

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	7
2.1.3	Pros and Cons of the Relational Model	8
2.2	Document Model	11
2.2.1	Vectorization of Documents	12
2.2.2	Extending the Document Model	17
2.2.3	Approximate String matching	17
2.2.4	Pros and Cons of the Document Model	19
3	Best of Both Worlds	21
3.1	Encoding Named Tuples into Documents	21
3.2	Mapping of Entity Groups to Documents	22
3.3	Encoding an Entity Group as a Document Group	22
3.4	Encoding Attribute Values into Searchable Documents	24
3.5	Iterative Search Using Document Encodings	24
4	Along Came Clojure	25

4.1	Basic Principles of Functional Programming	25
4.1.1	Features of Clojure	26
4.2	Search With Clojure	29
4.2.1	Full-Text Search Using Lucene	29
4.2.2	Indexing Relational Database	29
4.2.3	Indexing of relational objects (5 days, week 5)	31
4.2.4	Keyword Search in document space (5 days, week 6)	31
4.2.5	Graph Search in document space (5 days, week 7)	32
5	Experimental Evaluation	33
5.1	Implementation	33
5.2	The data set	33
5.3	Runtime Evaluation	33
5.3.1	Methodology	34
5.4	Lessons learned	38
6	Conclusion (0 days)	39
A	Source Code	40
A.1	molly	40
A.1.1	molly.core	40
A.2	molly.conf	43
A.2.1	molly.conf.config	43
A.2.2	molly.conf.mycampus	44
A.3	molly.datatypes	48
A.3.1	molly.datatypes.database	48
A.3.2	molly.datatypes.entity	49
A.3.3	molly.datatypes.schema	52
A.4	molly.index	53

A.4.1	molly.index.build	53
A.5	molly.util	54
A.5.1	molly.util.nlp	54
A.6	molly.search	55
A.6.1	molly.search.lucene	55
A.6.2	molly.search.query_builder	57
A.7	molly.server	58
A.7.1	molly.server.core	58
A.7.2	molly.server.remotes	59
A.7.3	molly.server.search	60
A.7.4	molly.server.util	62
A.8	molly.algo	63
A.8.1	molly.algo.common	63
A.8.2	molly.algo.bfs	65
A.8.3	molly.algo.bfs_atom	66
A.8.4	molly.algo.bfs_ref	67
A.9	molly.bench	68
A.9.1	molly.bench.benchmark	68

List of Tables

2.1	Course relation	3
2.2	Results of the query in Fig. 2.4 on page 7.	7
2.3	Properties of the Course titled Human-Mutant Relations.	8
2.4	Properties of the Course titled Human-Mutant Relations.	17
2.5	Course document for MATH 360.	18
3.1	Doc[<i>t</i>]	22
3.2	Document encoding of Fig. 3.1 on page 23	23
4.1	Comparison between Clojure's four systems for concurrency	27
4.2	Syntactic sugar for Java virtual machine (JVM) interoperability	28
4.3	Keys expected by EntitySchema records	30
4.4	Transformation from tuple to internal representation to document	31

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.	7
2.5	Human-Mutant Relations entity group	8
2.6	Comparison between n -grams of G and G'	19
3.1	Example entity group	23
4.1	Representation of how data structures are “changed” in Clojure (Source: [Hic09]) . .	27
4.2	Clojure code that, given a path, returns a Directory object	28
4.3	IConfig protocol all configurations must adhere to	29
5.1	Growth of graph search times based on number of hops, plotted separately	36
5.2	Growth of graph search times based on number of hops, combined plot	37

List of Algorithms

1	N-GRAM(S, n, s)	18
---	-------------------------------	----

Acronyms

CSV comma-separated values. 35

DSL domain-specific language. 30

FK foreign key. vii, 3–8

JDBC Java database connectivity. 28, 29

JSON JavaScript object notation. 35

JVM Java virtual machine. vi, 28, 29

RDBMS relational database management system. 9, 11

SQL structured query language. 4, 7, 19, 30

STM software transactional memory. 27

TF-IDF term frequency and inverse document frequency. 13

List of Symbols

- schema graph** (G) graph representation of schema. 4, 7, 22
- entity group** (T) forest of **named tuples**. 7, 22, 23
- document collection** (C) set of **documents**. x, 12–17, 23
- terms** (T) set of unique **terms** in a **document collection**. 12, 14, 17
- document** (d) set of fields. x, 12–17, 19–21
- search query** (q) special case of **document**. 15–17, 19, 20, 24
- field** (f) named sub-document in **document**. 21, 24
- term** (τ) unique term in **document collection**. x, 12–15, 19
- N number of documents in collection. 12
- database** (D) set of **relations**. 4, 7, 22
- relation** (r) set of named tuples. x, 2–4, 7, 22
- named tuple** (t) ordered set of values. vi, x, 2–4, 7, 21–23
- attribute** (α) named column. 3, 4, 7, 21
- 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 $t[\text{code}] =$

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

“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.

