

MOLLY

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

August 2013

Copyright © Richard J.I. Drake 2013

Abstract

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

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	RDBMS to Document Store	1
2	Search	10
3	Concurrency	14
4	Literature Review	18
5	Theoretical	19
5.1	Data Representation & Notation	20
5.2	Relational Database Abstraction	21
5.3	Data Corpus	21
5.4	Graph Search Algorithms	23
6	Implementation	25
6.1	Functional vs. Procedural Programming of Algorithms	25
6.1.1	Functional Data Structures	25
6.2	Tuneable Parameters	25
7	System Implementation	26
7.1	Technology Selection	26
7.1.1	Programming Language	26
7.1.2	Relational Database	27

7.1.3	Full-Text Search Database	28
7.1.4	Web Stack	28
7.2	System Design	28
7.2.1	Configuration	29
7.3	Implementation Issues	30
8	Performance & Evaluation	31
8.1	Environment	31
8.2	Methodology	31
8.2.1	Pitfalls of the Java Virtual Machine	32
8.2.2	Mitigating JVM Pitfalls	33
9	Conclusion	43

List of Tables

5.1	Courses entity schema	21
5.2	Instructors entity schema	22
5.3	Sections entity schema	22
5.4	Schedules entity schema	23
5.5	Teaches entity schema	23
8.1	Reported resolutions of <code>currentTimeMillis()</code> by platform [1].	32

List of Figures

1.1	core.clj	2
1.2	conf/mycampus.clj	3
1.3	datatypes/database.clj	4
1.4	datatypes/entity.clj	5
1.5	datatypes/index.clj	6
1.6	datatypes/schema.clj	7
1.7	index/build.clj	8
1.8	util/nlp.clj	9
2.1	search/lucene.clj	11
2.2	search/query_builder.clj	12
2.3	server/serve.clj	13
3.1	algo/common.clj	14
3.2	algo/bfs.clj	15
3.3	algo/bfs_atom.clj	16
3.4	algo/bfs_ref.clj	17
5.1	The structure of an entity	20
5.2	The structure of an entity group	21
8.1	core.clj	33
8.2	1 Hop	34

8.3	2 Hops	35
8.4	3 Hops	36
8.5	4 Hops	37
8.6	5 Hops	38
8.7	6 Hops	39
8.8	7 Hops	40
8.9	8 Hops	41
8.10	Comparison between single threaded and concurrent graph search	42

List of Algorithms

Chapter 1

RDBMS to Document Store

```

1  (ns molly.core
2    (:gen-class)
3    (:use clojure.pprint
4          molly.conf.config
5          molly.server.serve
6          molly.index.build
7          molly.search.lucene
8          molly.search.query-builder
9          [clojure.tools.cli :only (cli)]
10         [molly.algo.bfs-atom :only (bfs-atom)]
11         [molly.algo.bfs-ref :only (bfs-ref)]
12         [molly.algo.bfs :only (bfs)]))
13
14  (defn parse-args
15    [args]
16    (cli args
17      ["-c" "--config" "Path to configuration (properties) file"]
18      ["--serve"      "Start the server"
19       :default false
20       :flag true]
21      ["--port"       "Port to run the server on"
22       :default 8080
23       :parse-fn #(Integer. %)]
24      ["--algorithm"  "Algorithm to run"]
25      ["-s" "--source" "Source node"]
26      ["-t" "--target" "Target node"]
27      ["--index"      "Build an index of the database"
28       :default false
29       :flag true]
30      ["--benchmark"  "Run benchmarks"
31       :default false
32       :flag true]
33      ["-h" "--help"   "Show help"
34       :default false
35       :flag true]))
36
37  (def counter (atom 0))
38
39  (defn ns-to-ms
40    [nanosec]
41    (/ nanosec (* 1000.0 1000.0)))
42
43  (defn warmup
44    [func searcher source target seconds]
45    (let [start (System/nanoTime)
46          time-in-s (* seconds 1000 1000 1000)]
47      (while (< (- (System/nanoTime) start) time-in-s)
48        (func searcher source target)
49        (swap! counter inc)))
50

```

```

1  (ns molly.conf.mycampus
2    (:use molly.conf.config
3          molly.datatypes.database
4          molly.datatypes.index
5          molly.datatypes.schema
6          [clojureql.core :only (table project join where rename)]))
7  (:import (molly.datatypes.database Sqlite)
8           (molly.datatypes.index Lucene)
9           (molly.datatypes.schema EntitySchema)))
10
11 (def schedules-table
12   (->
13     (join
14       (->
15         (table :sections)
16         (project [[:id :as :sec_id] [:id :as :section_id] :actual
17                  :campus :capacity :credits :levels
18                  :registration_start :registration_end :semester
19                  :sec_code :sec_number :year :course]))
20       (->
21         (join
22           (->
23             (table :schedules)
24             (project [[:id :as :sch_id] :date_start :date_end :day
25                      :schedtype :hour_start :hour_end :min_start
26                      :min_end :classtype :location]))
27           (->
28             (table :teaches)
29             (project [:position :teaches.schedule_id
30                      :teaches.instructor_id [:id :as :teaches_id]]))
31           (where (= :schedules.id :teaches.schedule_id))))
32       (where (= :sections.id :section_id))))))
33
34 (def schedules-id
35   [[:schedules :schedule_id]
36    [:teaches :teaches_id]
37    [:sections :section_id]])
38
39 (def mycampus-schema
40   [(EntitySchema.
41     {:T      :entity
42      :C      :courses
43      :sql     (table :courses)
44      :ID      :code
45      :attrs   [:code :title :description]
46      :values  [:code :title]})
47    (EntitySchema.
48      3
49      {:T      :entity
50       :C      :instructors
51       :sql     (table :instructors)

