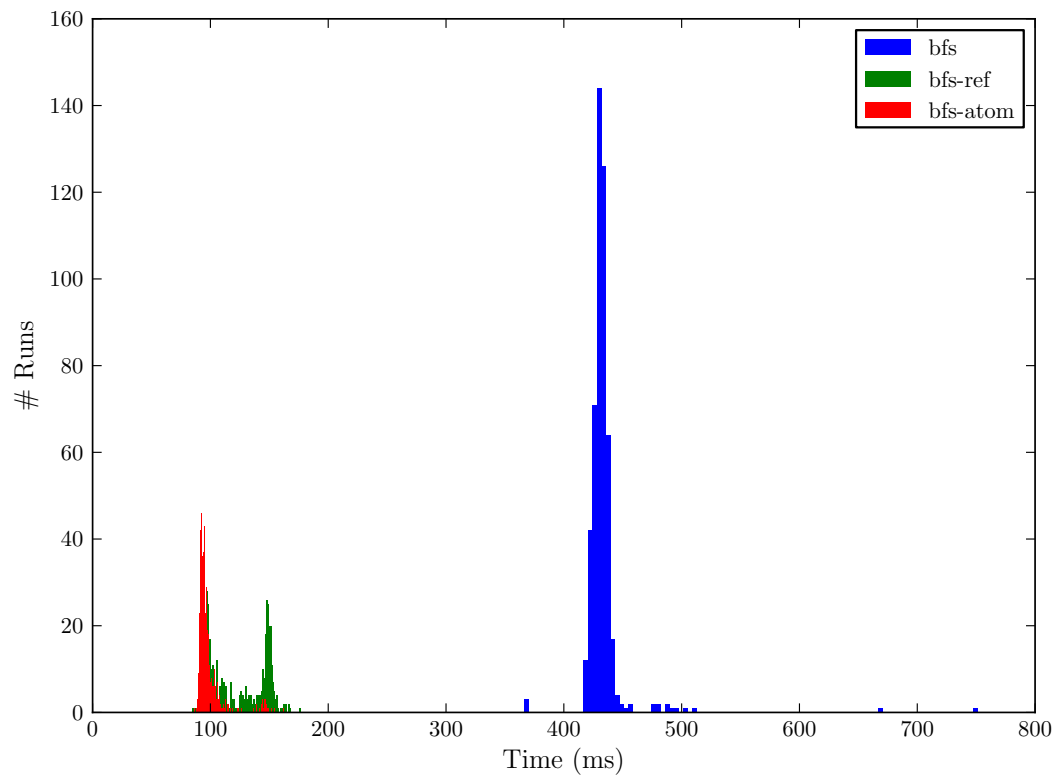


Figure 13.7: 6 Hops

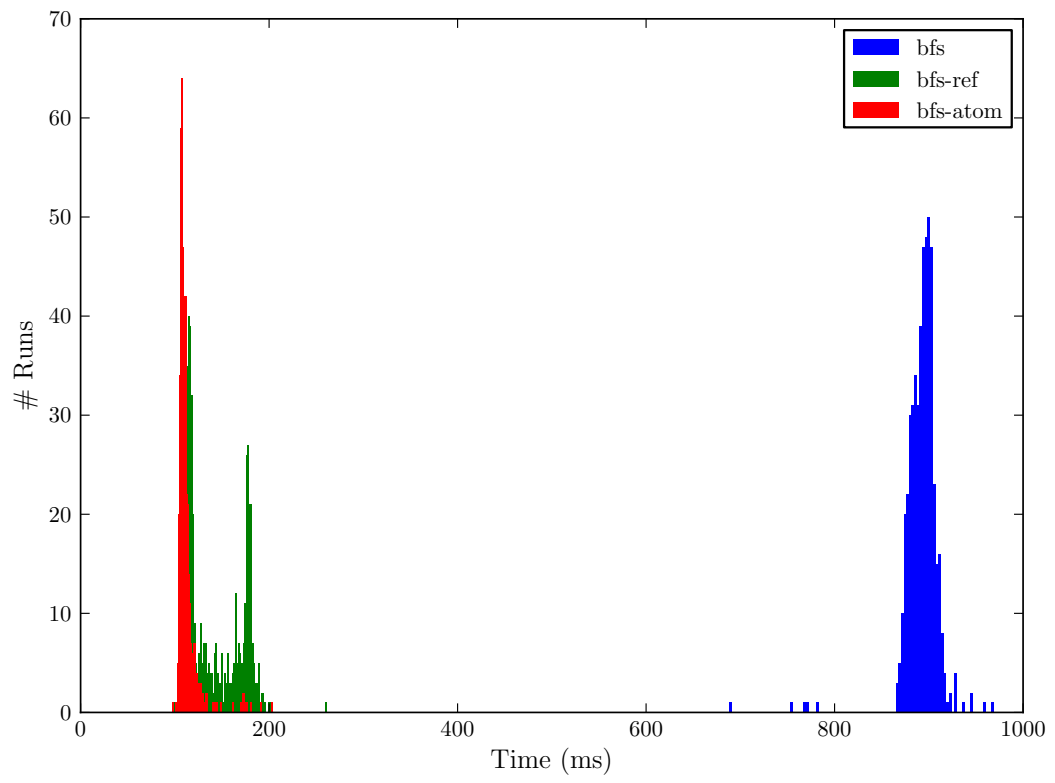


Figure 13.8: 7 Hops

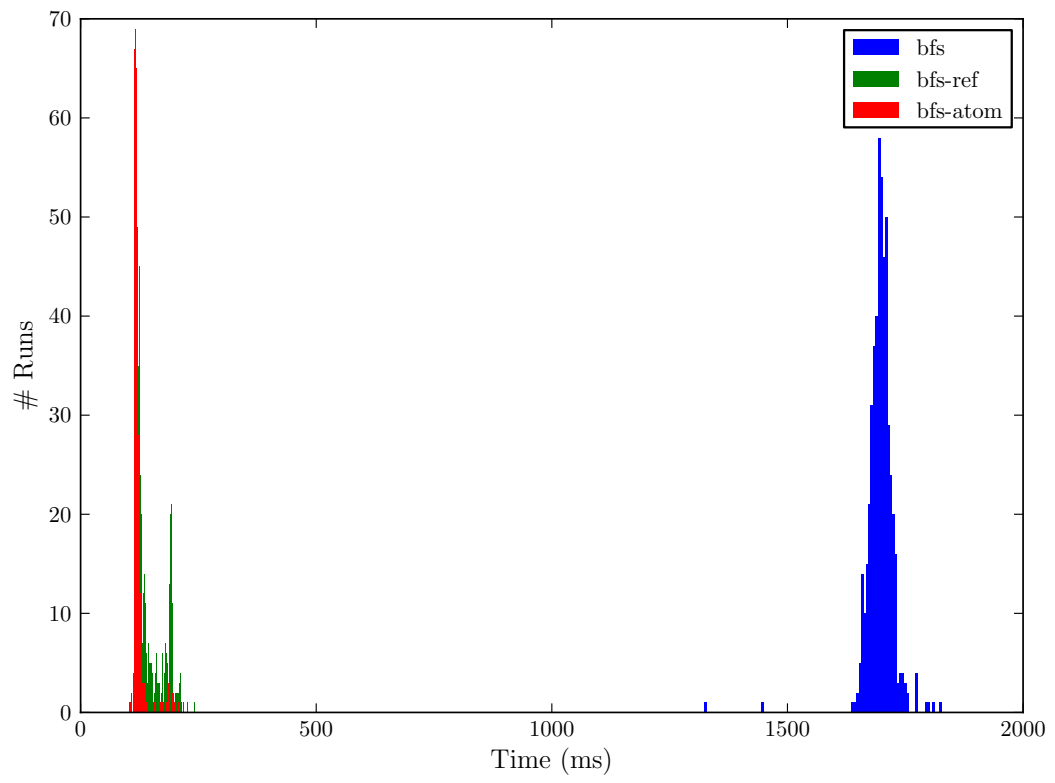


Figure 13.9: 8 Hops

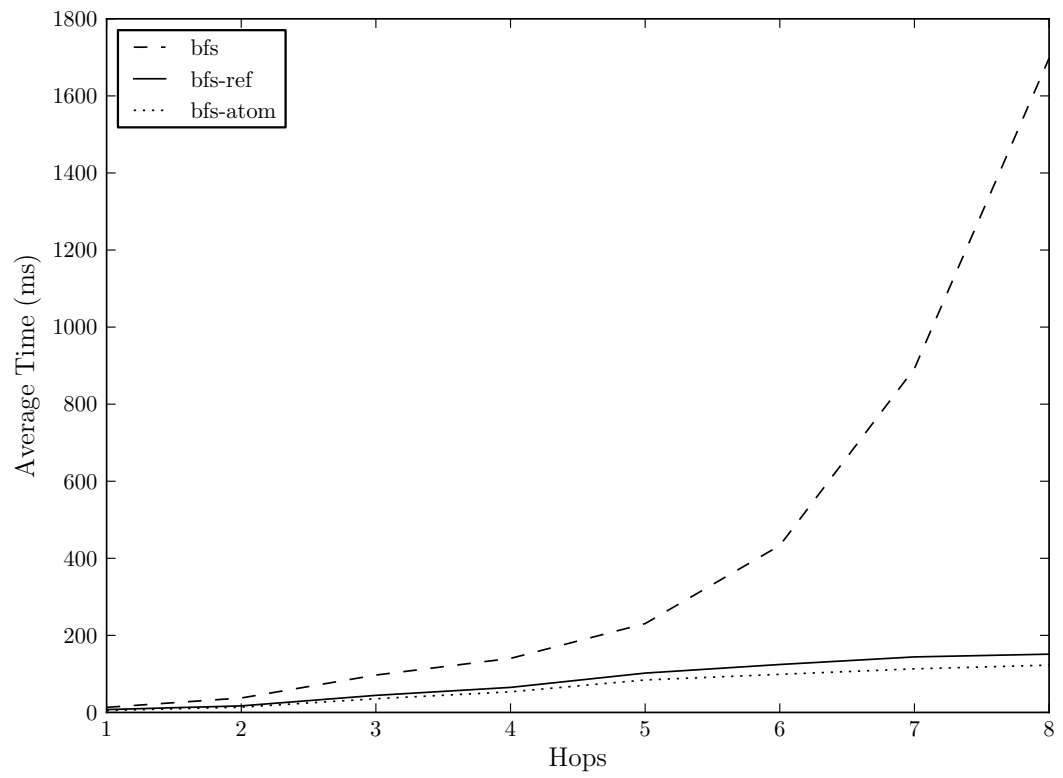


Figure 13.10: Comparison between single threaded and concurrent graph search

Chapter 14

Conclusion

[1]

Appendix A

Source Code

Each namespace in the code is divided into sections in the thesis document.

A.1 molly

A.1.1 molly.core

```
1 (ns molly.core
2   (:gen-class)
3   (:use molly.conf.config
4         molly.index.build
5         molly.search.lucene
6         [clojure.tools.cli :only (cli)]
7         [molly.algo.bfs-atom :only (bfs-atom)]
8         [molly.algo.bfs-ref :only (bfs-ref)]
9         [molly.algo.bfs :only (bfs)]
10        [molly.algo.ford-fulkerson :only (ford-fulkerson)]
11        [molly.bench.benchmark :only (benchmark-search)]))
12
13 (defn parse-args
14   [args]
15   (cli args
16     ["-c" "--config" "Path to configuration (properties) file"]
17     ["--algorithm" "Algorithm to run"]
18     ["-s" "--source" "Source node"]
19     ["-t" "--target" "Target node"]
20     ["--max-hops" "Maximum number of hops before stopping"]
21     ["--index" "Build an index of the database"]
22     :default false
23     :flag true]
24   ["--benchmark" "Run benchmarks"
25     :default false
26     :flag true])