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

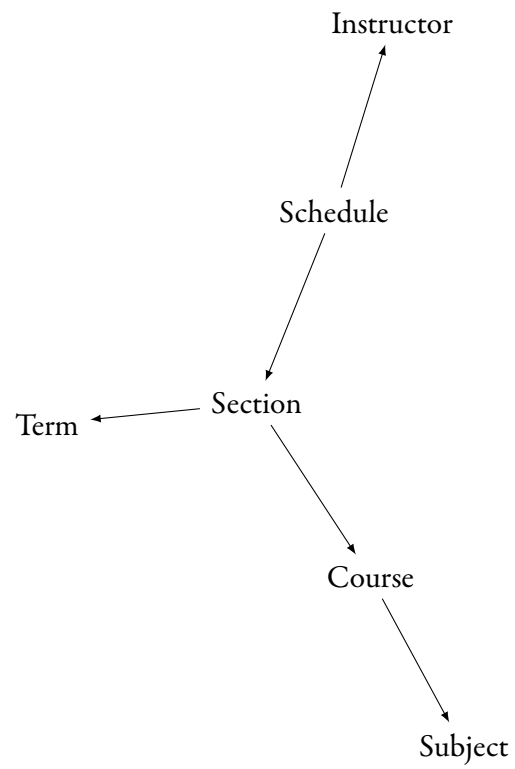


Figure 2.3: Graph representation of relations (Fig. 2.1) and FK (Fig. 2.2)

```

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.

crn
10000
10001
10002

Table 2.2: Results of the query in Fig. 2.4.

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

The **SQL** query in Fig. 2.4 results in Table 2.2.

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.

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.

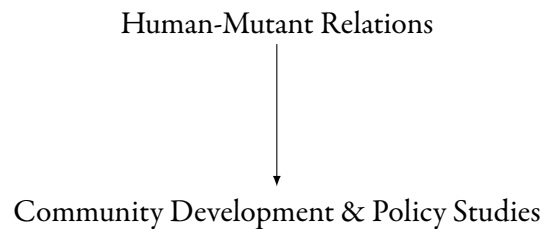


Figure 2.5: Human-Mutant Relations entity group

Example 6. The information in Table 2.4 on page 17 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.

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

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 7.

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 , \mathbf{T} denote all possible terms, and $\text{BAG}[\mathbf{T}]$ 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 \mathbf{C} is a set of documents, written $\mathbf{C} = \{d_1, d_2, \dots, d_k\}$. The cardinality of \mathbf{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 12 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 12). 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 $q = \{\text{math, lecture, was, great}\}$, compute the vector \vec{q} within

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

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

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

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

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.

Example 16. Given the document collection \mathcal{C} (from Example 12 on page 12) 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.

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 17) 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 17), users who substitute, delete, or insert characters from the desired term may still receive results for their intended term (see Example 19 on page 18 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 17), 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 7) 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.