```

```
1 (ns molly.datatypes.database
2   (:require
3     [clojure.java.jdbc :as sql]
4     [clojureql.core :as cql]))
5
6 (defprotocol Database
7   (execute-query [this query f]))
8
9 (deftype Sqlite [conn]
10  Database
11  (execute-query
12    [this query f]
13    (sql/with-connection conn (cql/with-results [rs query]
14                                              (doseq [res rs]
15                                                (f res))))))
```

Figure 1.3: datatypes/database.clj

```

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  (defn field
30    [field-name field-value]
31    (Field. field-name
32            field-value
33            Field$Store/YES
34            Field$Index/ANALYZED))
35
36  (defn document
37    [fields]
38    (let [doc (Document.)]
39      (do
40        (doseq [[field-name field-value] fields]
41          (.add doc (field (name field-name) (str field-value))))
42        doc)))
43
44  (defn row->data
45    ^{:doc "Transforms a row into the internal representation."}
46    [this schema]
47    (let [T (schema :T)      5
48          C (schema :C)
49          attr-cols (schema :attrs)
50          attrs (if (nil? attr-cols)

```

```

1 (ns molly.datatypes.index
2   (:import
3     (java.io File)
4     (org.apache.lucene.analysis WhitespaceAnalyzer)
5     (org.apache.lucene.index IndexReader
6                                   IndexWriter
7                                   IndexWriter$MaxFieldLength)
8     (org.apache.lucene.queryParser QueryParser)
9     (org.apache.lucene.search IndexSearcher)
10    (org.apache.lucene.store SimpleFSDirectory)
11    (org.apache.lucene.util Version)))
12
13 (def version
14   Version/LUCENE_35)
15 (def default-analyzer
16   (WhitespaceAnalyzer. version))
17 (def unlimited-fields IndexWriter$MaxFieldLength/UNLIMITED)
18
19 (defprotocol Index
20   (path [this])
21   (open-searcher [this])
22   (open-writer
23     [this]
24     [this analyzer])
25   (close-writer [this idx])
26   (search
27     [this searcher query]
28     [this searcher query analyzer])
29   (add-doc [this idx doc]))
30
31 (deftype Lucene [idx-path]
32   Index
33   (path
34     [this]
35     (-> idx-path File. SimpleFSDirectory.))
36   (open-searcher
37     [this]
38     (-> (IndexReader/open (path this)) IndexSearcher.))
39   (open-writer
40     [this]
41     (open-writer this default-analyzer))
42   (open-writer
43     [this analyzer]
44     (IndexWriter. (path this) analyzer unlimited-fields))
45   (close-writer
46     [this idx]
47     (doto idx
48       (.commit)
49       (.optimize)
50       (.close)))

```

```

1  (ns molly.datatypes.schema
2    (:use molly.datatypes.database
3          molly.datatypes.entity
4          molly.search.lucene
5          molly.util.nlp)
6    (:require [clojureql.core :as cql]))
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        (println sql)
19        (execute-query db-conn sql
20                       (fn [row]
21                         (add-doc idx-w
22                                (data->doc (row->data row S)))))))
23
24    (if (= (S :T) :entity)
25      (doseq [value (S :values)]
26        (let [query (->
27                  sql
28                  (cql/project [value])
29                  (cql/grouped [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                                                )))))))))
37    (klass
38      [this]
39      ((schema-map this) :C))
40    (schema-map
41      [this]
42      S))

```

Figure 1.6: datatypes/schema.clj

```

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 (class ent-def)) "...")
19       (crawl ent-def db-conn idx-w))
20
21     (close-idx-writer idx-w)))

```

Figure 1.7: index/build.clj

```

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) "$"))
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))))))

```

Figure 1.8: util/nlp.clj

Chapter 2

Search

```

1 (ns molly.search.lucene
2   (:import
3     (java.io File)
4     (org.apache.lucene.analysis WhitespaceAnalyzer)
5     (org.apache.lucene.index IndexReader IndexWriter
6       IndexWriter$MaxFieldLength)
7     (org.apache.lucene.queryParser QueryParser)
8     (org.apache.lucene.search IndexSearcher)
9     (org.apache.lucene.store SimpleFSDirectory)
10    (org.apache.lucene.util Version)))
11
12 (def version
13   Version/LUCENE_35)
14 (def default-analyzer
15   (WhitespaceAnalyzer. version))
16 (def unlimited-fields IndexWriter$MaxFieldLength/UNLIMITED)
17
18 (defn idx-path
19   [path]
20   (-> path File. SimpleFSDirectory.))
21
22 (defn idx-searcher
23   [idx-path]
24   (-> (IndexReader/open idx-path) IndexSearcher.))
25
26 (defn idx-writer
27   ([idx-path analyzer]
28    (IndexWriter. idx-path analyzer unlimited-fields))
29   ([idx-path]
30    (idx-writer idx-path default-analyzer)))
31
32 (defn close-idx-writer
33   [idx-writer]
34   (doto idx-writer
35     ; Lucene docs say not to use .optimize anymore.
36     (.commit)
37     (.close)))
38
39 (defn idx-search
40   [idx-searcher query topk]
41   (let [results (. (. idx-searcher search query topk) scoreDocs)]
42     (map (fn [result] (.doc idx-searcher (.doc result))) results)))
43
44 (defn add-doc
45   [idx doc]
46   (. idx addDocument doc))