```

```

27     ["-d" "--debug" "Displays additional information."
28     :default false
29     :flag true]
30     ["-h" "--help" "Show help"
31     :default false
32     :flag true]))
33
34 (defn -main
35   [& args]
36   (let [[opts arguments banner] (parse-args (flatten args))]
37     (when (or (opts :help) (not (opts :config)))
38       (println banner)
39       (System/exit 0)))
40
41     (let [properties (load-props (opts :config))
42           max-hops   (if (opts :max-hops)
43                         (Integer. (opts :max-hops))
44                         (properties :max-hops))]
45       (if (opts :index)
46         (let [database (properties :database)
47               index    (properties :index)]
48           (build database index))
49         nil)
50       (if (opts :algorithm)
51         (let [searcher (idx-searcher
52                        (idx-path
53                         (properties :index)))
54               source   (opts :source)
55               target   (opts :target)
56               f         (condp = (opts :algorithm)
57                            "bfs"          bfs
58                            "bfs-atom"     bfs-atom
59                            "bfs-ref"      bfs-ref
60                            "ford-fulkerson" ford-fulkerson
61                            (throw
62                             (Exception.
63                              "Not a valid algorithm choice.")))]
62           (if (opts :debug)
63             (let [[marked dist prev] (f searcher
64                                         source
65                                         target
66                                         max-hops)]
67               (println marked)
68               (println dist)
69               (println prev))
70             (println prev))
71

```

```
72         (benchmark-search f searcher source target max-hops))
73     (shutdown-agents))
74 nil))))
```

A.2 molly.conf

A.2.1 molly.conf.config

```
1 (ns molly.conf.config
2   (:require clojure.java.io))
3
4 (defn load-props
5   ([
6     (load-props ".properties"))
7     ([file-name]
8       (let [res (clojure.java.io/resource file-name)]
9         (if (nil? res)
10            (throw
11              (IllegalArgumentException. "Unable to load properties."))
12            (with-open [^java.io.Reader reader
13                        (clojure.java.io/reader res)]
14              (let [props (java.util.Properties.)]
15                (.load props reader)
16                (into {}
17                  (for [[k v] props] [(keyword k)
18                                     (read-string v)]))))))))))
19
20 (defprotocol IConfig
21   (connection [this])
22   (schema [this])
23   (index [this]))
```