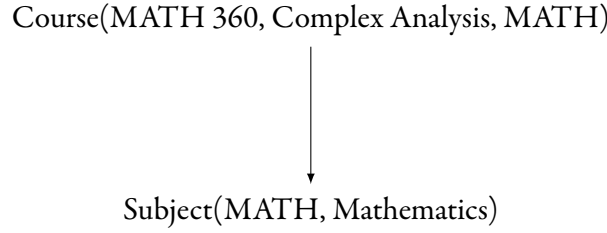


Figure 3.1: Example entity group

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

For every $t \in V(T)$, we construct a document $\text{Doc}[t]$ (Section 3.1 on page 21). 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

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

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, means that 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 previously.

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.

In Fig. 4.1 on the following page, we see what happens when a persistent data structure is “changed” in Clojure. The root of the left tree is the data structure before, and the root of the right tree is the data structure after. Note how the changed map retains pointers to all but the updated value; the newly created value is pointed to instead of the previous one.

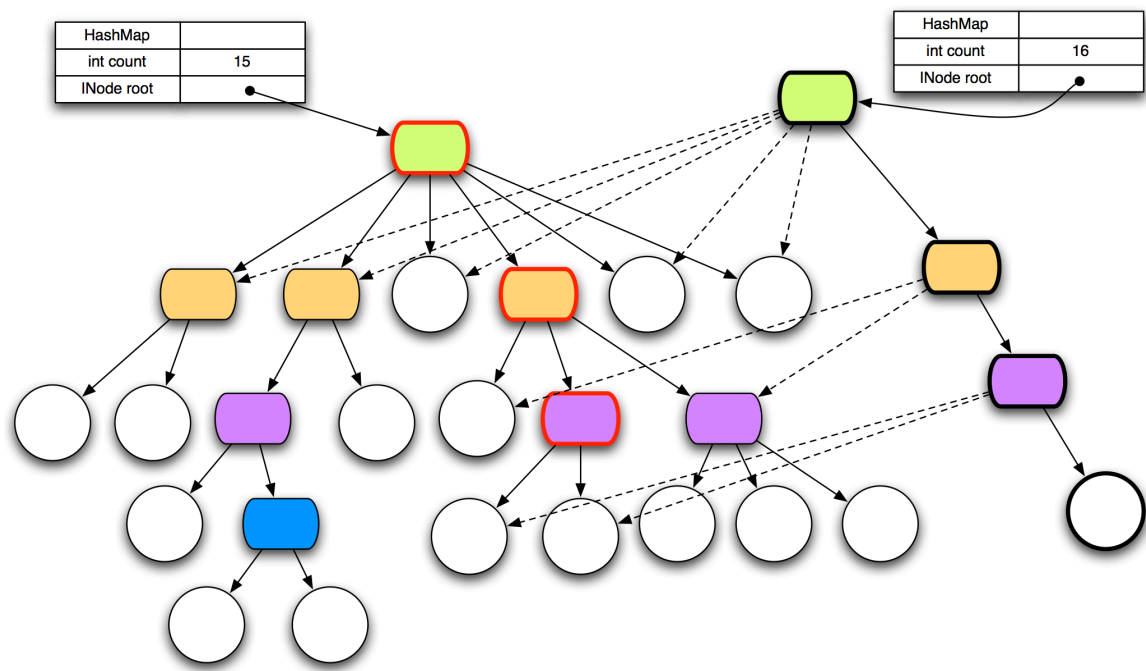


Figure 4.1: Representation of how data structures are “changed” in Clojure (Source: [Hic09])

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.

To do (2)

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	Form	Example
Member Access	(.<member> <obj> [args])	(.toString 5)
	(. <obj> <member> [args])	(. 5 toString)
	(<class>/<member> [args])	(Integer/parseInt "5")
Object Instantiation	(<class>. [args])	(Integer. 5)
	(new <class> [args])	(new Integer 5)
Multiple Operations	(doto <obj> [forms])	(doto (Vector.) (.add 1))

Table 4.2: Syntactic sugar for JVM interoperability

```

1 (defn ^Directory idx-path
2   [path]
3   (-> path File. SimpleFSDirectory.))

```

Figure 4.2: Clojure code that, given a path, returns a Directory object

Interoperability With the JVM

Traditionally, functional programming languages have been undesirable for numerous reasons: compatibility, libraries, portability, availability, packagability, and tools [Wad98]. 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 – features of a language added in order to simplify the language from a human perspective – to transparently call Java code, as well as make itself available to Java. This avoids the above issues.

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).

