

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 (M.Sc.)

in

The Faculty of Science

Computer Science

University of Ontario Institute of Technology

Supervisor: Dr. Ken Q. Pu

December 2013

© Richard J.I. Drake, 2013

Contents

1	Background (2 days)	1
2	A Tale of Two Data Models	2
2.1	Relational Model	2
2.1.1	Schema Group	4
2.1.2	Entity Group	6
2.1.3	Pros and Cons of the Relational Model	7
2.2	Document Model	11
2.2.1	Vectorization of Documents	12
2.2.2	Extending the Document Model	16
2.2.3	Approximate String matching	16
2.2.4	Pros and Cons of the Document Model	18
3	Best of Both Worlds	20
3.1	Encoding Named Tuples into Documents	20
3.2	Mapping of Entity Groups to Documents	21
3.3	Encoding an Entity Group as a Document Group	21
3.4	Encoding Attribute Values into Searchable Documents	23
3.5	Iterative Search Using Document Encodings	23
4	Along Came Clojure	24

4.1	Basic Principles of Functional Programming	24
4.1.1	Features of Clojure	25
4.2	Search w/ Clojure	27
4.2.1	Thirdparty libraries (1 day, week 4)	27
4.2.2	Indexing of relational objects (5 days, week 5)	28
4.2.3	Keyword Search in document space (5 days, week 6)	28
4.2.4	Graph Search in document space (5 days, week 7)	28
5	Experimental Evaluation	29
5.1	Implementation	29
5.2	The data set	29
5.3	Runtime Evaluation	29
5.3.1	Methodology	30
5.4	Lessons learned	32
6	Conclusion (0 days)	34
A	Source Code	35
A.1	molly	35
A.1.1	molly.core	35
A.2	molly.conf	38
A.2.1	molly.conf.config	38
A.2.2	molly.conf.mycampus	39
A.3	molly.datatypes	43
A.3.1	molly.datatypes.database	43
A.3.2	molly.datatypes.entity	44
A.3.3	molly.datatypes.schema	47
A.4	molly.index	48
A.4.1	molly.index.build	48

A.5	molly.util	49
A.5.1	molly.util.nlp	49
A.6	molly.search	50
A.6.1	molly.search.lucene	50
A.6.2	molly.search.query_builder	52
A.7	molly.server	53
A.7.1	molly.server.serve	53
A.8	molly.algo	54
A.8.1	molly.algo.common	54
A.8.2	molly.algo.bfs	56
A.8.3	molly.algo.bfs_atom	57
A.8.4	molly.algo.bfs_ref	58
A.9	molly.bench	59
A.9.1	molly.bench.benchmark	59

List of Tables

2.1	Course relation	3
2.2	Results of the query in Fig. 2.4 on page 6.	5
2.3	Properties of the Course titled Human-Mutant Relations.	7
2.4	Properties of the Course titled Human-Mutant Relations.	16
2.5	Course document for MATH 360.	17
3.1	Doc[<i>t</i>]	21
3.2	Document encoding of Fig. 3.1 on page 22	22
4.1	Comparison between Clojure's four systems for concurrency	26
4.2	Syntactic sugar for Java virtual machine (JVM) interoperability	27

List of Figures

2.1	Subset of mycampus dataset schema	5
2.2	foreign key (FK) constraints on schema in Fig. 2.1 on page 5	5
2.3	Graph representation of relations (Fig. 2.1) and FK (Fig. 2.2)	6
2.4	Query to find section CRNs for a subject name.	6
2.5	Human-Mutant Relations entity group	7
2.6	Comparison between n -grams of G and G'	18
3.1	Example entity group	22
4.1	Representation of how data structures are “changed” in Clojure (Source: [Hic09]) . .	26
5.1	Growth of graph search times based on number of hops, plotted separately	32
5.2	Growth of graph search times based on number of hops, combined plot	33

List of Algorithms

1	N-GRAM(S, n, s)	17
---	-------------------------------	----

Acronyms

CSV comma-separated values. 31

FK foreign key. vii, 3–6, 8

JDBC Java database connectivity. 27

JSON JavaScript object notation. 31

JVM Java virtual machine. vi, 27

RDBMS relational database management system. 8, 10

SQL structured query language. 4, 5, 18

STM software transactional memory. 25, 26

TF-IDF term frequency and inverse document frequency. 12

List of Symbols

- schema graph** (G) graph representation of schema. 4, 6, 21
- entity group** (T) forest of **named tuples**. 6, 21, 22
- document collection** (C) set of **documents**. 11–16, 22
- terms** (T) set of unique **terms** in a **document collection**. 11, 13, 16
- document** (d) set of fields. x, 11–16, 18–20
- search query** (q) special case of **document**. 14–16, 18, 19, 23
- field** (f) named sub-document in . 20, 23
- term** (τ) unique term in **document collection**. 11–14, 18
- N number of documents in collection. 11
- database** (D) set of **relations**. 4, 6, 21
- relation** (r) set of named tuples. 2–4, 6, 21
- named tuple** (t) ordered set of values. vi, 2–4, 6, 7, 20–22
- attribute** (α) named column. 3, 4, 6, 7, 20
- key** (K) uniquely identifies a **named tuple** in a **relation**. 3

Chapter 1

Background (2 days)

Literature search on:

- DBExplore
- XRank
- BANKS
- ...

Chapter 2

A Tale of Two Data Models

The term “data model” refers to a notation for describing data and/or information. It consists of the data structure, operations that may be performed on the data, as well as constraints placed on the data [GUW09].

In this chapter we provide a formal definition of the relational data model, discuss its merits, its shortcomings, and contrast it to the document data model. Contrary to the relational model, the document model permits fast and flexible keyword search without requiring explicit domain knowledge of the data. In addition, we demonstrate the feasibility of encoding a relational model into a document model in a lossless manner.

2.1 Relational Model

In its most basic form, the relational data model is built upon sets and tuples. Each of these sets consist of a set of finite possible values. Tuples are constructed from these sets to form relations.

Definition 1 (Named Tuple). A named tuple t is an instance of a relation r , consisting of values corresponding to the attributes of r . For example,

Example 1. Given a tuple $t = \{\text{code} : \text{“CDPS 101”}, \text{title} : \text{“Human-Mutant Relations”}, \text{subject} : \text{“CDPS”}\}$, we denote the attributes of t as $\text{ATTR}[t] = \{\text{code}, \text{title}, \text{subject}\}$. The values are $t[\text{code}] =$

“CDPS 101”, $t[\text{title}] = \text{“Human-Mutant Relations”}$, and $t[\text{subject}] = \text{“CDPS”}$.