A.2.2 molly.conf.mycampus

```
1 (ns molly.conf.mycampus
2   (:use molly.conf.config
3         molly.datatypes.database
4         molly.datatypes.schema
5         korma.core
6         korma.db)
7   (:import (molly.datatypes.database Sqlite)
8             (molly.datatypes.schema EntitySchema)))
9
10 (declare Campus Course Subject Term Section Schedule
11       Location Instructor db-conn)
12
13 (defentity Campus
14   (has-many Location))
15
16 (defentity Location
17   (belongs-to Campus))
18
19 (defentity Subject
20   (has-many Course))
21
22 (defentity Course
23   (pk :code)
24   (belongs-to Subject)
25   (has-many Section))
26
27 (defentity Instructor
28   (has-many Schedule))
29
30 (defentity Term
31   (has-many Section))
32
33 (defentity Section
34   (pk :crn)
35   (has-many Schedule)
36   (belongs-to Term)
37   (belongs-to Course {:fk :course_code}))
38
39 (defentity Schedule
40   (belongs-to Section)
41   (belongs-to Instructor)
42   (belongs-to Location))
```

```

43
44 (def mycampus-schema
45   [(EntitySchema.
46     {:T      :entity
47      :C      :course
48      :sql    Course
49      :ID     :code
50      :attrs  [:code :title]
51      :values [:code :title]})
52    (EntitySchema.
53      {:T      :entity
54       :C      :instructor
55       :sql    Instructor
56       :ID     :id
57       :attrs  [:name]
58       :values [:name]})
59    (EntitySchema.
60      {:T      :entity
61       :C      :location
62       :sql    Location
63       :ID     :id
64       :attrs  [:name]
65       :values [:name]})
66    (EntitySchema.
67      {:T      :entity
68       :C      :subject
69       :sql    Subject
70       :ID     :id
71       :attrs  [:id :name]
72       :values [:id :name]})
73    (EntitySchema.
74      {:T      :entity
75       :C      :campus
76       :sql    Campus
77       :ID     :id
78       :attrs  [:name]
79       :values [:name]})
80    (EntitySchema.
81      {:T      :entity
82       :C      :term
83       :sql    Term
84       :ID     :id
85       :attrs  [:id :name]
86       :values [:id :name]})
87    (EntitySchema.

```

```

88     {:T      :entity
89      :C      :section
90      :sql    Section
91      :ID     :crn
92      :attrs  [:crn :reg_start :reg_end :credits
93              :section_num :levels]
94      :values [:crn]})
95 (EntitySchema.
96  {:T      :entity
97   :C      :schedule
98   :sql    Schedule
99   :ID     :id
100  :attrs  [:days :sch_type :date_start :date_end
101           :time_start :time_end :week]
102  :values []})
103 (EntitySchema.
104  {:T      :group
105   :C      "Instructor schedule"
106   :sql    (->
107            (select* Schedule)
108            (with Instructor))
109   :ID     [[:instructor :instructor_id "Instructor ID"]
110           [:schedule :id "Schedule ID"]]
111   :attrs  []
112   :values []})
113 (EntitySchema.
114  {:T      :group
115   :C      "Course schedule"
116   :sql    (->
117            (select* Schedule)
118            (with Section
119             (with Course)))
120   :ID     [[:section :crn "CRN"]
121           [:course :code "Code"]
122           [:schedule :id "Schedule ID"]]
123   :attrs  []
124   :values []})
125 (EntitySchema.
126  {:T      :group
127   :C      "Schedule location"
128   :sql    (->
129            (select* Schedule)
130            (with Location
131             (with Campus)))
132   :ID     [[:campus :campus_id "Campus ID"]

```

```

133         [:location :location_id "Location ID"]
134         [:schedule :id "Schedule ID"]]]
135     :attrs []
136     :values []})
137 (EntitySchema.
138   {:T :group
139    :C "Course subject"
140    :sql (->
141          (select* Course)
142          (with Subject))
143    :ID [[[:course :id "Course"]
144          [:subject :subject_id "Subject"]]]
145    :attrs []
146    :values []})
147 (EntitySchema.
148   {:T :group
149    :C "Section term"
150    :sql (->
151          (select* Section)
152          (with Term))
153    :ID [[[:section :id "Section"]
154          [:term :term_id "Term"]]]})
155   ])
156
157 (deftype Mycampus [db-path idx-path]
158   IConfig
159   (connection
160     [this]
161     (defdb db-conn (sqlite3 {:db db-path}))
162     (Sqlite. db-conn))
163   (schema
164     [this]
165     mycampus-schema)
166   (index
167     [this]
168     idx-path))

```

A.3 molly.datatypes

A.3.1 molly.datatypes.database

```
1 (ns molly.datatypes.database
2   (:use korma.core
3         korma.db))
4
5 (defprotocol Database
6   (execute-query [this query f]))
7
8 (deftype Sqlite [conn]
9   Database
10  (execute-query
11    [this query f]
12    (with-db conn
13      (doseq [result (-> query (select))]
14        (f result)))))
```

