

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

April 2014

© Richard J.I. Drake, 2014

Contents

1	Background & Introduction	1
1.1	Structured Data and Structured Language: 1970 – 2000	1
1.2	Text data and keyword search: 1970 – 2014	2
1.3	Semi-structured Data and Query Languages: 1990 – 2010	2
1.4	Hybrid data models and query languages: 2010 – Present	3
1.5	Outline of Thesis	4
2	A Tale of Two Data Models	5
2.1	Relational Model	5
2.1.1	Schema Group	7
2.1.2	Entity Group	10
2.1.3	Pros and Cons of the Relational Model	11
2.2	Document Model	14
2.2.1	Vectorization of Documents	15
2.2.2	Extending the Document Model	19
2.2.3	Approximate String matching	19
2.2.4	Pros and Cons of the Document Model	21
3	Best of Both Worlds	23
3.1	Encoding Named Tuples into Documents	23
3.2	Mapping of Entity Groups to Documents	24

3.3	Encoding an Entity Group as a Document Group	24
3.4	Encoding Attribute Values into Searchable Documents	25
3.5	Iterative Search Using Document Encodings	26
4	Along Came Clojure	27
4.1	Basic Principles of Functional Programming	27
4.1.1	Features of Clojure	28
4.2	Search With Clojure	31
4.2.1	Full-Text Search Using Lucene	31
4.2.2	Indexing Relational Database	31
4.2.3	Keyword Search in Document Space	35
4.2.4	Graph Search in Document Space	36
5	Experimental Evaluation	38
5.1	Implementation	38
5.1.1	Code Base Statistics	38
5.2	The Data Corpus	39
5.3	Runtime Evaluation	39
5.3.1	Methodology	40
5.3.2	Performance	41
6	Conclusion	45
6.1	Summary	45
6.2	Lessons Learned	46
A	Source Code	47
A.1	molly	47
A.1.1	molly.core	47
A.2	molly.conf	50

A.2.1	molly.conf.config	50
A.2.2	molly.conf.mycampus	51
A.3	molly.datatypes	55
A.3.1	molly.datatypes.database	55
A.3.2	molly.datatypes.entity	56
A.3.3	molly.datatypes.schema	59
A.4	molly.index	61
A.4.1	molly.index.build	61
A.5	molly.util	62
A.5.1	molly.util.nlp	62
A.6	molly.search	63
A.6.1	molly.search.lucene	63
A.6.2	molly.search.query_builder	65
A.7	molly.server	66
A.7.1	molly.server.core	66
A.7.2	molly.server.remotes	67
A.7.3	molly.server.search	68
A.8	molly.algo	70
A.8.1	molly.algo.common	70
A.8.2	molly.algo.bfs	72
A.8.3	molly.algo.bfs_atom	73
A.8.4	molly.algo.bfs_ref	74
A.9	molly.bench	75
A.9.1	molly.bench.benchmark	75
B	Data Corpus	76

List of Tables

2.1	Course relation	6
2.2	Results of the query in Fig. 2.4 on page 10.	10
2.3	Properties of the Course titled Human-Mutant Relations.	11
2.4	Properties of the Course titled Human-Mutant Relations.	19
2.5	Course document for MATH 360.	20
3.1	Doc[<i>t</i>]	24
3.2	Document encoding of Fig. 3.1 on page 25	26
4.1	Comparison between Clojure’s four systems for concurrency	29
4.2	Java virtual machine (JVM) interoperability	30
4.3	Keys expected by EntitySchema records	32
4.4	Metadata associated with each type	33
4.5	Document of internal representation from Fig. 4.4 on page 34	34
5.1	Number of objects in data corpus, grouped by class	39
5.2	Indexing time growth by number of entity groups	42
B.1	Structure of data corpus	77