Definition 2 (Relation). A relation r is a set of named tuples, $r = \{t_1, t_2, \dots, t_n\}$, such that all the named tuples share the same attributes.

$$\forall t, t' \in r, \text{ATTR}[t] = \text{ATTR}[t'] \quad (2.1)$$

Example 2. An example Course relation, r , would be

$$r = \left\{ \begin{array}{lll} \{\text{code} : \text{“CDPS 101”}, & \text{title} : \text{“Human-Mutant Relations”}, & \text{subject} : \text{“CDPS”}\}, \\ \{\text{code} : \text{“CDPS 201”}, & \text{title} : \text{“Humans and You”}, & \text{subject} : \text{“CDPS”}\}, \\ \{\text{code} : \text{“MATH 360”}, & \text{title} : \text{“Complex Analysis”}, & \text{subject} : \text{“MATH”}\} \end{array} \right\}$$

Relations are typically represented as tables.

code	title	subject
CDPS 101	Human-Mutant Relations	CDPS
CDPS 201	Humans and You	CDPS
MATH 360	Complex Analysis	MATH

Table 2.1: Course relation

Definition 3 (Keys). Keys are constraints imposed on relations. A key constraint K on a relation r is a subset of $\text{ATTR}[r]$ which may uniquely identify a tuple. Formally, we say r satisfies the key constraint K , denoted as $r \models K$, subject to

$$\forall t, t' \in r, t \neq t' \implies t[K] \neq t'[K]$$

For example, in Table 2.1, the relation satisfies the key constraint $\{\text{code}\}$ or $\{\text{title}\}$, but not $\{\text{subject}\}$.

Definition 4 (Foreign Keys). A FK constraint applies to two relations, r_1, r_2 . It asserts that values of certain attributes of r_1 must appear as values of some corresponding attributes of r_2 . A FK constraint is written as

$$\theta = r_1(\alpha_{11}, \alpha_{12}, \dots, \alpha_{1k}) \rightarrow r_2(\alpha_{21}, \alpha_{22}, \dots, \alpha_{2k})$$

where $\alpha_{1i} \subseteq \text{ATTR}[r_1]$ and $\alpha_{2i} \subseteq \text{ATTR}[r_2]$. We say (r_1, r_2) satisfies θ , denoted as $(r_1, r_2) \models \theta$, if for every tuple $t \in r_1$, there exists a tuple $t' \in r_2$ such that $t[\alpha_{11}, \alpha_{12}, \dots, \alpha_{1k}] = t'[\alpha_{21}, \alpha_{22}, \dots, \alpha_{2k}]$.

We say r_1 is the source, while r_2 is the target.

Example 3. Suppose we have a relation $\text{Course}(\text{code}, \text{title}, \text{subject})$. We impose a **FK** constraint of

$$\theta = \text{Course}(\text{subject}) \rightarrow \text{Subject}(\text{id}) \quad (2.2)$$

which asserts $(\text{Course}, \text{Subject}) \models \theta$. Therefore, if

$$t = \{\text{code} : \text{"CDPS 101"}, \text{title} : \text{"Human-Mutant Relations"}, \text{subject} : \text{"CDPS"}\}$$

then $\exists! t' \in \text{Subject}$ such that $t'[\text{id}] = \text{"CDPS"}$.

Definition 5 (Relational Database). A relational database, D , is a named collection of relations (as defined by Definition 2 on the preceding page), keys (as defined by Definition 3 on the previous page), and foreign key constraints (as defined by Definition 4 on the preceding page).

We use $\text{NAME}[D]$ to denote the name of D , $\text{REL}[D]$ the list of relations in D , $\text{KEY}[D]$ the list of key constraints of D , and $\text{FK}[D]$ the list of foreign key constraints of D .

2.1.1 Schema Group

Definition 6 (Schema Graph). If we view relations as vertices, and foreign key constraints as edges, a database D can be viewed as a *schema graph* G , formally defined as

$$\text{vertices} : V(G) = \text{REL}[D] \quad (2.3)$$

$$\text{edges} : E(G) = \text{FK}[D] \quad (2.4)$$

Example 4. Given the schema in Fig. 2.1 on the next page and the **FK** constraints in Fig. 2.2 on the following page we produce the schema graph in Fig. 2.3 on page 6

The relational data model is particularly powerful for analytic queries. Given the schema graph in Fig. 2.3 on page 6, one can formulate the following analytic queries in a query language known as **structured query language (SQL)**.

Subject(id, name)
 Course(code, title, subject)
 Term(id, name)
 Section(crn, term, course)
 Schedule(id, days, sch_type, time_start, time_end, location, section, instructor)
 Instructor(id, name)

Figure 2.1: Subset of mycampus dataset schema

Course(subject) \rightarrow Subject(id)
 Section(term) \rightarrow Term(id)
 Section(course) \rightarrow Course(code)
 Schedule(section) \rightarrow Section(crn)
 Schedule(instructor) \rightarrow Instructor(id)

Figure 2.2: FK constraints on schema in Fig. 2.1

Example 5. Using **SQL**, find all section CRNs for the subject titled “Community Development & Policy Studies.”

The **SQL** query in Fig. 2.4 on the next page results in Table 2.2.

crn
10000
10001
10002

Table 2.2: Results of the query in Fig. 2.4 on the following page.

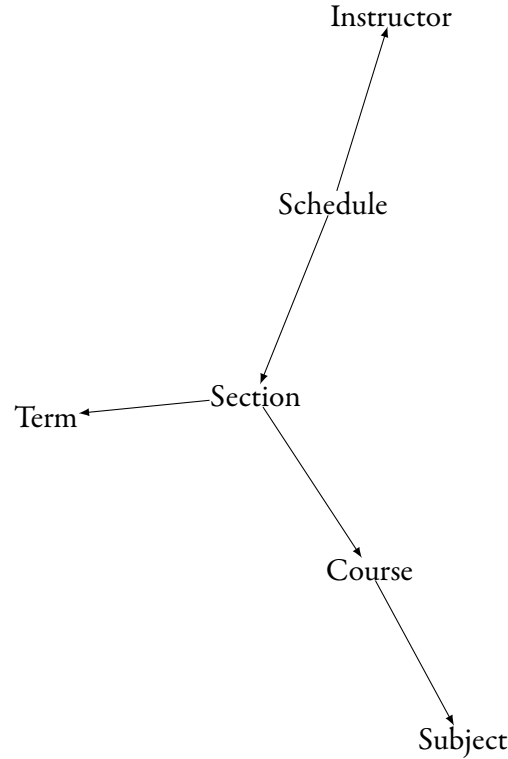


Figure 2.3: Graph representation of relations (Fig. 2.1) and FK (Fig. 2.2)

```

1 SELECT section.crn
2 FROM   section
3       JOIN course
4         ON section.course_code = course.code
5       JOIN subject
6         ON subject.id = course.subject_id
7 WHERE  subject.name = 'Community Development & Policy Studies';

```

Figure 2.4: Query to find section CRNs for a subject name.

2.1.2 Entity Group

Definition 7 (Entity Group). An entity group is a forest, T , of tuples interconnected by join conditions defined by the FK constraints in the schema graph G .

Given two vertices $t, t' \in V(T)$, $\exists r_1, r_2 \in \text{REL}[D]$ such that $t \in r_1$, $t' \in r_2$, and $(r_1, r_2) \in G$.

That is, t and t' belong to two relations that are connected by the schema graph.

Let $r_1(\alpha_{11}, \alpha_{12}, \dots, \alpha_{1k}) \rightarrow r_2(\alpha_{21}, \alpha_{22}, \dots, \alpha_{2k})$ be the FK that connects r_1, r_2 . We further assert

that $t[\alpha_{11}, \alpha_{12}, \dots, \alpha_{1k}] = t'[\alpha_{21}, \alpha_{22}, \dots, \alpha_{2k}]$.

Entity groups define complex, structured objects that include more information than individual tuples in the relations.

Example 6. The information in Table 2.4 on page 16 all relates to the Course titled Human-Mutant Relations, however no single tuple in the database has all of this information as a result of database normalization.