A.3.2 molly.datatypes.entity

```
1 (ns molly.datatypes.entity
2   (:use molly.util.nlp)
3   (:import
4     [clojure.lang IPersistentMap IPersistentList]
5     [org.apache.lucene.document
6       Document Field
7       Field$Index Field$Store]))
8
9 (defn special?
10   [field-name]
11   (and (.startsWith field-name "__") (.endsWith field-name "__")))
12
13 (defn uid
14   "Possible inputs include:
15   row :T :ID
16   row [[:T :ID] [:T :ID]]
17   row [[:T :ID :desc] [:T :ID :desc]]"
18   ([row C id]
19     (if (nil? (row id))
20       (throw
21         (Exception.
22           (str "ID column " id " does not exist in row " row ".")))
23       (str (name C)
24            "|")
25       (clojure.string/replace (row id) #"\\s+" "__")))
26   ([row Tids]
27     (clojure.string/join " " (for [[C id] Tids]
28                                   (uid row C id)))))
29
30 (defn field
31   [field-name field-value]
32   (Field. field-name
33           field-value
34           Field$Store/YES
35           Field$Index/ANALYZED))
36
37 (defn document
38   [fields]
39   (let [doc (Document.)]
40     (do
41       (doseq [[field-name field-value] fields]
42         (.add doc (field (name field-name) (str field-value))))
```

```

43         doc)))
44
45 (defn row->data
46   ^{:doc "Transforms a row into the internal representation."}
47   [this schema]
48   (let [T      (schema :T)
49         C      (schema :C)
50         attr-cols (schema :attrs)
51         attrs    (if (nil? attr-cols)
52                     this
53                     (select-keys this attr-cols))
54         meta-data {:type T :class C}
55         id-col    (schema :ID)]
56     (with-meta (if (= T :group)
57                  (conj attrs {:entities (uid this id-col)})
58                  attrs)
59              (condp = T
60                  :value   (assoc meta-data
61                                   :class
62                                   (clojure.string/join "|"
63                                   (map name
64                                        [C (first attr-cols)])))
65                  :entity (assoc meta-data :id
66                                   (if (coll? id-col)
67                                       (uid this id-col)
68                                       (uid this C id-col)))
69                  :group  (assoc meta-data
70                                   :entities
71                                   (uid this id-col)))
72              (throw
73               (IllegalArgumentException.
74                "I only know how to deal with types :value,
75                :entity, and :group"))))))
76
77 (defn doc->data
78   ^{:doc "Transforms a Document into the internal representation."}
79   [this]
80   (let [fields      (.getFields this)
81         extract     (fn [x] [(keyword (clojure.string/replace
82                                         (.name x) "_" ""))
83                               (.stringValue x)])
84         check-special (fn [x] (special? (.name x)))
85         filter-fn    (fn [f] (apply hash-map
86                                     (flatten
87                                      (map extract

```

```

88                                     (filter f fields))))))]]
89   (with-meta (filter-fn (fn [x] (not (check-special x))))
90               (filter-fn check-special))))
91
92   (defn data->doc
93     ^{:doc "Transforms the internal representation into a Document."}
94     [this]
95     (let [int-meta (meta this)
96           T        (int-meta :type)
97           all       (clojure.string/lower-case
98                     (clojure.string/join " "
99                                           (if (= T :entity)
100                                             (conj (vals this)
101                                                  (name
102                                                    (int-meta :class)))
103                                             (vals this))))
104           luc-meta  [[:__type__ (name T)]
105                     [[:__class__ (name (int-meta :class))]]
106                     [[:__all__ (if (= T :value)
107                                   (q-gram all)
108                                   all))]]
109           raw-doc   (concat luc-meta
110                             this
111                             (condp = (int-meta :type)
112                               :value  [[:value all]]
113                               :entity [[:__id__ (int-meta :id)]]
114                               :group  [[]]))
115           (document raw-doc)))

```

A.3.3 molly.datatypes.schema

```
1 (ns molly.datatypes.schema
2   (:use molly.datatypes.database
3         molly.datatypes.entity
4         molly.search.lucene
5         molly.util.nlp
6         korma.core))
7
8 (defprotocol Schema
9   (crawl [this db-conn idx-w])
10  (klass [this])
11  (schema-map [this]))
12
13 (deftype EntitySchema [S]
14   Schema
15   (crawl
16     [this db-conn idx-w]
17     (let [sql (S :sql)]
18       (execute-query db-conn sql
19         (fn [row]
20           (add-doc idx-w
21             (data->doc (row->data row S)))))))
22
23     (if (= (S :T) :entity)
24       (doseq [value (S :values)]
25         (let [query (->
26                   sql
27                   (modifier "DISTINCT")
28                   (fields value)
29                   (group value))]
30           (execute-query db-conn query
31             (fn [row]
32               (add-doc idx-w (data->doc
33                             (row->data row
34                               (assoc S
35                                 :T :value)))))))))))
36   (klass
37     [this]
38     ((schema-map this) :C))
39   (schema-map
40     [this]
41     S))
```