List of Figures

2.1	Subset of mycampus dataset schema	8
2.2	foreign key (FK) constraints on schema in Fig. 2.1 on page 8	8
2.3	Graph representation of relations (Fig. 2.1) and FK (Fig. 2.2)	9
2.4	Query to find section CRNs for a subject name.	10
2.5	Human-Mutant Relations entity group	11
2.6	Comparison between n -grams of G and G'	21
3.1	Example entity group	25
4.1	Representation of how data structures are “changed” in Clojure (Source: [Hic09]) . .	29
4.2	Clojure code that, given a path, returns a <code>Directory</code> object	31
4.3	Building the index	32
4.4	Internal representation of named tuple from Table 2.4 on page 19	34
4.5	Boolean query in Clojure	36
4.6	Human-Mutant Relations entity group	36
5.1	Partial JavaScript object notation (JSON) output from Criterium.	41
5.2	Comparison of time taken per hop between all methods (top- k entities: 25), target found	43
5.3	Comparison of time taken per hop between all methods (top- k entities: 25), target not found	44

List of Algorithms

1	N-GRAM(S, n, s)	20
---	-------------------------------	----

Acronyms

API application programming interface. 3, 31, 35, 36

BFS breadth-first search. 36, 42, 43, 46

DSL domain-specific language. 31

FF Ford-Fulkerson. 36

FK foreign key. vii, 6–11

HTML HyperText Markup Language. 2, 39

HTTP HyperText Transfer Protocol. 66–68

JAR Java archive. 47

JDBC Java database connectivity. 30

JSON JavaScript object notation. vii, 3, 40, 41

JVM Java virtual machine. vi, 30, 31, 38, 40, 41

MDX MultiDimensional eXpressions. 2

OLAP online analytical processing. 2

RDBMS relational database management system. 1, 2, 12, 14

SQL structured query language. 1–3, 7, 10, 21, 32, 59

STM software transactional memory. 29, 46

TF-IDF term frequency and inverse document frequency. 16

UOIT University of Ontario Institute of Technology. i, 39

XML Extensible Markup Language. 2, 3

List of Symbols

- schema graph** (G) graph representation of schema. 7, 10, 24
- entity group** (T) forest of named tuples. 10, 24, 25
- document collection** (C) set of documents. xi, 15–19, 25
- terms** (T) set of unique terms in a document collection. 15, 16, 19
- document** (d) set of fields. xi, 14–19, 21–23
- search query** (q) special case of document. 17–19, 21, 22, 26
- field** (f) named sub-document in document. 23, 26, 35
- term** (τ) unique term in document collection. xi, 14–17, 21
- N number of documents in collection. 15
- database** (D) set of relations. 7, 10, 24
- relation** (r) set of named tuples. xi, 5–7, 10, 24
- named tuple** (t) ordered set of values. vi, xi, 5–7, 10, 23–25
- attribute** (α) named column. 6, 7, 10, 23
- key** (K) uniquely identifies a named tuple in a relation. 6

Chapter 1

Background & Introduction

In this chapter, we provide a background to the motivation of this thesis. We will discuss the evolution of data sets, their logical models, and the corresponding query languages. We feel that modern day data sets call for a new data model with a new query language. This thesis is an attempt to make an incremental advance toward a new data model and new way of querying data.

1.1 Structured Data and Structured Language: 1970 – 2000

The proliferation of database system research and development started with the disruptive invention of the relational data model by Edgar F. Codd [Cod79].

The invention of the relational data model was a significant achievement as it decoupled the task of data analytics from any one language, precisely and accurately described data sets across a variety of storage and analytical systems, and lead to the creation of the structured data analytics language known as structured query language (SQL).

SQL itself deserves further discussion; the relational data model provided a foundation upon which languages for data manipulation were designed. One can describe their data set and operations using first-order logic and relational algebra, then realize it using SQL.

The success of the relational model can only be appreciated when one looks at the continuous success of relational database management system (RDBMS) such as Oracle and IBM DB2 which

span over 3 decades without any significant decline.

Since the 90s, the emergence of Business Intelligence [CKKS02] furthered the development of RDBMS by specializing the relational data model to multidimensional data model [Col96]. The family of databases known as online analytical processing (OLAP) produced a new query language known as MultiDimensional eXpressions (MDX).

Both SQL and MDX are highly structured query languages: they are completely described by their respective grammar and operational semantics. Users who wish to harness the power of SQL and MDX must be well versed in the languages themselves, and understand precisely the semantics of each syntactic constructs.