```

```

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

```

Figure 2.2: search/query_builder.clj

```

1  (ns molly.server.serve
2    (:use noir.core
3      molly.datatypes.entity
4      molly.search.lucene
5      molly.search.query-builder
6      molly.util.nlp
7      [molly.algo.bfs :only (bfs)]
8      [molly.algo.bfs-atom :only (bfs-atom)]
9      [molly.algo.bfs-ref :only (bfs-ref)])
10   (:require [noir.server :as server]
11             [noir.response :as response]))
12
13  (def runtime (Runtime/getRuntime))
14
15  (defn start!
16    [props]
17    (let [searcher (idx-searcher (idx-path (props :index)))]
18      (defn dox
19        [q1 field S op topk]
20        (let [bq (boolean-query
21                  (concat [[q1 :and]]
22                           (for [s S]
23                               [(query field s) op])))
24              result (map doc->data (idx-search searcher bq topk))
25              wat (fn [data] {:meta (meta data) :results data})]
26          (map wat result)))
27
28      (defn entities
29        [field q topk]
30        (dox (query :type :entity)
31              field
32              (clojure.string/split q #"\\s{1}")
33              :and
34              topk))
35
36      (defpage "/" []
37        (response/redirect "/index.html"))
38
39      (defpage "/value" {:keys [q]}
40        (response/json
41          {:result
42            (dox (query :type :value)
43                  :text
44                  (clojure.string/split (q-gram q) #"\\s{1}")
45                  :or
46                  (props :topk_value))}))
47
48      (defpage "/entities" {:keys [q]}
49        (response/json
50          {:result

```

Chapter 3

Concurrency

```
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   ;(find-group-for-id G v))
24   (remove #{v} (find-group-for-id G v)))
```

Figure 3.1: algo/common.clj

```

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

```

Figure 3.2: algo/bfs.clj


```

1 (ns molly.algo.bfs-atom
2   (use molly.algo.common))
3
4 (defn initial-state
5   [s]
6   (atom {:Q      (-> (clojure.lang.PersistentQueue/EMPTY) (conj s))
7             :marked #{}
8             :dist   {s 0}
9             :prev   {}
10            :done   false}))
11
12 (defn update-state
13   [state u v]
14   (let [Q      (state :Q)
15         marked (state :marked)
16         dist   (state :dist)
17         prev   (state :prev)]
18     (assoc state
19           :Q      (conj Q v)
20           :marked (conj marked v)
21           :dist   (assoc dist v (inc (dist u)))
22           :prev   (assoc prev v u))))
23
24 (defn deref-future
25   [dfd]
26   (if (future? dfd)
27       (deref dfd)
28       dfd))
29
30 (defn update-adj
31   [state-ref G u]
32   (let [marked? (@state-ref :marked)
33         deferred (doall
34                     (for [v (find-adj G u)]
35                       (if (marked? v)
36                           nil
37                           (future
38                             (swap! state-ref update-state u v))))))]
39     (doall (map deref-future deferred))))
40
41 (defn bfs-atom
42   [G s t]
43   (let [state-ref (initial-state s)]
44     (while (and (not (empty? (@state-ref :Q)))
45                (not (@state-ref :done)))
46       (let [u      (first (@state-ref :Q))
47             Q'     (pop (@state-ref :Q))
48             (swap! state-ref assoc :Q Q')
49             (if (some (fn [node] (= node t)) (@state-ref :marked))
50                 (swap! state-ref assoc :done true)

```

```

1  (ns molly.algo.bfs-ref
2    (use molly.algo.common))
3
4  (defn initial-state
5    [s]
6    (ref {:Q      (-> (clojure.lang.PersistentQueue/EMPTY) (conj s))
7           :marked #{s}
8           :dist   {s 0}
9           :prev   {}
10          :done   false}))
11
12  (defn update-state
13    [state u v]
14    (let [Q      (state :Q)
15          marked (state :marked)
16          dist   (state :dist)
17          prev   (state :prev)]
18      (assoc state
19             :Q      (conj Q v)
20             :marked (conj marked v)
21             :dist   (assoc dist v (inc (dist u)))
22             :prev   (assoc prev v u))))
23
24  (defn deref-future
25    [dfd]
26    (if (future? dfd)
27        (deref dfd)
28        dfd))
29
30  (defn update-adj
31    [state-ref G u]
32    (let [marked? (@state-ref :marked)
33          deferred (doall
34                     (for [v (find-adj G u)]
35                       (if (marked? v)
36                           nil
37                           (future (dosync (alter
38                                       state-ref
39                                       update-state
40                                       u
41                                       v))))))]
42      (doall (map deref-future deferred))))
43
44  (defn bfs-ref
45    [G s t]
46    (let [state-ref (initial-state s)]
47      (while (and (not (empty? (@state-ref :Q)))
48                  (not (@state-ref :done)))
49        (let [u (first (@state-ref :Q))
50              o' (pop (@state-ref :Q))]

```

Chapter 4

Literature Review

Chapter 5

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 5.1 provides an overview of how the data is represented in the system. It also defines notation to represent this data. Section 5.2 outlines how the relational database is related to the data in the system. Section 5.3 provides an example data corpus which is utilized throughout this thesis. It details the "mycampus" dataset. Finally, Section 5.4 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.