```
16 (defprotocol IConfig
17   (connection [this])
18   (schema [this])
19   (index [this]))
```

Figure 4.3: IConfig protocol all configurations must adhere to

4.2 Search With Clojure

Clojure’s excellent **JVM** interoperability permits the use of countless third-party libraries. The most extensively used was Lucene.

4.2.1 Full-Text Search Using Lucene

“Apache Lucene™ is a high-performance, full-featured text search engine library written entirely in Java.” [Fou]

4.2.2 Indexing Relational Database

The process of indexing a relational database is a multi-step one. It begins with the declaration of the database connection information, the path to the index, and the schema definition.

In our implementation, this information is specified by a record that adheres to the protocol in Fig. 4.3. The record which defines the Mycampus dataset uses SQLite for its database engine, so it accepts two strings; one specifies the path to the database file, while the other specifies the path to the index.

The first component, `connection`, returns a **JDBC**-compatible object. The second component, `schema`, returns a list of `EntitySchema` records. The `EntitySchema` record is defined in Section 4.2.2 on the following page. The final component, `index`, specifies the path to the index.

Key	Description	Type(s)
<code>:T</code>	Entity (<code>:entity</code>) or entity group (<code>:entity</code>)	Symbol
<code>:C</code>	Table name for entities, brief description for entity groups	Symbol or String
<code>:sql</code>	SQL query used to construct the entity or entity schema	Expression
<code>:ID</code>	Attribute or attributes that comprise the key (Definition 3 on page 3)	Symbol or list of symbols
<code>:attrs</code>	List of attributes to analyze to fields	List of symbols
<code>:values</code>	List of attributes to index as values, must be subset of <code>:attrs</code>	List of symbols

Table 4.3: Keys expected by `EntitySchema` records

Schema Graph Definition