1.2 Text data and keyword search: 1970 – 2014

Parallel to the development of the relational data base technology, research in the information retrieval has been focusing on text data [SB88; Jon72]. Unlike the relational data, text data does not have much complex structure to its schema. Thus, it's not immediately possible to design a rich set of data operators (as was the case for relational algebra). Consequently, for text data, there is no structured query language like SQL.

The research, thus, has been focusing on pattern matching queries using keyword search. Though the semantics of keyword search is very simple, the statistical methods developed by the information retrieval community [SB88; RZ09; DFLDH88] have been extremely effective. In fact, one can argue that the modern World Wide Web and its related commercial successes founded on the ideas of text databases and keyword search over Internet scale data sets.

1.3 Semi-structured Data and Query Languages: 1990 – 2010

The growth of the World Wide Web popularized the usage of markup languages (such as HyperText Markup Language (HTML) and Extensible Markup Language (XML)) as the underlying Web content description. Thus, researchers have designed data models [Suc98] to formalize the semantics of

XML and related data formats. Subsequently, the logical characterization of XML led the to design and implementation of XQuery [10], a navigational based query language for analysis of XML data sets.

Interestingly enough, XML has proven to be inefficient as a data description format. Nonetheless, a semi-structured data description language is highly sought after for message passing in Internet scale software systems. Modern Web services are built on the concept of RESTful application programming interface (API) [Sch11], with semi-structured data message passing. In order to minimize network overhead, XML based message passing has been replaced by the more efficient data encoding standard of JSON [ECM13].

The query language community responded to the popularization of JSON encoded data sets with several query languages [BCNS13] (for example Jaql []) for JSON data sets.

1.4 Hybrid data models and query languages: 2010 – Present

With the explosive growth of social networks, we are witnessing the emergence of a new type of data sets. These data sets exhibit the following properties:

- The data has relational characteristics: such as relationships of friends on Facebook, their preferences over different Web sites, and their account information.
- The data also has many text attributes: such as blog articles, or tweets on Twitter.
- The volume of data is often Internet scale.

The mixture of relational structure and rich text components of such data sets make them challenging for the purpose of data management and data analytics. There has been several attempts to integrate keyword search from information retrieval with SQL [BHNC02; LJLF11; HGP03]. However, these methods, thus far, suffer from scalability and restricted search capabilities.

In this thesis, we are motivated by the following problems:

- Define a formal framework to describe data sets with relational structures and text components.

- Design a collection of expressive query operators for analyzing text relational data sets.
- Investigate implementation techniques to make the query operators performant for modern multi-core computers.

1.5 Outline of Thesis

Chapter 2 provides a formal definition of the relational data model and the document data model. Rather than introducing a new hybrid model, our thesis is to provide efficient mappings between the two models, and thus allow the data to freely be exchanged between query operators of the domains. We demonstrate that complex relational data can be encoded into a linked document space (and back).^{To do (1)}

Chapter 4 describes the details of our implementation the data transformation and query operators for graph search in the linked document space. Our choice of utilizing a modern functional programming language for our implementation makes high degree of concurrency possible.

Chapter 5 provides further justification to our choice of data model mapping and the choice of programming language used. Through a series of experiments, we see that our proposal allows a tight integration of the relational database engines and keyword search libraries. Furthermore, the implementation enjoys a linear speed up with respect to the number of processors available.

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 page 6), keys (as defined by Definition 3 on page 6), and foreign key constraints (as defined by Definition 4 on page 6).

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 9

The relational data model is particularly powerful for analytic queries. Given the schema graph in Fig. 2.3 on page 9, one can formulate the following analytic queries in a query language known as 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