A.4 molly.index

A.4.1 molly.index.build

```
1 (ns molly.index.build
2   (:use molly.conf.config
3         molly.conf.mycampus
4         molly.datatypes.database
5         molly.datatypes.entity
6         molly.datatypes.schema
7         molly.search.lucene)
8   (:import (molly.conf.mycampus Mycampus)))
9
10 (defn build
11   [db-path path]
12   (let [conf      (Mycampus. db-path path)
13         db-conn   (connection conf)
14         ft-path   (idx-path (index conf))
15         idx-w     (idx-writer ft-path)
16         schemas   (schema conf)]
17     (doseq [ent-def schemas]
18       (println "Indexing" (name (klass ent-def)) "...")
19       (crawl ent-def db-conn idx-w))
20
21     (close-idx-writer idx-w)))
```

A.5 molly.util

A.5.1 molly.util.nlp

```
1 (ns molly.util.nlp)
2
3 (defn q-gram
4   ([S]
5    (q-gram S 3 "$"))
6   ([S n]
7    (q-gram S n "$"))
8   ([S n s]
9    (let [padding (clojure.string/join "" (repeat (dec n) s))
10          padded-S (str padding
11                        (clojure.string/replace S " " padding)
12                        padding)]
13      (clojure.string/join " "
14                            (for [i (range
15                                    (+ 1 (- (count padded-S) n)))]
16                              (. padded-S substring i (+ i n)))))))
```

A.6 molly.search

A.6.1 molly.search.lucene

```
1 (ns molly.search.lucene
2   (:import
3     (java.io File)
4     (org.apache.lucene.analysis.core WhitespaceAnalyzer)
5     (org.apache.lucene.index IndexReader IndexWriter
6                                   IndexWriterConfig)
7     (org.apache.lucene.search IndexSearcher)
8     (org.apache.lucene.store Directory SimpleFSDirectory)
9     (org.apache.lucene.util Version)))
10
11 (def version
12   Version/LUCENE_44)
13 (def default-analyzer
14   (WhitespaceAnalyzer. version))
15
16 (defn ^Directory idx-path
17   [path]
18   (-> path File. SimpleFSDirectory.))
19
20 (defn idx-searcher
21   [^IndexSearcher idx-path]
22   (-> (IndexReader/open idx-path) IndexSearcher.))
23
24 (defn ^IndexWriter idx-writer
25   ([^Directory idx-path analyzer]
26    (IndexWriter. idx-path (IndexWriterConfig. version analyzer)))
27   ([^Directory idx-path]
28    (idx-writer idx-path default-analyzer)))
29
30 (defn close-idx-writer
31   [^IndexWriter idx-writer]
32   (doto idx-writer
33     (.commit)
34     (.close)))
35
36 (defn idx-search
37   [idx-searcher query topk]
38   (let [results (. (. idx-searcher search query topk) scoreDocs)]
39     (map (fn [result] (.doc idx-searcher (.doc result))) results)))
40
```

```
41 (defn add-doc  
42   [idx doc]  
43   (. idx addDocument doc))
```

A.6.2 molly.search.query_builder

```
1 (ns molly.search.query-builder
2   (:import (org.apache.lucene.index Term)
3             (org.apache.lucene.search BooleanClause$Occur
4                                         BooleanQuery
5                                         PhraseQuery)))
6
7 (defn query
8   [kind & args]
9   (let [field-name (condp = kind
10                      :type    "__type__"
11                      :class   "__class__"
12                      :id      "__id__"
13                      :text     "__all__"
14                      ; Assume "kind" is an attribute name.
15                      (condp = (type kind)
16                            clojure.lang.Keyword (name kind)
17                            java.lang.String      kind))
16     phrase-query (PhraseQuery.)]
17     (doseq [arg args]
18       (. phrase-query add (Term. field-name (name arg)))))
19     phrase-query))
20
21
22
23
24 (defn boolean-query
25   [args]
26   (let [query (BooleanQuery.)]
27     (doseq [[q op] args]
28       (. query add q (condp = op
29                        :and BooleanClause$Occur/MUST
30                        :or  BooleanClause$Occur/SHOULD
31                        :not BooleanClause$Occur/MUST_NOT)))
32     query))
```