The schema graph is defined using Korma, which “is a **domain-specific language (DSL)** for Clojure” [Gra]. Each schema component, whether an entity or entity group, is defined by `EntitySchema` records. Each record accepts a map which specifies how each class of document should be indexed and identified. The keys of this map are given in Table 4.3.

The `EntitySchema` records contain not only the information required to construct them, but also the required behaviour. Every record, given the database and index connections, is capable of retrieving the set of all named tuples it represents in the database. It iterates through every tuple and constructs a document for each tuple, as well as any value documents, if applicable.

Indexing Process

With the database, index, and schema graph defined, the system is able to transform the data from the relational model into the document model.

The first step is to retrieve the list of named tuples from the database. Each `EntitySchema` record is iterated over; the **SQL** defined for each record is executed against the database. A function is exe-

				Field	Terms
				__type__	:entity
				__class__	:course
		Key	Value	__id__	course cdps_101
		:type	:entity	__all__	“cdps 101 hu-
		:class	:course		man mutant re-
Attribute	Value	:id	“course cdps_101”		lations cdps”
code	CDPS 101	:code	“CDPS 101”	code	{cdps, 101}
title	Human-	:title	“Human-Mutant	title	{human, mu-
	Mutant Rela-		Relations”		tant, relations}
	tions				
subject	CDPS	:subject	“CDPS”	subject	{cdps}

(a) Tuple
(b) Internal Representation
(c) Document

Table 4.4: Transformation from tuple to internal representation to document

cuted on every tuple that is retrieved.

The second step is to transform these tuples into the internal representation. The function executed on every tuple is responsible for this conversion.

The third step is to convert the internal representation into a document. Again this is performed by the function operating on every tuple.

Finally, every document is indexed.

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

- Fuzzy indexing of values (typed by classes)

4.2.4 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

4.2.5 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