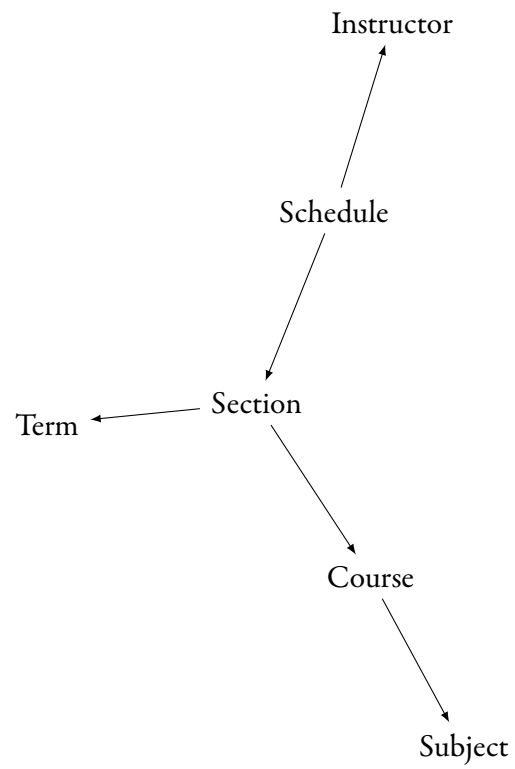


Figure 2.3: Graph representation of relations (Fig. 2.1) and FK (Fig. 2.2)

```

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.

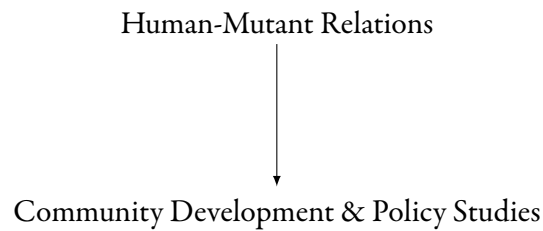


Figure 2.5: Human-Mutant Relations entity group

Example 6. The information in Table 2.4 on page 19 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. 4.6 on page 36) 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 6) 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 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 10.

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 page 15, 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 C with M unique terms $T = [\tau_1, \tau_2, \dots, \tau_n]$,

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 15 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 14). 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 C (as defined in Example 12 on page 15).

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

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 C (from Example 12 on page 15) and search query q , compute the similarity between q and every document $d \in C$.

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

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

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

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

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

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 page 21.

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

While none of the above queries required domain knowledge, it is possible to use the extended document model (Section 2.2.2 on page 19) 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 19), users who substitute, delete, or insert characters from the desired term may still receive results for their intended term (see Example 19 on page 20 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 6). 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 19), 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 10) 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 (2)

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.

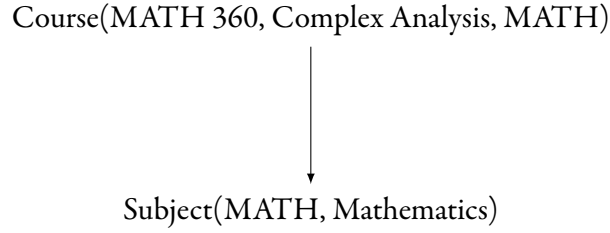


Figure 3.1: Example entity group

For every $t \in V(T)$, we construct a document $\text{Doc}[t]$ (Section 3.1 on page 23). 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 9 is shown in Fig. 3.1.

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

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

By Claim 1 on page 24, 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.

Field	Terms	Field	Terms	Field	Terms
__id__	{subject math_360}	__id__	{subject math}		
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 on page 25

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

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

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”.

To do (3)

Chapter 4

Along Came Clojure

To do (4)

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 next 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.

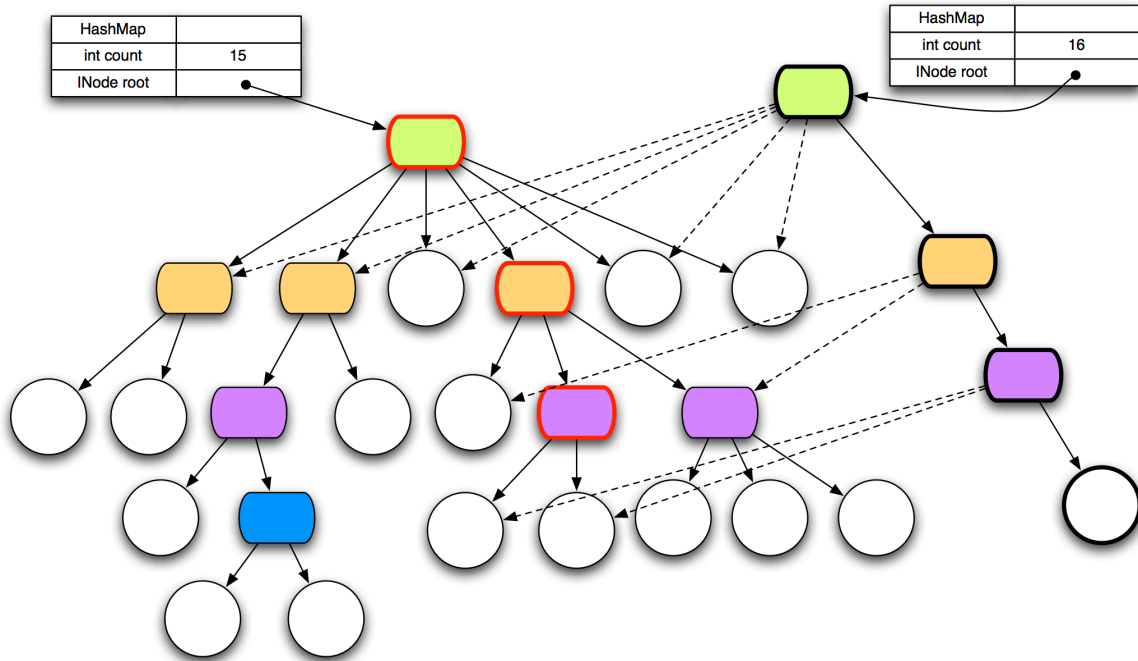


Figure 4.1: Representation of how data structures are “changed” in Clojure (Source: [Hic09])

Concurrency

Clojure supports four systems for concurrency: software transactional memory (STM), agents, atoms, and dynamic vars. The differences between these systems – including whether or not they are synchronous, coordinated, and what scope they encompass – are summarized in Table 4.1.

In this context, synchronous refers to the fact that each operation waits for the proceeding operation to complete before continuing. When a system is coordinated, operations occur in a transaction.

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

Table 4.2: JVM interoperability

If one operation fails, the entire transaction fails and is rolled back. This results in additional overhead, but is considered safer as the system shall not be left in an inconsistent state as a result of a failed transaction.

To do (5)

Interoperability With the JVM

To do (6) 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).

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

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

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] Lucene implements the Document Model (Section 2.2), providing a simple yet powerful API to perform full-text search. Among these features is the ability to vectorize documents according to the Vector Space Model (Section 2.2.1), utilize the extended document model to provide semi-structured documents (Section 2.2.2), and issue search queries against the index (Definition 12).

4.2.2 Indexing Relational Database

The indexing of relational objects is a complicated process. The objects must be retrieved from the relational database, transformed from named tuples into documents, then placed in the index. Additionally, all foreign keys (Definition 4) must be encoded as documents (Section 3.3) in order to satisfy Section 3.2.

The schema graph (Definition 6) must be defined before the relational database may be crawled.

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

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 6)	Symbol or list of symbols
<code>:attrs</code>	List of attributes to analyze to fields	List of symbols
<code>:values</code> <code>:attrs</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