Attribute	Value
code	CDPS 101
title	Human-Mutant Relations
subject	Community Development & Policy Studies

Table 2.3: Properties of the Course titled Human-Mutant Relations.

We require an entity group (Fig. 2.5) to join together all pieces of information related to this course.

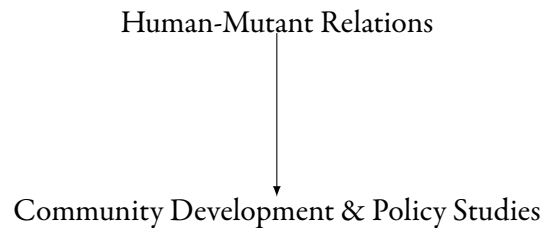


Figure 2.5: Human-Mutant Relations entity group

2.1.3 Pros and Cons of the Relational Model

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

Pros

The enforcement of constraints is essential to the relational model. There are several types of constraints, including uniqueness and **FKs**. The first constraint maintains uniqueness.

The Course relation (Table 2.1 on page 3) has the attribute `code` as its primary key. In order for other relations to reference a specific named tuple, the `code` attribute must be unique.

Example 7 (Unique Constraint). Attempt to insert another course with a `code` of “CDPS 101.”

```
1 INSERT INTO course
2 VALUES      ('CDPS 101',
3              'Mutant-Human Relations',
4              'CDPS');
```

The **relational database management system (RDBMS)** enforces the primary key constraint on the `code` attribute, rejecting the insertion.

Error: column code is not unique

With the uniqueness of named tuples guaranteed (as demonstrated in Example 7), we must ensure that any named tuples that are referenced actually exist. If they do not, the database must not permit the operation to continue. Doing so would lead to dangling references.

Example 8 (Referential Integrity). Attempt to insert the tuple (“CHEM 101”, “Introductory Chemistry”, “CHEM”) in the Course relation.

```
1 INSERT INTO course
2 VALUES      ('CHEM 101',
3              'Introductory Chemistry',
4              'CHEM');
```

Again we see the **RDBMS** protecting the integrity of the data.

Error: foreign key constraint failed

In addition to enforcing consistency, the relational model is capable of providing higher-level views of the data through aggregation.

Example 9 (Aggregation). Find the number of sections offered for the subject named “Community Development & Policy Studies.”

```
1 SELECT Count(*)
2 FROM    section
3          JOIN course
4          ON section.course = course.code
5          JOIN subject
6          ON subject.id = course.subject
7 WHERE    subject.name = 'Community Development & Policy Studies';
```

Information stored within a properly designed database is normalized. That is, no information is repeated.

Example 10 (Normalization). For example, suppose Emma Frost became headmistress and the subject named “Community Development & Policy Studies” was renamed to “Community Destruction & Policy Studies.” If this information were not normalized, each course in this subject would need to be updated. Since this information is normalized, the following query will suffice.

```
1 UPDATE subject
2 SET    name = 'Community Destruction & Policy Studies'
3 WHERE  id = 'CDPS';
```

The above examples are some of the most important reasons for choosing the relational model over others. Unfortunately, the relational model is not without its downsides.

Cons

While the relational model excels at ensuring data consistency, aggregation, and reporting; it is not suitable for every task. In order to issue queries, a user must be familiar with the schema. This requires specific domain knowledge of the data.

An example of a complicated query involving two joins is give in Fig. 2.4 on page 6.

A casual user is unlikely to determine the correct join path, name of the tables, name of the attributes, etc. This is in contrast to the document model, where the data is semi-structured or unstructured, requiring minimal domain knowledge.

The relational model is also rigid in structure. If a relation is modified, every query referencing said relation may require a rewrite. Even a simple attribute being renamed (e.g. $\rho_{\text{name/alias}}(\text{Person})$) is capable of modifying the join paths. This rigidity places additional cognitive burden on users.

In addition to having a rigid structure, most relational database management systems lack flexible string matching options. Assuming basic SQL-92 compliance, a **RDBMS** only supports the **LIKE** predicate [ISO11].

Example 11 (LIKE Predicate). Find all courses with a title that contains “man.”

```
1 SELECT *  
2 FROM   course  
3 WHERE  title LIKE '%man%';
```

There are a couple of limitations to the **LIKE** predicate. First, it only supports basic substring matching. If a user accidentally searches for all courses with a title containing “men,” nothing would be found.

Second, unless the predicate is applied to the end of the string and the column is indexed, performance will be poor. The database must scan the entire table in order to answer the query, resulting in performance of $\mathcal{O}(n)$, where n is the number of named tuples in the relation.

2.2 Document Model

In contrast to the relational model, the document model represents semi-structured as well as unstructured data. Examples of information suitable to the document model includes emails, memos, book chapters, etc.

These pieces, or units, of information are broken into documents. Groups of related documents (for example, a library catalogue) are referred to as a document collection.

Definition 8 (Terms and Document). A term, τ , is an indivisible string (e.g. a proper noun, word, or a phrase). A document, d , is a bag of words; order is irrelevant.

Let $\text{freq}(\tau, d)$ be the frequency of term τ in d , T denote all possible terms, and $\text{BAG}[T]$ be all possible bag of terms.

Remark 1. We use the bag-of-words model for documents. This means that position information of terms in a document is irrelevant, but the frequency of terms are kept in the document. Documents are non-distinct sets.

Definition 9 (Document Collection). A document collection C is a set of documents, written $C = \{d_1, d_2, \dots, d_k\}$. The cardinality of C is denoted by N .

Example 12. Consider the following short phrases

1. math 360 is a math class
2. cdps 101 is a boring lecture
3. mathematics lecture was great

Each sentence phrase produces a document, giving us the following

$$d_1 = \{\text{"math"} : 2, \text{"a"} : 1, \text{"is"} : 1, \text{"360"} : 1, \text{"class"} : 1\} \quad (2.5)$$

$$d_2 = \{\text{"a"} : 1, \text{"boring"} : 1, \text{"is"} : 1, \text{"cdps"} : 1, \text{"lecture"} : 1, \text{"101"} : 1\} \quad (2.6)$$

$$d_3 = \{\text{"mathematics"} : 1, \text{"great"} : 1, \text{"was"} : 1, \text{"lecture"} : 1\} \quad (2.7)$$

2.2.1 Vectorization of Documents

One of the most fundamental approaches for searching documents is to treat documents as high-dimensional vectors, and the document collection as a subset in a vector space. Search queries become a nearest neighbour search in a vector space using a distance metric.

The first step is to convert a bag of terms into vectors. The standard technique [MRS08] uses a scoring function that measures the relative importance of terms in documents.

Definition 10 (**Term frequency and inverse document frequency (TF-IDF) Score**). The term frequency is the number of times a term τ appears in a document d , as given by $\text{freq}(\tau, d)$. The document frequency of a term τ , denoted by $\text{df}(\tau)$, is the number of documents in C that contains τ . It is defined as

$$\text{df}(\tau) = |\{d \in C : \tau \in d\}|$$

The combined **TF-IDF** score of τ in a document d is given by

$$\text{tf-idf}(C, \tau, d) = \frac{\text{freq}(\tau, d)}{|d|} \cdot \log \frac{N}{\text{df}(\tau)}$$

The first component, $\frac{\text{freq}(\tau, d)}{|d|}$, measures the importance of a term within a document. It is normalized to account for document length. The second component, $\log \frac{N}{\text{df}(\tau)}$, is a measure of the rarity of the term within the document collection C .