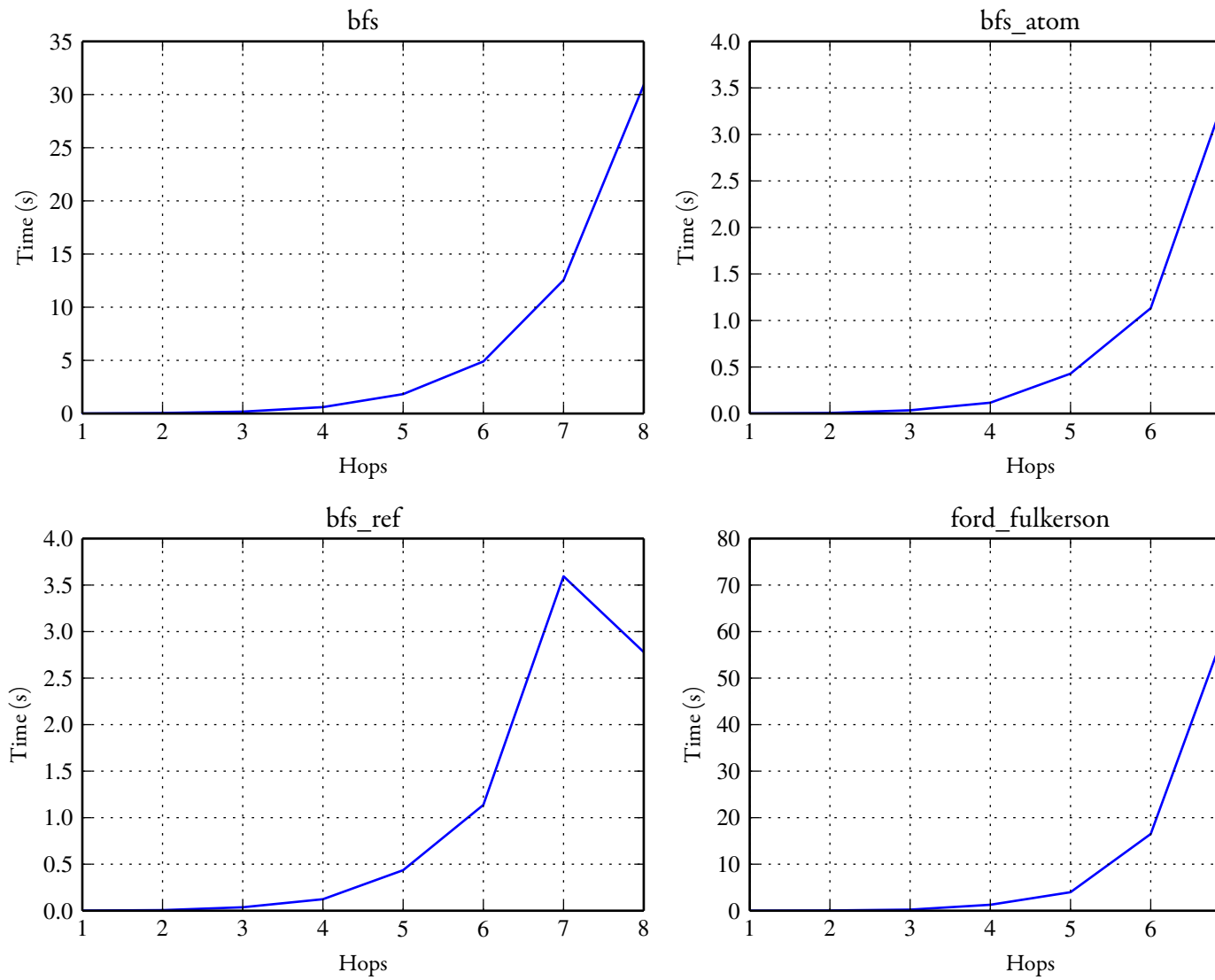


Figure 5.1: Growth of graph search times based on number of hops, plotted separately

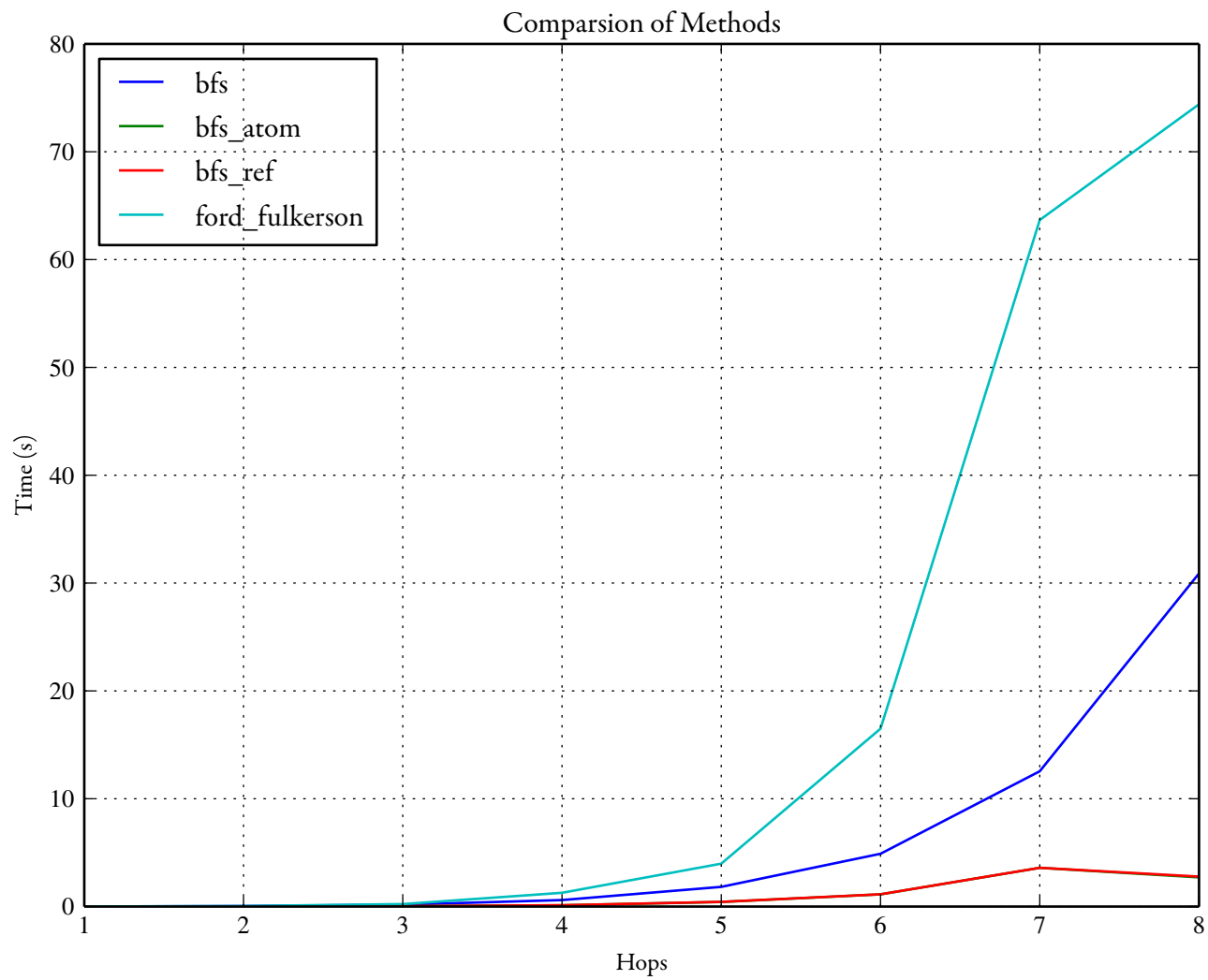


Figure 5.2: Growth of graph search times based on number of hops, combined plot

5.4 Lessons learned

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

Chapter 6

Conclusion (0 days)

Survived Clojure.

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   (:require [clojure.tools.cli :refer [cli]]
3             [molly.algo.bfs :refer [bfs]]
4             [molly.algo.bfs-atom :refer [bfs-atom]]
5             [molly.algo.bfs-ref :refer [bfs-ref]]
6             [molly.algo.ford-fulkerson :refer [ford-fulkerson]]
7             [molly.bench.benchmark :refer [benchmark-search]]
8             [molly.conf.config :refer [load-props]]
9             [molly.index.build :refer [build]]
10            [molly.search.lucene :refer [idx-path idx-searcher]])
11   (:gen-class))
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       (when (opts :index)
47         (let [database (properties :db.path)
48               index     (properties :idx.path)]
49           (build database index)))
50       (when (opts :algorithm)
51         (let [searcher (idx-searcher
52                        (idx-path
53                          (properties :idx.path)))
54               source    (opts :source)
55               target    (opts :target)
56               f          (condp = (opts :algorithm)
57                             "bfs"          bfs
58                             "bfs-atom"     bfs-atom
59                             "bfs-ref"      bfs-ref
60                             "ford-fulkerson" ford-fulkerson)
61                       (throw
62                        (Exception.
63                         "Not a valid algorithm choice.")))]
62
63         (if (opts :debug)
64           (let [[marked dist prev] (f searcher
65                                       source
66                                       target
67                                       max-hops)]
68             (println marked)
69             (println dist))
70