5.1 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 1. 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 2. 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 5.1: 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 3. 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 5.2: The structure of an entity group

5.2 Relational Database Abstraction

5.3 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 5.1). 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 5.1: Courses entity schema

The **Instructors** entity represents an individual instructor. Each instructor has a unique identifier, as well as a name (See Table 5.2). 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 5.2: 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 5.3).

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 5.3: 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 5.4).

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 5.4: 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 5.5).

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

Table 5.5: Teaches entity schema

5.4 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 6

Implementation

6.1 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)

6.1.1 Functional Data Structures

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

6.2 Tuneable Parameters

Tuneable parameters

Chapter 7

System Implementation

7.1 Technology Selection

7.1.1 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.

7.1.2 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.

7.1.3 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).

7.1.4 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.

7.2 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 (1)

7.2.1 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.

7.3 Implementation Issues

Chapter 8

Performance & Evaluation

8.1 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

8.2 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.

8.2.1 Pitfalls of the Java Virtual Machine

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 8.1: 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...

8.2.2 Mitigating JVM Pitfalls

JIT Compilation

```
33         ["-h" "--help" "Show help"]
34         :default false
35         :flag true]))
36
37 (def counter (atom 0))
38
39 (defn ns-to-ms
40   [nanosec]
41   (/ nanosec (* 1000.0 1000.0)))
```