Example 13. Using the documents from Example 12 on the preceding page, the **TF-IDF** scores are

as follows.

	d_1	d_2	d_3
τ_1 : “101”	0.0000	0.2642	0.0000
τ_2 : “360”	0.3170	0.0000	0.0000
τ_3 : “a”	0.1170	0.0975	0.0000
τ_4 : “boring”	0.0000	0.2642	0.0000
τ_5 : “cdps”	0.0000	0.2642	0.0000
τ_6 : “class”	0.3170	0.0000	0.0000
τ_7 : “great”	0.0000	0.0000	0.3962
τ_8 : “is”	0.1170	0.0975	0.0000
τ_9 : “lecture”	0.0000	0.0975	0.1462
τ_{10} : “math”	0.6340	0.0000	0.0000
τ_{11} : “mathematics”	0.0000	0.0000	0.3962
τ_{12} : “was”	0.0000	0.0000	0.3962

Definition 11 (Document Vector). Given a document collection \mathcal{C} with M unique terms $\mathbf{T} = [\tau_1, \tau_2, \dots, \tau_n]$, each document d can be represented by an M -dimensional vector.

$$\vec{d} = \begin{bmatrix} \text{tf-idf}(\tau_1, d) \\ \text{tf-idf}(\tau_2, d) \\ \vdots \\ \text{tf-idf}(\tau_n, d) \end{bmatrix}$$

Example 14. The documents in Example 12 on page 11 would produce the following vectors.

$$\vec{d}_n = \begin{bmatrix} \text{tf-idf}(\tau_1, d_n) \\ \text{tf-idf}(\tau_2, d_n) \\ \text{tf-idf}(\tau_3, d_n) \\ \text{tf-idf}(\tau_4, d_n) \\ \text{tf-idf}(\tau_5, d_n) \\ \text{tf-idf}(\tau_6, d_n) \\ \text{tf-idf}(\tau_7, d_n) \\ \text{tf-idf}(\tau_8, d_n) \\ \text{tf-idf}(\tau_9, d_n) \\ \text{tf-idf}(\tau_{10}, d_n) \\ \text{tf-idf}(\tau_{11}, d_n) \\ \text{tf-idf}(\tau_{12}, d_n) \end{bmatrix}, \vec{d}_1 = \begin{bmatrix} 0.0000 \\ 0.3170 \\ 0.1170 \\ 0.0000 \\ 0.0000 \\ 0.3170 \\ 0.0000 \\ 0.1170 \\ 0.0000 \\ 0.6340 \\ 0.0000 \\ 0.0000 \end{bmatrix}, \vec{d}_2 = \begin{bmatrix} 0.2642 \\ 0.0000 \\ 0.0975 \\ 0.2642 \\ 0.2642 \\ 0.0000 \\ 0.0000 \\ 0.0975 \\ 0.0975 \\ 0.0000 \\ 0.0000 \\ 0.0000 \end{bmatrix}, \vec{d}_3 = \begin{bmatrix} 0.0000 \\ 0.0000 \\ 0.0000 \\ 0.0000 \\ 0.0000 \\ 0.0000 \\ 0.3962 \\ 0.0000 \\ 0.1462 \\ 0.0000 \\ 0.3962 \\ 0.3962 \end{bmatrix}$$

Definition 12 (Search Query). A search query q is simply a document (as defined by Definition 8 on page 11). The top- k answers to q with respect to a collection C is defined as the k documents, $\{d_1, d_2, \dots, d_k\}$ in C , such that $\{\vec{d}_1, \vec{d}_2, \dots, \vec{d}_k\}$ are the closest vectors to \vec{q} using a Euclidean distance measure in \mathbb{R}^N .

Example 15. Given the search query $q = \{\text{math}, \text{lecture}, \text{was}, \text{great}\}$, compute the vector \vec{q} within

the document collection \mathcal{C} (as defined in Example 12 on page 11).

$$\vec{q} = \begin{bmatrix} 0.0000 \\ 0.0000 \\ 0.0000 \\ 0.0000 \\ 0.0000 \\ 0.0000 \\ 0.2500 \\ 0.1038 \\ 0.1038 \\ 0.2500 \\ 0.0000 \\ 0.0000 \end{bmatrix}$$

In order to determine the top- k documents for search query q , we need a way of measuring the similarity between documents.

Definition 13 (Cosine Similarity). Given two document vectors, \vec{d}_1 and \vec{d}_2 , the cosine similarity is the dot product $\vec{d}_1 \cdot \vec{d}_2$, normalized by the product of the Euclidean distance of \vec{d}_1 and \vec{d}_2 in \mathbb{R}^N . It is denoted as $\text{similarity}(\vec{d}_1, \vec{d}_2)$.

$$\text{similarity}(\vec{d}_1, \vec{d}_2) = \frac{\vec{d}_1 \cdot \vec{d}_2}{\|\vec{d}_1\| \cdot \|\vec{d}_2\|} \quad (2.8)$$

$$= \frac{\sum_{i=1}^N \vec{d}_{1,i} \times \vec{d}_{2,i}}{\sqrt{\sum_{i=1}^N (\vec{d}_{1,i})^2} \times \sqrt{\sum_{i=1}^N (\vec{d}_{2,i})^2}} \quad (2.9)$$

Recall we may represent search queries as documents and thus document vectors. Therefore we may compute the score of a document d for a search query q as

$$\text{similarity}(\vec{d}, \vec{q})$$

Example 16. Given the document collection \mathcal{C} (from Example 12 on page 11) and search query q , compute the similarity between q and every document $d \in \mathcal{C}$.

$$\text{similarity}(\vec{d}_1, \vec{q}) = 0.390890 \quad (2.10)$$

$$\text{similarity}(\vec{d}_2, \vec{q}) = 0.061592 \quad (2.11)$$

$$\text{similarity}(\vec{d}_3, \vec{q}) = 0.252789 \quad (2.12)$$

2.2.2 Extending the Document Model

In the extended document model, documents have fields, denoted as $\text{FIELD}[d]$, and each field has a value. Thus

$$d : \text{FIELD}[d] \rightarrow \text{BAG}[\mathbf{T}]$$

Example 17 (Semi-Structured Document). We see that d_1 is about MATH 360. The document contents are semi-structured, containing both a course code and the subject ID. By adding fields to the document, we are left with Table 2.5 on the next page.

Attribute	Value
code	CDPS 101
title	Human-Mutant Relations
subject	Community Development & Policy Studies

Table 2.4: Properties of the Course titled Human-Mutant Relations.

which is similar in structure to Table 2.4.

2.2.3 Approximate String matching

Definition 14 (N-Gram). An n -gram is a contiguous sequence of substrings of string S of length n . An algorithm for computing the n -gram of S is given in Algorithm 1 on the next page.

Field	Value
code	MATH 360
subject	MATH
body	math 360 is a math class

Table 2.5: Course document for MATH 360.