```

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

A.2 molly.conf

A.2.1 molly.conf.config

```
1 (ns molly.conf.config
2   (:require [proptea.core :refer [read-properties]]))
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   (:require [korma.core :refer [belongs-to
3               defentity
4               has-many
5               pk
6               select*
7               with]]
8             [korma.db :refer [defdb sqlite3]])
9   (:import (molly.conf.config IConfig)
10            (molly.datatypes.database Sqlite)
11            (molly.datatypes.schema EntitySchema)))
12
13 (declare Campus Course Subject Term Section Schedule
14         Location Instructor db-conn)
15
16 (defentity Campus
17   (has-many Location))
18
19 (defentity Location
20   (belongs-to Campus))
21
22 (defentity Subject
23   (has-many Course))
24
25 (defentity Course
26   (pk :code)
27   (belongs-to Subject)
28   (has-many Section))
29
30 (defentity Instructor
31   (has-many Schedule))
32
33 (defentity Term
34   (has-many Section))
35
36 (defentity Section
37   (pk :crn)
38   (has-many Schedule)
39   (belongs-to Term)
40   (belongs-to Course {:fk :course_code}))
41
42 (defentity Schedule

```

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

```

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

```

```

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

```

A.3 molly.datatypes

A.3.1 molly.datatypes.database

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

A.3.2 molly.datatypes.entity

```

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

```

```

43 [this schema]
44 (let [T      (schema :T)
45       C      (schema :C)
46       attr-cols (schema :attrs)
47       attrs    (if (nil? attr-cols)
48                   this
49                   (select-keys this attr-cols))
50       meta-data {:type T :class C}
51       id-col    (schema :ID)]
52 (with-meta (if (= T :group)
53               (conj attrs {:entities (uid this id-col)})
54               attrs)
55 (condp = T
56   :value   (assoc meta-data
57                   :class
58                   (clojure.string/join "|"
59                     (map name
60                        [C (first attr-cols)])))
61   :entity  (assoc meta-data :id
62                   (if (coll? id-col)
63                       (uid this id-col)
64                       (uid this C id-col)))
65   :group   (assoc meta-data
66                   :entities
67                   (uid this id-col))
68   (throw
69     (IllegalArgumentException.
70      "I only know how to deal with types :value,
71      :entity, and :group")))))
72
73 (defn doc->data
74   ^{:doc "Transforms a Document into the internal representation."}
75   [this]
76   (let [fields      (.getFields this)
77         extract     (fn [x] [(keyword (clojure.string/replace
78                                         (.name x) "_" ""))
79                               (.stringValue x)])
80         check-special (fn [x] (special? (.name x)))
81         filter-fn    (fn [f] (apply hash-map
82                                     (flatten
83                                      (map extract
84                                           (filter f fields))))))]
85     (with-meta (filter-fn (fn [x] (not (check-special x))))
86               (filter-fn check-special))))
87

```

```

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

```

A.3.3 molly.datatypes.schema

```

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

```

A.4 molly.index

A.4.1 molly.index.build

```
1 (ns molly.index.build
2   (:require [molly.datatypes.schema :refer [crawl klass]]
3             [molly.search.lucene :refer [close-idx-writer
4                                           idx-path
5                                           idx-writer]]
6             [molly.conf.mycampus])
7   (:import [molly.conf.mycampus Mycampus]))
8
9 (defn build
10  [db-path path]
11  (let [conf (Mycampus. db-path path)
12        db-conn (.connection conf)
13        ft-path (idx-path (.index conf))
14        idx-w (idx-writer ft-path)
15        schemas (.schema conf)]
16    (doseq [ent-def schemas]
17      (println "Indexing" (name (klass ent-def)) "...")
18      (crawl ent-def db-conn idx-w)
19      (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   ^{:doc "Given a string S, an integer n (optional), and a character
5         s (optional), returns the n-gram of S using s as the
6         padding character."}
7   ([S]
8    (q-gram S 3 "$"))
9   ([S n]
10    (q-gram S n "$"))
11   ([S n s]
12    (let [padding (clojure.string/join "" (repeat (dec n) s))
13          padded-S (str padding
14                        (clojure.string/replace S " " padding)
15                        padding)]
16      (clojure.string/join " "
17                            (for [i (range
18                                    (inc (- (count padded-S) n)))]
19                              (.substring padded-S i (+ i n)))))))
```

A.6 molly.search

A.6.1 molly.search.lucene

```

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

```



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

```

A.7 molly.server

A.7.1 molly.server.core

```
1 (ns molly.server.core
2   (:require [compojure.core :refer [GET defroutes]]
3             [compojure.handler :refer [site]]
4             [compojure.route :refer [not-found resources]]
5             [molly.conf.config :refer [load-props]]
6             [molly.search.lucene :refer [idx-path idx-searcher]]))
7
8 (defroutes app-routes
9   (GET "/" [] "root")
10  (resources "/")
11  (not-found "Can't find that one."))
12
13 (def config (load-props))
14 (def searcher (idx-searcher (idx-path (config :idx.path))))
15
16 (def handler
17   (site app-routes))
```

A.7.2 molly.server.remotes

```
1 (ns molly.server.remotes
2   (:require [compojure.handler :refer [site]]
3             [molly.server.core :refer [handler]]
4             [molly.server.search :refer [compute-span
5                                         find-entities
6                                         find-entity
7                                         find-value]]
8             [shoreleave.middleware.rpc :refer [defremote wrap-rpc]]))
9
10 (defremote get-value [q]
11   (find-value q))
12
13 (defremote get-entities [q]
14   (find-entities q))
15
16 (defremote get-entity [id]
17   (find-entity id))
18
19 (defremote get-span [s t method]
20   (compute-span s t method))
21
22 (def app (->
23   (var handler)
24   (wrap-rpc)
25   (site)))
```

A.7.3 molly.server.search