Figure 8.1: core.clj

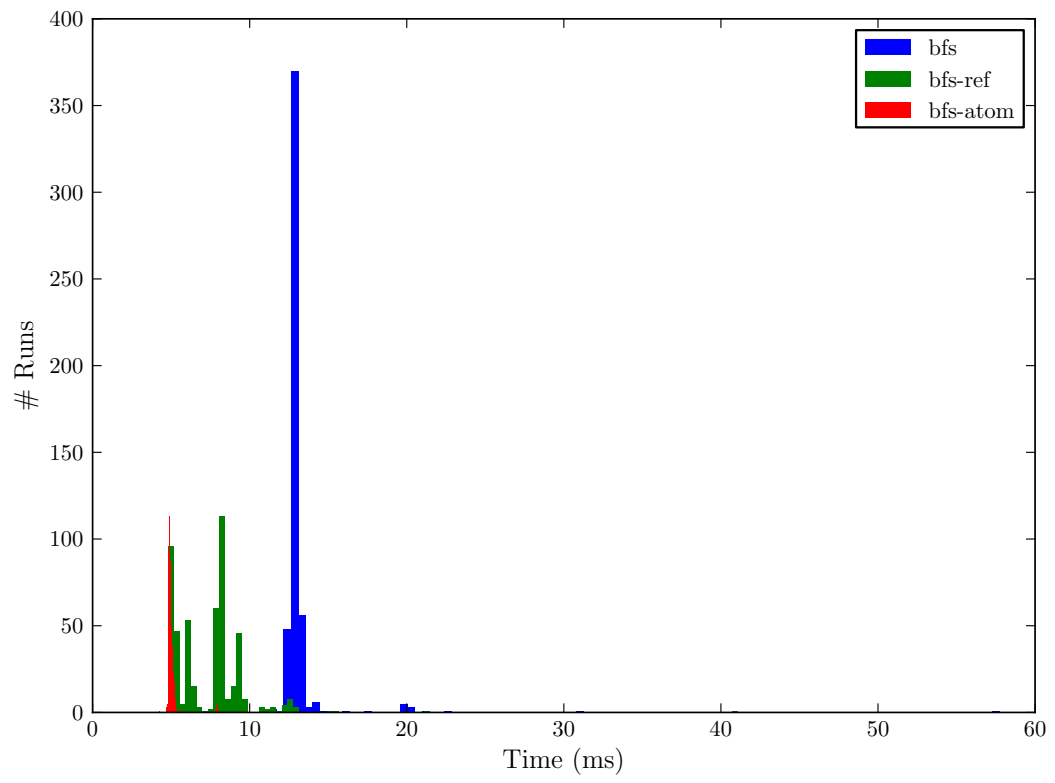


Figure 8.2: 1 Hop

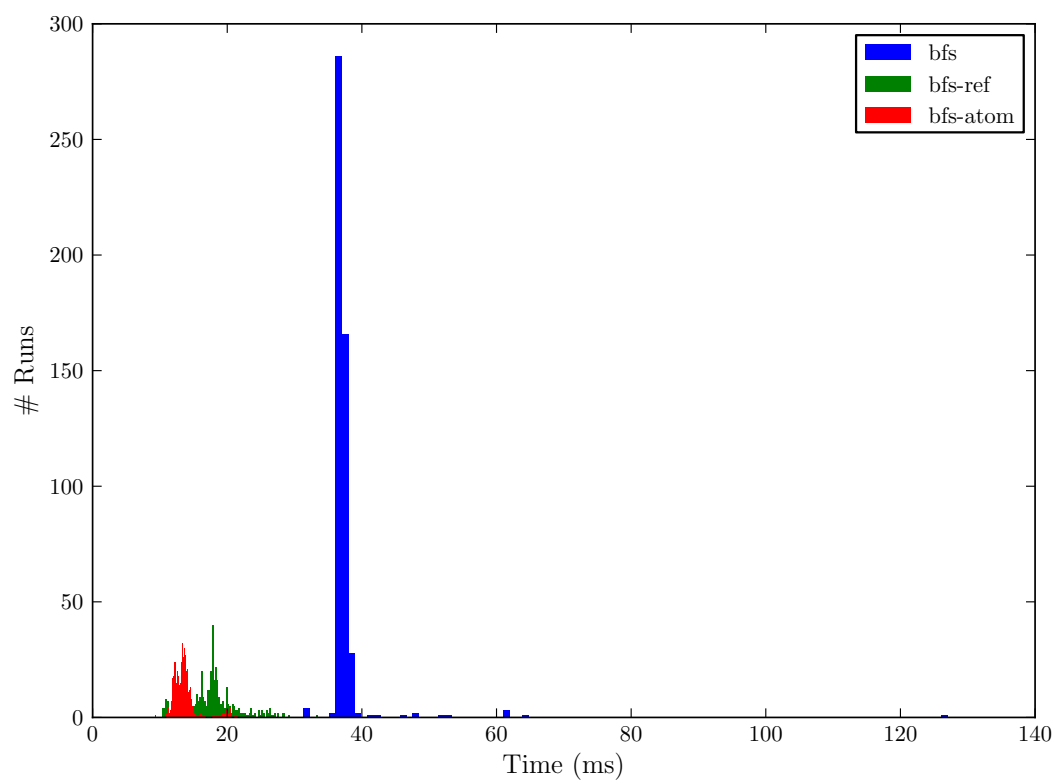


Figure 8.3: 2 Hops

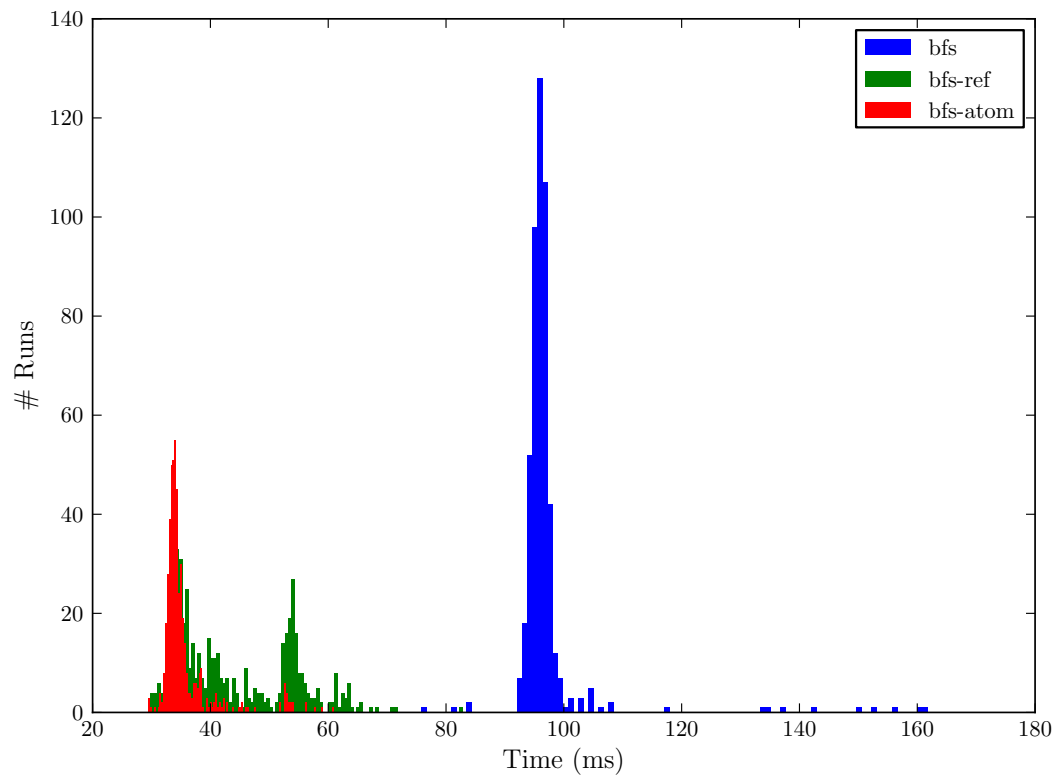