Algorithm 1 N-GRAM(S, n, s)**Require:** S is a string, $n \geq 1$, and s is a character**Ensure:** the list of n -grams of S

```

1:  $G \leftarrow []$ 
2:  $p \leftarrow \text{REPEAT}(s, n - 1)$ 
3:  $S \leftarrow \text{PAD}(S, p)$ 
4:  $S \leftarrow \text{REPLACE}(S, ' ', p)$ 
5: for  $i = 0$  to  $l - n + 1$  do
6:   append  $S[i, i + n]$  to  $G$ 
7: end for
8: return  $G$ 

```

Where l is the length of S , $\text{REPEAT}(S, n)$ repeats s character n times, $\text{PAD}(S, p)$ prefixes and postfixes S with p , and $\text{REPLACE}(S, s, p)$ replaces character s with p in string S .

Example 18. Given a string $S = \text{"human"}$, compute the trigram of S using Algorithm 1.

$$G = \{ \text{"$$h"}, \text{"$hu"}, \text{"hum"}, \text{"uma"}, \text{"man"}, \text{"an$"}, \text{"n$$"} \}$$

We use n -grams in order to permit approximate string matching.

Example 19. Given a string S (Example 18), let $S' = \text{"humans"}$. Compute the trigram of S' and compare it to S .

$$G' = \{ \text{"$$h"}, \text{"$hu"}, \text{"hum"}, \text{"uma"}, \text{"man"}, \text{"ans"}, \text{"ns$"}, \text{"s$$"} \}$$

Comparing G to G' results in the following matrix

As Fig. 2.6 on the following page shows, using n -grams yield a similarity of $\frac{5}{10}$.

	G	G'
$\tau_1 : \text{"ns\$"}$	0	1
$\tau_2 : \text{"n\$\$"}$	1	0
$\tau_3 : \text{"s\$\$"}$	0	1
$\tau_4 : \text{"ans"}$	0	1
$\tau_5 : \text{"man"}$	1	1
$\tau_6 : \text{"uma"}$	1	1
$\tau_7 : \text{"\$\$h"}$	1	1
$\tau_8 : \text{"hum"}$	1	1
$\tau_9 : \text{"\$hu"}$	1	1
$\tau_{10} : \text{"an\$"}$	1	0

Figure 2.6: Comparison between n -grams of G and G' .

2.2.4 Pros and Cons of the Document Model

There are numerous reasons to use the document model. The most significant reason is that it allows users without domain knowledge and working knowledge of a complex query language such as **SQL** to find information.

Example 20 (Simple Queries). Find all documents related to “mathematics” or “lecture”. The result of the query q would be

$$\text{query}(\text{"mathematics"}) \cup \text{query}(\text{"lecture"}) \rightarrow \{d_2, d_3\}$$

Users can also modify queries to require certain terms be present or not present.

Example 21 (AND Query). Find all documents containing both “mathematics” and “lecture”. This query would return the following set of documents

$$\text{query}(\text{"mathematics"}) \cap \text{query}(\text{"lecture"}) \rightarrow \{d_3\}$$

as only d_3 contains both terms.

Example 22 (NOT Query). Find all documents containing “mathematics” but not “lecture”. This query would return different results than Example 21 on the previous page.

$$\text{query}(\text{“mathematics”}) \neg \text{query}(\text{“lecture”}) \rightarrow \emptyset$$

While none of the above queries required domain knowledge, it is possible to use the extended document model (Section 2.2.2 on page 16) to search specific fields. Doing so permits users to leverage their existing domain knowledge in order to achieve finer control over what documents are retrieved.

Example 23 (Extended Query). Find all documents with a subject of “MATH” that contain the term “class”.

$$\text{query}(\text{“subject”, “MATH”}) \cap \text{query}(\text{“class”}) \rightarrow \{d_1\}$$

Not only does the document model provide a familiar interface to search for information with, it also ranks the results. In the relational model a search for “mathematics” would return all named tuples that contained that term. In the document model, documents are ranked against the query q and the top- k documents are returned.

The advantage is that users have the result of q already ranked so only the most relevant documents may be explored. As the number of documents matching q for a large corpus can be high, showing only the top- k relevant documents may save the user a substantial amount of time.

The relational model does not permit approximate string matching. By utilizing the document model with n -grams (Section 2.2.3 on page 16), users who substitute, delete, or insert characters from the desired term may still receive results for their intended term (see Example 19 on page 17 for a demonstration of how n -grams overcome character insertion).

Unfortunately the document model does not support the concept of foreign keys (Definition 4 on page 3). While information is easily accessible due to flexible search, each document is a discrete unit of information. Aggregate queries are unsupported, as these units are not linked amongst one another.

Chapter 3

Best of Both Worlds

3.1 Encoding Named Tuples into Documents

Recall in the extended document model (Section 2.2.2 on page 16), a document d consists of fields f_1, f_2, \dots, f_n . Using the extended document model, we are left with a straight forward mapping of a tuple t to document d .

For tuple t , every attribute $\alpha \in \text{ATTR}[t]$ maps to a field f in document d . Every attribute value must be analyzed into an indexable form in order to store it in a field.

$$\text{ATTR}[t] \xrightarrow{\text{analyzed}} \text{FIELD}[d] \quad (3.1)$$

$$\alpha_1, \alpha_2, \dots, \alpha_n \xrightarrow{\text{analyzed}} f_1, f_2, \dots, f_n \quad (3.2)$$

We denote the document encoding of t as $\text{DOC}[t]$.

Example 24. Given the tuple

$$t = \{\text{code} : \text{“CDPS 101”}, \text{title} : \text{“Human-Mutant Relations”}, \text{subject} : \text{“CDPS”}\}$$

produce the document encoding $\text{DOC}[t]$.

Field	Terms
code	{cdps, 101}
title	{human, mutant, relations}
subject	{cdps}

Table 3.1: Doc[t]

3.2 Mapping of Entity Groups to Documents

Recall that an entity group (Definition 7 on page 6) is a forest T of tuples t such that for every $(t, t') \in T$, where $t \neq t'$, implies $\text{REL}[t] \neq \text{REL}[t']$. That is, every distinct tuple is from a distinct relation.

Given the restriction

$$\forall (r, r') \in G, \exists! (r, r') \models \theta$$

we assert that if t and t' are in the entity group T , then there is a foreign key constraint between t and t' . We denote the vertices of T as $V(T)$, and the edges of T as $E(T)$.

Example 25. Using the schema graph...

Claim 1. Given $V(T)$, we are always able to reconstruct T .

Proof. Given $V(T)$, we must reconstruct $E(T)$ in order to complete T .

Choose any $(t, t') \in V(T)$. If $(\text{REL}[t], \text{REL}[t']) \in GD$, then (t, t') is an edge in T .

Recall our earlier assertion that GD is cycle-free and foreign keys must be unique. □

To do (1)

3.3 Encoding an Entity Group as a Document Group

Given a entity group T , we construct two or more documents in order to represent the entity group in the document model.

For every $t \in V(T)$, we construct a document $\text{Doc}[t]$ (Section 3.1 on page 20). With each tuple t stored in the document collection C , we construct an additional document which stores the association information.