```

1 (ns molly.server.search
2   (:require [molly.algo.bfs :refer [bfs]]
3             [molly.algo.bfs-atom :refer [bfs-atom]]
4             [molly.algo.bfs-ref :refer [bfs-ref]]
5             [molly.algo.ford-fulkerson :refer [ford-fulkerson]]
6             [molly.datatypes.entity :refer [doc->data]]
7             [molly.search.lucene :refer [idx-search]]
8             [molly.search.query-builder :refer [boolean-query query]]
9             [molly.server.core :refer [config searcher]]
10            [molly.util.nlp :refer [q-gram]]))
11
12 (def runtime (Runtime/getRuntime))
13
14 (defn dox
15   [q field S op topk]
16   (let [bq (boolean-query
17           (concat [[q :and]]
18                   (for [s S]
19                     [(query field s) op])))]
20     result (map doc->data (idx-search searcher bq topk))
21     fmt (fn [data] {:meta (meta data) :results data})]
22     (map fmt result)))
23
24 (defn entities
25   [field q topk]
26   (dox (query :type :entity)
27        field
28        (clojure.string/split q #"\\s{1}")
29        :and
30        topk))
31
32 (defn find-value [q]
33   (dox (query :type :value)
34        :text
35        (clojure.string/split (q-gram q) #"\\s{1}")
36        :or
37        (config :idx.topk.value)))
38
39 (defn find-entities [q]
40   (entities
41    :text (clojure.string/lower-case q)
42    (config :idx.topk.entities)))

```

```

43
44 (defn find-entity [id]
45   (entities :id id (config :idx.topk.entity)))
46
47 (defn compute-span [s t method]
48   (let [max-hops (config :idx.search.max-hops)
49         start      (System/nanoTime)
50         [visited dist prev]
51         (condp = method
52           "bfs" (bfs searcher s t max-hops)
53           "atom" (bfs-atom searcher s t max-hops)
54           "ref" (bfs-ref searcher s t max-hops)
55           "ff" (ford-fulkerson searcher s t max-hops))
56         time-taken (- (System/nanoTime) start)
57         eids      (conj (for [[k v] prev] k) s)
58         get-entities (fn [eid]
59                       {(keyword eid)
60                        (entities :id eid
61                                (config :idx.topk.entity)))})
62         entities    (into {} (map get-entities eids))]
63   {:from      s
64    :to        t
65    :prev      prev
66    :entities  entities
67    :debug     {:time      time-taken
68                :mem_total (.totalMemory runtime)
69                :mem_free  (.freeMemory runtime)
70                :mem_used   (- (.totalMemory runtime)
71                               (.freeMemory runtime))
72                :properties config}}))

```

A.7.4 molly.server.util

```
1 (ns molly.server.util)
2
3 (println "loading properties")
```

A.8 molly.algo

A.8.1 molly.algo.common

```

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

```



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

A.8.2 molly.algo.bfs

```

1 (ns molly.algo.bfs
2   (:require [molly.algo.common :refer [find-adj]]))
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                 frontier' (conj frontier v)]
20             (recur adj' marked dist' prev' frontier'))))))))
21
22 (defn bfs
23   [G s t max-hops]
24   (loop [Q      (conj (clojure.lang.PersistentQueue/EMPTY) s)
25          marked #{s}
26          dist   {s 0}
27          prev   {s nil}]
28     ; Terminate when nothing is left to explore.
29     (if (seq Q)
30       (let [u      (first Q)
31             Q'      (rest Q)]
32         ; Terminate when the target is found, or we reach a limit.
33         (if (or (= u t)
34                 (>= (dist u) max-hops))
35           [marked dist prev]
36           (let [[marked' dist' prev' frontier]
37                 (update-adj G marked dist prev u)]
38             (recur (concat Q' frontier) marked' dist' prev')))))
39     [marked dist prev]))

```

A.8.3 molly.algo.bfs_atom

```

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

```

A.8.4 molly.algo.bfs_ref

```

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

- [Fou] Apache Software Foundation. URL: <http://lucene.apache.org/core/>.
- [Gra] Chris Granger. Korma: tasty sql for clojure. URL: <http://sqlkorma.com/>.
- [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/d1/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. 22): Above proof could use some love
- ☐ 2 (p. 27): More concurrency