A.7 molly.server

A.7.1 molly.server.serve

```
1 (ns molly.server.serve
2   (:use compojure.core
3         molly.conf.config
4         molly.datatypes.entity
5         molly.search.lucene
6         molly.search.query-builder
7         molly.util.nlp
8         [molly.algo.bfs :only (bfs)]
9         [molly.algo.bfs-atom :only (bfs-atom)]
10        [molly.algo.bfs-ref :only (bfs-ref)])
11   (:require [noir.response :as response]
12             [compojure.handler :as handler]
13             [compojure.route :as route]))
14
15 (def runtime (Runtime/getRuntime))
16 (def props (load-props ".properties"))
17 (def searcher (idx-searcher (idx-path (props :index))))
18
19 (defn dox
20   [q1 field S op topk]
21   (let [bq (boolean-query
22           (concat [[q1 :and]]
23                   (for [s S]
24                       [(query field s) op])))
25         result (map doc->data (idx-search searcher bq topk))
26         fmt (fn [data] {:meta (meta data) :results data})]
27     (map fmt result)))
28
29 (defn entities
30   [field q topk]
31   (dox (query :type :entity)
32        field
33        (clojure.string/split q #"\\s{1}")
34        :and
35        topk))
36
37 (defn home-page []
38   (response/redirect "/index.html"))
39
40 (defn get-value [q]
```

```

41 (response/json
42   {:result
43     (dox (query :type :value)
44           :text
45           (clojure.string/split (q-gram q) #"\\s{1}")
46           :or
47           (props :topk_value)))))
48
49 (defn get-entities [q]
50   (response/json
51     {:result
52       (entities
53         :text (clojure.string/lower-case q) (props :topk_entities)))))
54
55 (defn get-entity [q]
56   (response/json
57     {:result
58       (entities :id q (props :topk_entity))}))
59
60 (defn get-span [e0 eL method]
61   (let [start (System/nanoTime)
62         [visited dist prev]
63         (condp = method
64           "atom" (bfs-atom searcher e0 eL)
65           "ref" (bfs-ref searcher e0 eL)
66           (bfs searcher e0 eL))
67         t (- (System/nanoTime) start)
68         eids (for [[k v] prev] k)
69         get-entities (fn [eid]
70                       {(keyword eid)
71                        (entities :id eid
72                                (props :topk_entity))})
73         entities (into {} (map get-entities eids))]
74     (response/json
75       {:from e0
76        :to eL
77        :prev prev
78        :entities entities
79        :debug {:time t
80                :mem_total (.totalMemory runtime)
81                :mem_free (.freeMemory runtime)
82                :mem_used (- (.totalMemory runtime)
83                             (.freeMemory runtime))
84                :properties props}})))
85

```



```
86 (defroutes app-routes
87   (GET "/" [] (home-page))
88   (GET "/value" [q] (get-value q))
89   (GET "/entities" [q] (get-entities q))
90   (GET "/entity" [q] (get-entity q))
91   (GET "/span" [e0 eL method] (get-span e0 eL method))
92   (route/files "/" {:root "resources/public"}))
93
94 (def app (handler/site app-routes))
```

A.8 molly.algo

A.8.1 molly.algo.common

```
1 (ns molly.algo.common
2   (:use molly.datatypes.entity
3         molly.search.lucene
4         molly.search.query-builder))
5
6 (defn find-entity-by-id
7   [G id]
8   (let [query (boolean-query [[(query :type :entity) :and]
9                                [(query :id id) :and]])]
10     (map doc->data (idx-search G query 10))))
11
12 (defn find-group-for-id
13   [G id]
14   (let [query (boolean-query [[(query :type :group) :and]
15                                [(query :entities id) :and]])
16         results (map doc->data (idx-search G query 10))
17         big-str (clojure.string/join " "
18                                       (map #(% :entities) results))]
19     (distinct (clojure.string/split big-str #"s{1}"))))
20
21 (defn find-adj
22   [G v]
23   (remove #{v} (find-group-for-id G v)))
24
25 (defn initial-state
26   [s]
27   {:Q      (-> (clojure.lang.PersistentQueue/EMPTY) (conj s))
28    :marked #{s}
29    :dist   {s 0}
30    :prev   {}
31    :done   false})
32
33 (defn update-state
34   [state u v max-hops]
35   (let [Q      (state :Q)
36         marked (state :marked)
37         dist   (state :dist)
38         prev   (state :prev)
39         done   (> (dist u) max-hops)]
40     (assoc state
```

```
41         :Q      (if done
42                   Q
43                   (conj Q v))
44         :marked  (conj marked v)
45         :dist    (assoc dist v (inc (dist u)))
46         :prev    (assoc prev v u)
47         :done    done)))
48
49 (defn deref-future
50   [dfd]
51   (if (future? dfd)
52       (deref dfd)
53       dfd))
```