Let x be the indexing document of T .

$$x[\text{“entities”}] = \bigcup_{t \in V(T)} \text{UID}[t] \quad (3.3)$$

Thus, the encoding of T is defined as

$$T \xrightarrow{\text{encode}} \{\text{Doc}[t] : t \in V(T)\} \cup \{x\} \quad (3.4)$$

Example 26. An entity group produced from the schema in Fig. 2.3 on page 6 would be as follows

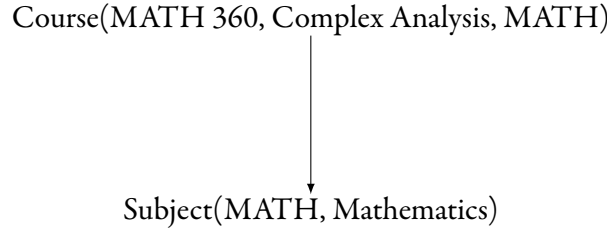


Figure 3.1: Example entity group

Transforming the example entity group in Fig. 3.1 would produce documents shown in Table 3.2

Field	Terms	Field	Terms	Field	Terms
code	{math, 360}	id	{math}	entities	{course math_360,
title	{complex, analysis}	name	{mathematics}		subject math}
subject	{math}				

(a) Course
(b) Subject
(c) Indexing document

Table 3.2: Document encoding of Fig. 3.1

It’s easy to see that from $\text{encode}(T)$ we can recover $V(T)$, the tuples in T .

By Claim 1 on the preceding page, this is sufficient to recover T entirely.

3.4 Encoding Attribute Values into Searchable Documents

Each value for user selected attributes are converted into n -grams, and stored in special documents.

3.5 Iterative Search Using Document Encodings

A document database supports fast and flexible keyword search queries. A search query is characterized by $q = (f, w)$, where f is an optional field name, and w is a search phrase.

$\text{query}(q)$ is the set of documents returned by the text index. The query function, combined with the extended document model, permits powerful search queries to be issued. Our implementation supports approximate string matching using n -grams (Section 2.2.3) for values, searching for entities containing keywords (see Example 27), and the discovery of intermediate entities given two known entities.

Example 27 (Entity Search). Find all entities that match the keyword “math”.

Let $q = \text{“math”}$ be the search query. The results are

$$\begin{aligned}\text{query}(q) &= \text{query}(\text{“math”}) \\ &= \{\text{subject|math}, \text{course|math_360}\}\end{aligned}$$

which are coincidentally related. The results of an entity search query are not necessarily related.

Example 28 (Entity Graph Search). Find the shortest path between the two entities with the unique identities of “subject|math” and “instructor|5”.

Chapter 4

Along Came Clojure

Talk about how great Clojure is.

4.1 Basic Principles of Functional Programming

The functional programming paradigm follows a handful of basic tenets; values are immutable, and functions must be free of side-effects [[Hug89](#)].

The first tenet, that values are immutable, refers to the fact that once a value is bound, this value may not change. In procedural programming there is the concept of assignment, whereas in functional programming, a value is bound. Assignment allows a value to change, binding does not.

Immutable values are advantageous as they remove a common source of bugs; state must explicitly be changed. This removes the ability for different areas of a program to modify the state (i.e. global variables).

Unfortunately immutable values can also lead to inefficiency. For example, in order to add a key-value pair to a map, an entirely new map must be created with the existing key-value pairs copied to it. In practice this is avoided through the use of persistent data structures with multi-versioning.

The second tenet, that functions must be free of side-effects, meaning the output of a function must be predictable for any given input. This purity reduces a large source of bugs, and allows out-of-order execution. [[Hug89](#)].

4.1.1 Features of Clojure

The creator of Clojure, Rich Hickey, describes his language as follows:

Clojure is a dialect of Lisp, and shares with Lisp the code-as-data philosophy and a powerful macro system. Clojure is predominantly a functional programming language, and features a rich set of immutable, persistent data structures. When mutable state is needed, Clojure offers a software transactional memory system and reactive Agent system that ensure clean, correct, multithreaded designs. ([Hica])

As the above quote describes, Clojure follows the basic tenets of functional programming.

Immutable, Persistent Data Structures

Clojure supports a rich set of data structures. These are immutable, satisfying the first tenet, as well as persistent, in order to overcome the inefficiency described in the first tenet.

The provided data structures range from scalars (numbers, strings, characters, keywords, symbols), to collections (lists, vectors, maps, array maps, sets) [Hicb]. These data structures are sufficient enough to allow us to use the Universal Design Pattern [Yeg08].

Clojure also has the concept of persistent data structures. These are used in order to avoid the inefficiency of creating a new data structure and copying over the contents of the old data structure simply to make a change. Clojure creates a skeleton of the existing data structure, inserts the value into the data structure, then retains a pointer to the old data structure. If an old property is accessed on the new data structure, Clojure follows the pointers until the property is found on a previous data structure. See Fig. 4.1 on the following page for an example.

To do (2)

Concurrency

Clojure supports four systems for concurrency: **software transactional memory (STM)**, agents, atoms, and dynamic vars. The differences between these systems are summarized in Table 4.1 on the next page.

To do (3)

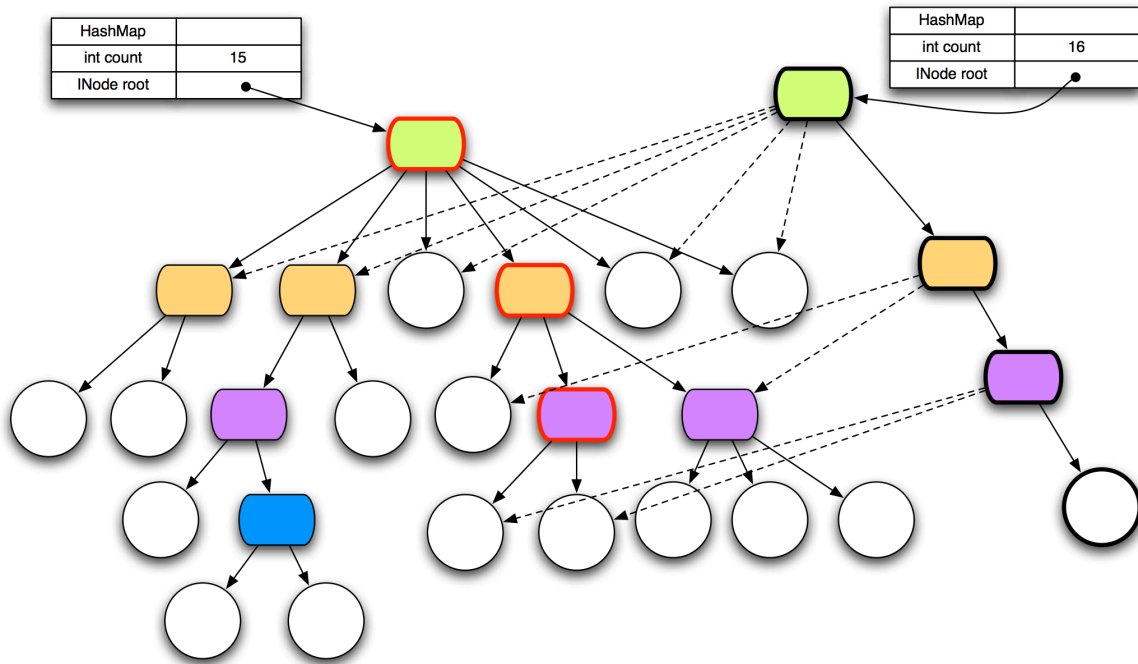


Figure 4.1: Representation of how data structures are “changed” in Clojure (Source: [Hic09])

System Name	Synchronous	Coordinated	Scope
STM	Yes	Yes	Application
Agents	No	No	Application
Atoms	Yes	No	Application
Dynamic Vars	Not Applicable	Not Applicable	Thread

Table 4.1: Comparison between Clojure’s four systems for concurrency

Operation	Description	Example
Member Access		

Table 4.2: Syntactic sugar for JVM interoperability

Interoperability With the JVM

Traditionally functional programming languages have been undesirable for numerous reasons, many of which are detailed in the classic paper titled “Why No One Uses Functional Languages”. Clojure attempts to avoid many of these reasons by running on the JVM. The JVM allows Clojure to both call and be called by Java and other languages. It includes syntactic sugar to transparently call Java code, as well as make itself available to Java. This avoids the issues of compatibility, libraries, portability, availability, packagability, and tools.