Figure 8.4: 3 Hops

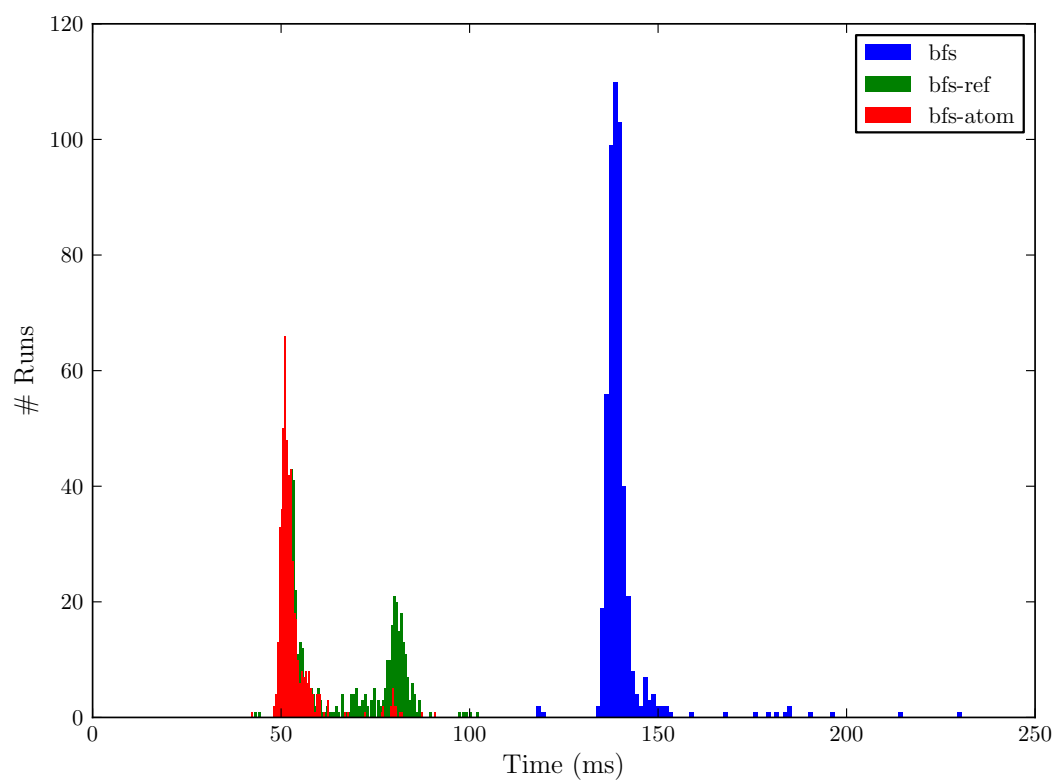


Figure 8.5: 4 Hops

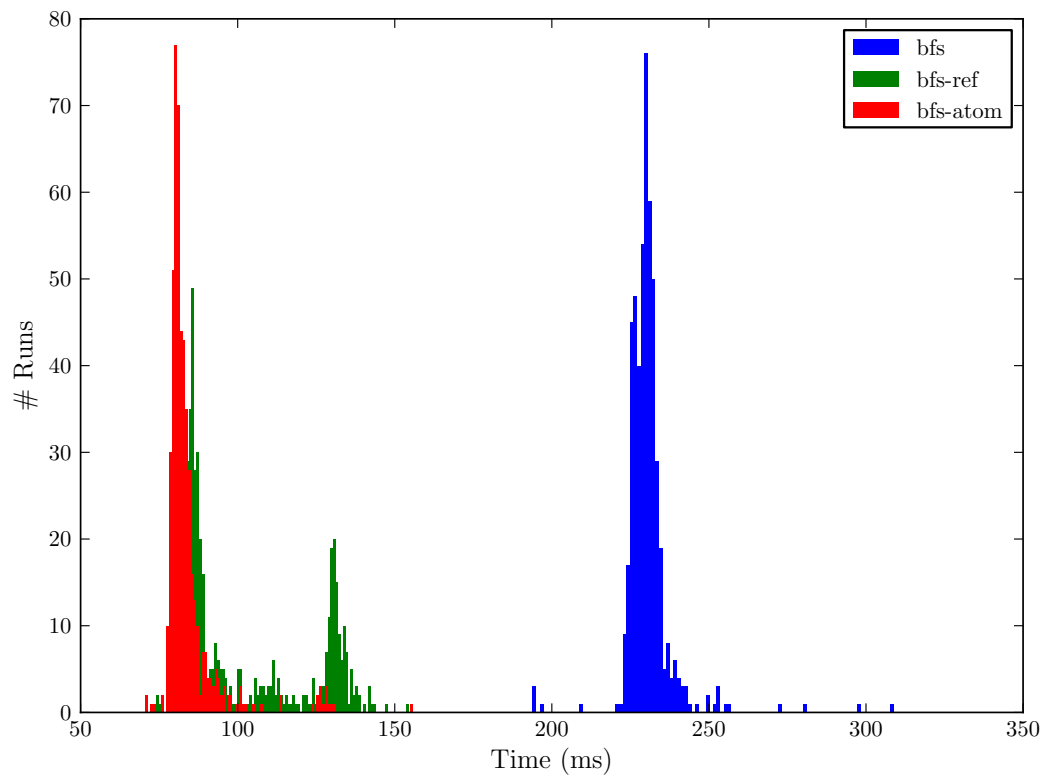


Figure 8.6: 5 Hops

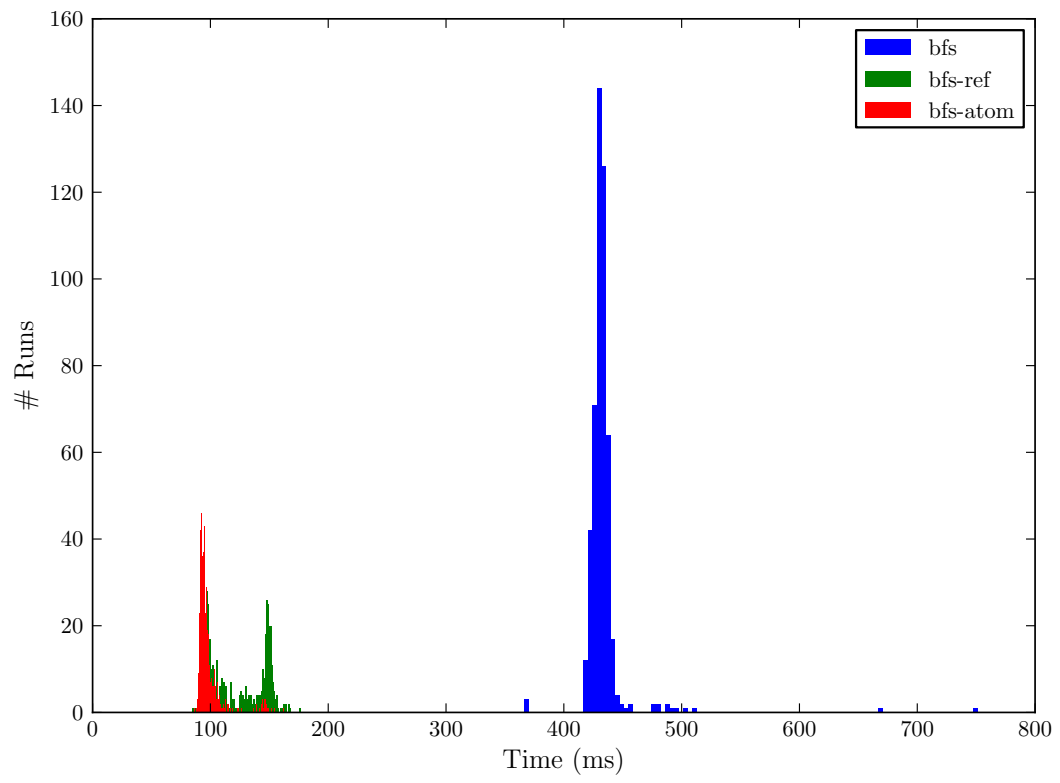