A.8.2 molly.algo.bfs

```
1 (ns molly.algo.bfs
2   (use molly.algo.common))
3
4 (defn update-adj
5   [G marked dist prev u max-hops]
6   (loop [adj      (find-adj G u)
7          marked    marked
8          dist      dist
9          prev      prev
10         frontier []]
11     (if (or (empty? adj) (>= (dist u) max-hops))
12         [(conj marked u) dist prev frontier]
13         (let [v      (first adj)
14               adj'    (rest adj)]
15             (if (marked v)
16                 (recur adj' marked dist prev frontier)
17                 (let [dist'      (assoc dist v (inc (dist u)))
18                       prev'      (assoc prev v u)
19                       frontier'  (conj frontier v)]
20                     (recur adj' marked dist' prev' frontier'))))))))
21
22 (defn bfs
23   [G s t max-hops]
24   (loop [Q      (-> (clojure.lang.PersistentQueue/EMPTY) (conj s))
25          marked #{ }
26          dist   {s 0}
27          prev   {s nil}]
28       (if (or (empty? Q)
29               (some (fn [node] (= node t)) marked))
30           [marked dist prev]
31           (let [u      (first Q)
32                 Q'     (rest Q)
33                 [marked' dist' prev' frontier]
34                 (update-adj G marked dist prev u max-hops)]
35             (recur (concat Q' frontier) marked' dist' prev')))))
```

A.8.3 molly.algo.bfs_atom

```
1 (ns molly.algo.bfs-atom
2   (use molly.algo.common))
3
4 (defn update-adj
5   [state-ref G u max-hops]
6   (let [marked? (@state-ref :marked)
7         deferred (if (>= ((@state-ref :dist) u) max-hops)
8                       []
9                       (doall
10                        (for [v (find-adj G u)]
11                          (if (marked? v)
12                              nil
13                              (future
14                               (swap!
15                                state-ref
16                                update-state
17                                u
18                                v
19                                max-hops))))))]
20     (doall (map deref-future deferred))))
21
22 (defn bfs-atom
23   [G s t max-hops]
24   (let [state-ref (atom (initial-state s))]
25     (while (and (not (empty? (@state-ref :Q)))
26                 (not (@state-ref :done)))
27       (let [u (first (@state-ref :Q))
28             Q' (pop (@state-ref :Q))]
29         (swap! state-ref assoc :Q Q')
30         (if (some (fn [node] (= node t)) (@state-ref :marked))
31             (swap! state-ref assoc :done true)
32             (update-adj state-ref G u max-hops)))
33     [(@state-ref :marked) (@state-ref :dist) (@state-ref :prev)]))
```

A.8.4 molly.algo.bfs_ref

```
1 (ns molly.algo.bfs-ref
2   (use molly.algo.common))
3
4 (defn update-adj
5   [state-ref G u max-hops]
6   (let [marked? (@state-ref :marked)
7         deferred (if (>= ((@state-ref :dist) u) max-hops)
8                       []
9                       (doall
10                        (for [v (find-adj G u)]
11                          (if (marked? v)
12                              nil
13                              (future (dosync (alter
14                                             state-ref
15                                             update-state
16                                             u
17                                             v
18                                             max-hops)))))))]
19     (doall (map deref-future deferred))))
20
21 (defn bfs-ref
22   [G s t max-hops]
23   (let [state-ref (ref (initial-state s))]
24     (while (and (not (empty? (@state-ref :Q)))
25                 (not (@state-ref :done)))
26       (let [u (first (@state-ref :Q))
27             Q' (pop (@state-ref :Q))]
28         (dosync (alter state-ref assoc :Q Q'))
29         (if (some (fn [node] (= node t)) (@state-ref :marked))
30             (dosync (alter state-ref assoc :done true))
31             (update-adj state-ref G u max-hops)))
32     [(@state-ref :marked) (@state-ref :dist) (@state-ref :prev)]))
```

A.9 molly.bench

A.9.1 molly.bench.benchmark

```
1 (ns molly.bench.benchmark
2   (use criterium.core)
3   (require [clojure.data.json :as json]))
4
5 (defn benchmark-search
6   [f G s t max-hops]
7   (let [method (last (clojure.string/split (str (class f)) #"\\$"))
8         result
9         (dissoc
10          (benchmark (f G s t max-hops) {:verbose false})
11          :results)]
12     (println
13      (json/write-str
14       {:method      method
15        :max-hops    max-hops
16        :results     result}))))
```

Bibliography

- [1] Brent Boyer. *Robust Java benchmarking, Part 1: Issues*. June 2008. URL: <http://www.ibm.com/developerworks/java/library/j-benchmark1/index.html>.
- [2] Jonathan Corbet. *Counting on the time stamp counter*. Nov. 2006. URL: <http://lwn.net/Articles/209101/>.
- [3] Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. *Database systems - the complete book (2. ed.)*. Pearson Education, 2009.
- [4] David Holmes. *Inside the Hotspot VM: Clocks, Timers and Scheduling Events - Part I - Windows*. Oct. 2006. URL: https://blogs.oracle.com/dholmes/entry/inside_the_hotspot_vm_clocks.
- [5] *The Xapian Project : Features*. URL: <http://xapian.org/features>.

To do...

- ☐ 1 (p. 4): Diagram of schema-level graph (FKs, relations)
- ☐ 2 (p. 4): More examples of queries
- ☐ 3 (p. 5): Add overview diagram of system design.