The syntactic sugar provided by Clojure allows for the accessing of object members, the creation of objects, the calling of methods on an instance or class, etc. Clojure also includes shortcuts to perform multiple operations on the same object. The syntax is given in Table 4.2.

We utilize Clojure’s JVM interoperability to make use of Apache Lucene and Java database connectivity (JDBC).

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

4.2 Search w/ Clojure

4.2.1 Thirdparty libraries (1 day, week 4)

- Lucene

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

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

4.2.3 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

4.2.4 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

5.1 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

Scripts were written to coordinate the execution, collection, and transformation of the performance data of our implementation.

5.3.1 Methodology

We used Criterium¹ to handle the execution of the benchmarks as it handles unique concerns stemming from benchmarking on the JVM. These include:

- Statistical processing of multiple evaluations
- Inclusion of a warm-up period, designed to allow the JIT compiler to optimize its code
- Purging of the garbage collector before testing, to isolate timings from GC state prior to testing
- A final forced GC after testing to estimate impact of cleanup on the timing results

Unfortunately this requires a much longer runtime as each function must be invoked numerous times. In extreme cases (Ford-Fulkerson, 8 hops) this can take upwards of 4 hours in our test environment.

Data Collection

Criterium provides us with a Clojure map with performance data. It performs analysis, presenting us with outliers, samples, etc. As this data collection process can take several hours or more, this data is collected and stored for offline analysis.

In order to utilize the Clojure output in Python, a data interchange format (JSON) is used. The benchmark function writes the Criterium performance analysis out as a JSON string to stdout and Python captures the output, JSONifies it, and stores it in an array. This array is written to disk in JSON as well so it can be loaded into the data transformation script.

For example:

```
[{"results": {...}, "method": "bfs", "max-hops": 1}, ...]
```

¹<http://hugoduncan.org/criterium/>

Data Processing

Several scientific computing libraries are used in the processing and visualization.

There are two forms the data takes:

- comma-separated values (CSV)
- JavaScript object notation (JSON)

The CSV data is generated from the JSON data which is generated as described in Section 5.3.1 on the previous page.

With the data loaded, we're interested in a handful of pieces of data per each entry.

- Max hops
- Method
- Mean execution time

We can easily load and parse the JSON data.

- Index speed
- Keyword search speed
- Graph search speed:
 - Ford Fulkerson
 - BFS
 - Concurrent BFS using refs
 - Concurrent BFS using atoms

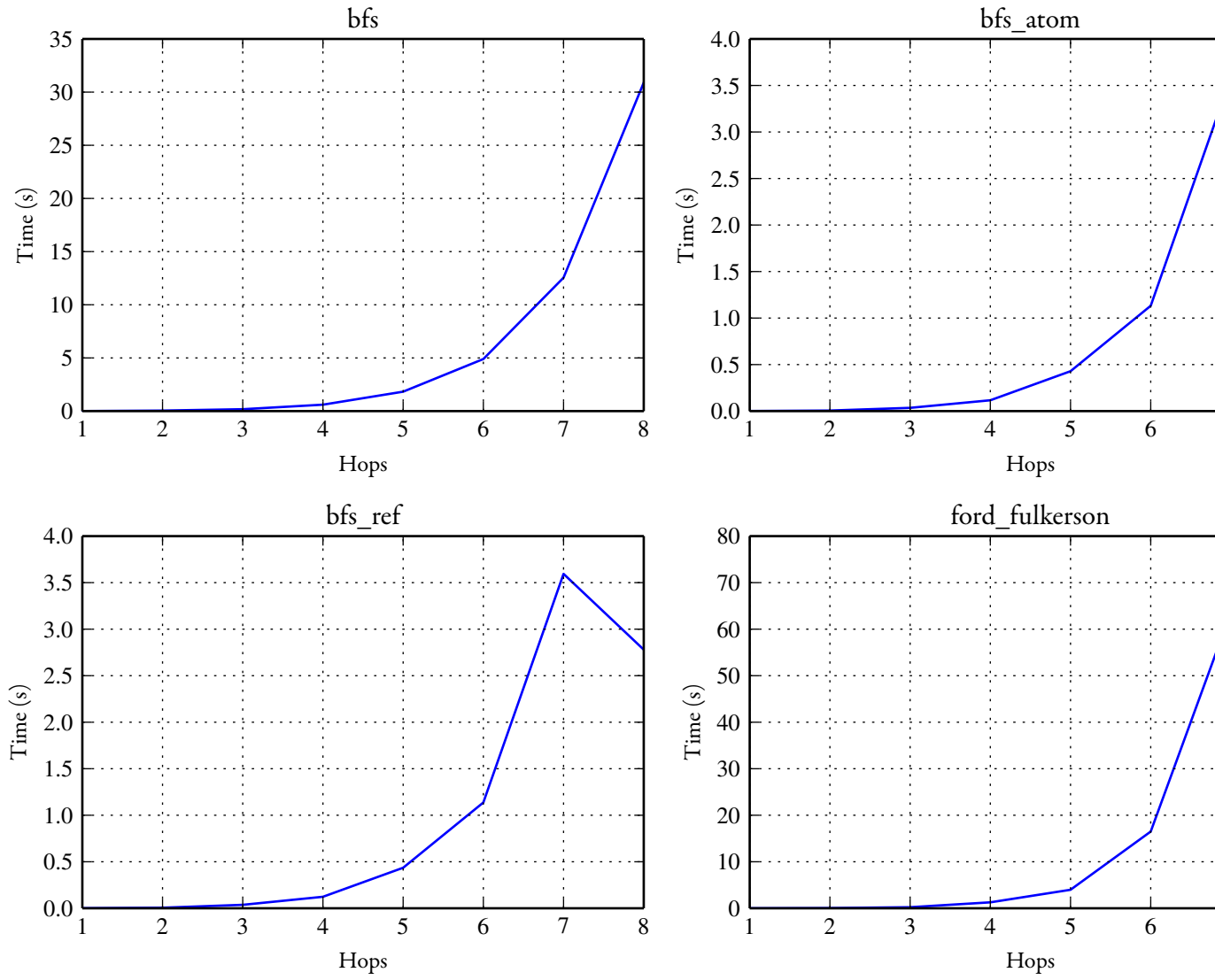


Figure 5.1: Growth of graph search times based on number of hops, plotted separately

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.