Figure 8.7: 6 Hops

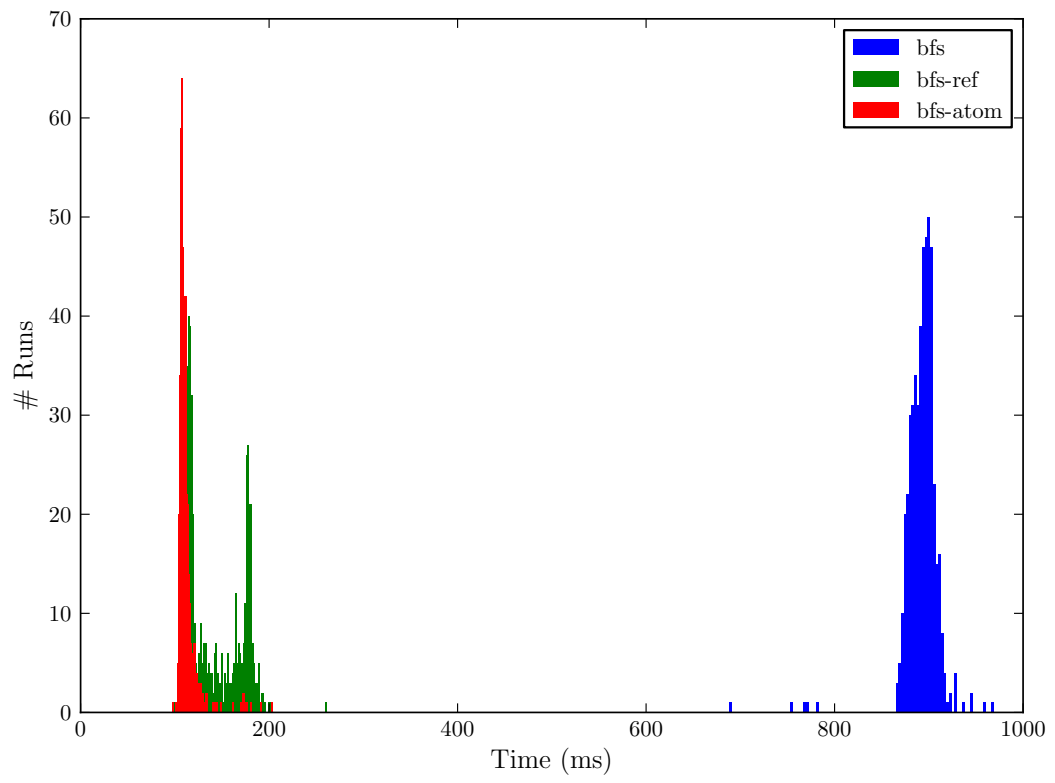


Figure 8.8: 7 Hops

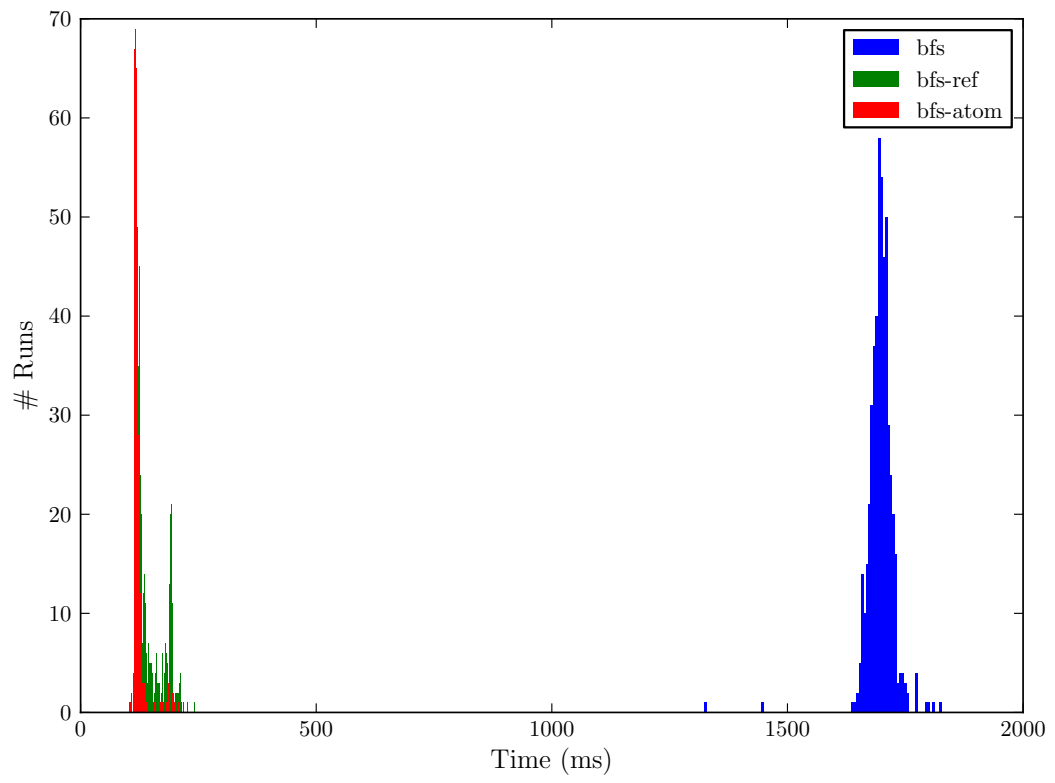


Figure 8.9: 8 Hops

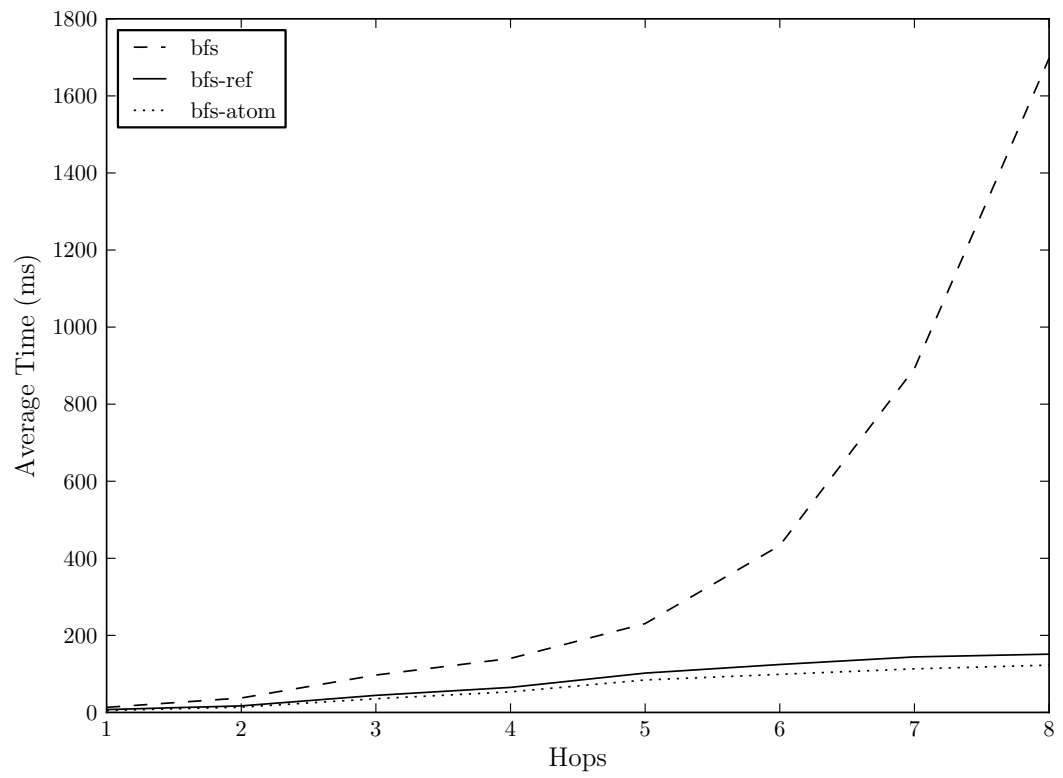


Figure 8.10: Comparison between single threaded and concurrent graph search

Chapter 9

Conclusion

[1]

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. 29): Add overview diagram of system design.