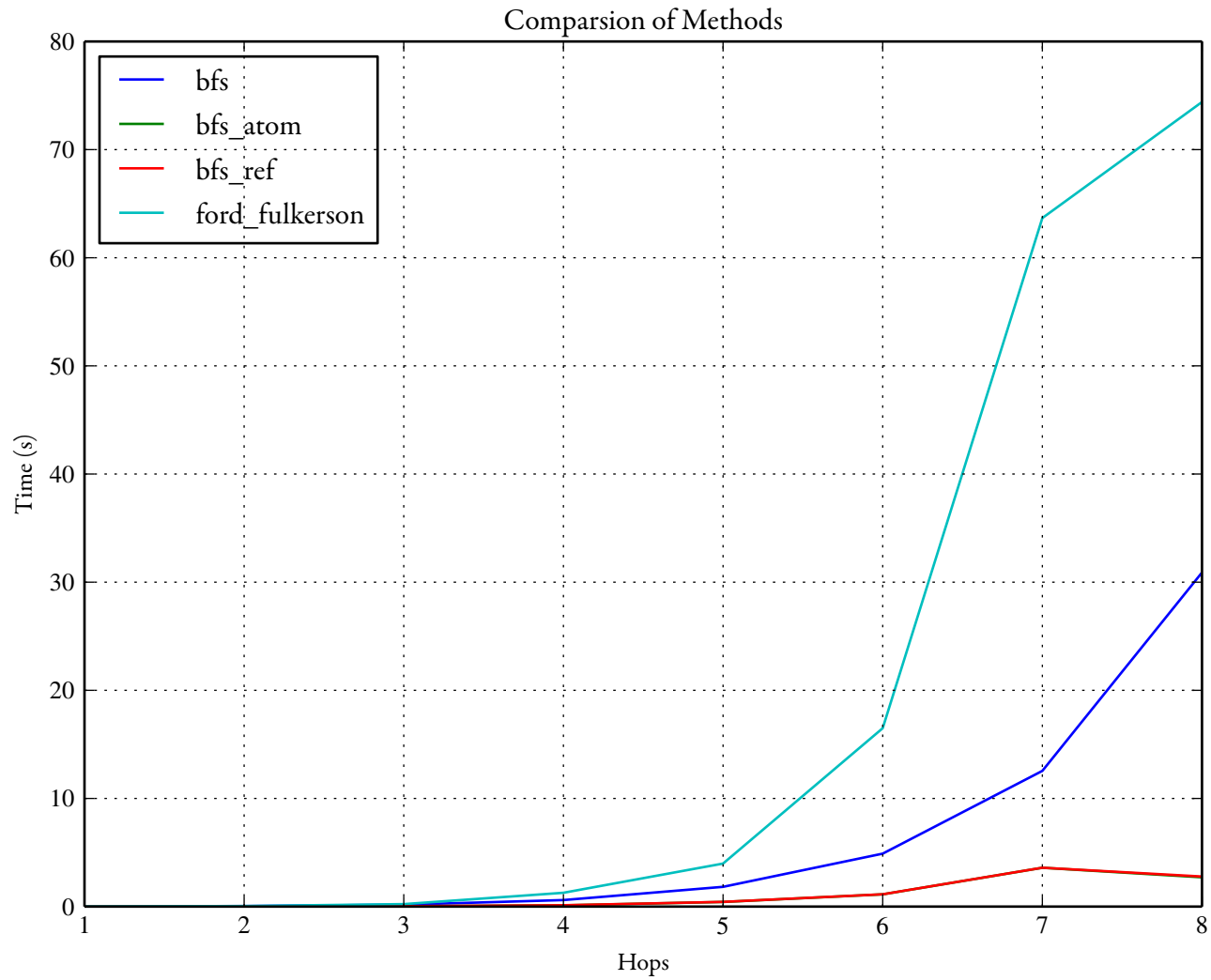


Figure 5.2: Growth of graph search times based on number of hops, combined plot

- The clojure way: correctness first, runtime optimization latter (ref to atom is natural).

Chapter 6

Conclusion (0 days)

Survived Clojure.

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      :parse-fn #(Integer. %)]
22     ["--index" "Build an index of the database"
23      :default false]
24     [:flag true]
25     ["--benchmark" "Run benchmarks"]
```

```

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

```

```
71         (println dist)
72         (println prev))
73     (benchmark-search f searcher source target max-hops))
74     (shutdown-agents))
75     nil))))
```

A.2 molly.conf

A.2.1 molly.conf.config

```
1 (ns molly.conf.config
2   (:use propertea.core))
3
4 (defn load-props
5   ([ ]
6     (load-props "config/molly.properties"))
7   ([file-name]
8     (read-properties file-name
9                      :parse-int [:idx.topk.value
10                                :idx.topk.entities
11                                :idx.topk.entity
12                                :idx.search.max-hops]
13                      :required [:db.path
14                                :idx.path])))
15
16 (defprotocol IConfig
17   (connection [this])
18   (schema [this])
19   (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  [[:__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$Occur
4                                         BooleanQuery
5                                         PhraseQuery)))
6
7 (defn query
8   [kind & args]
9   (let [field-name (condp = kind
10                      :type    "__type__"
11                      :class    "__class__"
12                      :id       "__id__"
13                      :text     "__all__"
14                      ; Assume "kind" is an attribute name.
15                      (condp = (type kind)
16                            clojure.lang.Keyword (name kind)
17                            java.lang.String      kind))
18       phrase-query (PhraseQuery.)]
19     (doseq [arg args]
20       (. phrase-query add (Term. field-name (name arg))))
21
22     phrase-query))
23
24 (defn boolean-query
25   [args]
26   (let [query (BooleanQuery.)]
27     (doseq [[q op] args]
28       (. query add q (condp = op
29                          :and BooleanClause$Occur/MUST
30                          :or  BooleanClause$Occur/SHOULD
31                          :not BooleanClause$Occur/MUST_NOT)))
32
33     query))

```

A.7 molly.server

A.7.1 molly.server.serve

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

- [GUW09] Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. *Database systems - the complete book (2. ed.)*. Pearson Education, 2009.
- [Hica] Rich Hickey. Clojure. URL: <http://clojure.org/>.
- [Hicb] Rich Hickey. Data structures. URL: http://clojure.org/data_structures.
- [Hic09] Rich Hickey. Persistent data structures and managed references. London, United Kingdom: QCon London, March 2009. URL: http://qconlondon.com/dl/qcon-london-2009/slides/RichHickey_PersistentDataStructuresAndManagedReferences.pdf.
- [Hug89] J. Hughes. Why Functional Programming Matters. *Computer journal*, 32(2):98–107, 1989.
- [ISO11] ISO. Information technology – database languages – sql – part 2: foundation (sql/foundation). ISO (ISO/IEC 9075-2:2011). Geneva, Switzerland: International Organization for Standardization, 2011.
- [MRS08] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to information retrieval*. Cambridge University Press, New York, NY, USA, 2008. ISBN: 0521865719, 9780521865715.
- [Wad98] Philip Wadler. Why no one uses functional languages. *Sigplan not.*, 33(8):23–27, August 1998. ISSN: 0362-1340. DOI: [10.1145/286385.286387](https://doi.org/10.1145/286385.286387). URL: <http://doi.acm.org/10.1145/286385.286387>.

- [Yeg08] Steven Yegge. The universal design pattern. October 2008. URL: <http://steveyegge.blogspot.ca/2008/10/universal-design-pattern.html>.

To do...

- ☐ 1 (p. 21): Above proof could use some love
- ☐ 2 (p. 25): Ensure this is sourced properly
- ☐ 3 (p. 25): More concurrency