```

1 (doseq [ent-def schemas]
2   (crawl ent-def db-conn idx-w))

```

Figure 4.3: Building the index

and identified. The keys of this map are given in Table 4.3.

Crawling the Relational Database

With the schema graph defined, the system is able to crawl the relational database, yielding a sequence of named tuples. It iterates through every `EntitySchema` record, instructing every record to crawl itself given a database connection and index writer (Fig. 4.3).

The `Database` protocol provides an `execute-query` method that permits access to the database. In the current implementation, the `SQLite3` record implements the `Database` protocol. This record issues the query as-is, applying a given function to every result.

Each record issues a SQL query against the database that retrieves all named tuples that it represents. This query is given by the `:sql` key of the record. For every symbol defined by the `:values` key, an additional query is issued. These queries retrieve all distinct (within the context of that relation and

	Key	Value		Key	Value
Key	Value		:type	:entity	
:type	:value		:class	<rel>	
:class	<rel> <attr>		:ID	<rel> <pk>	
					[<rel> <pk> ...]
(a) Value			(b) Entity		
			(c) Entity Group		

Table 4.4: Metadata associated with each type

attribute) values.

Transformation

For every named tuple, a document is constructed (see Section 3.1). In addition to the attributes, several other fields may be added to the document. These special fields contain additional meta information about the document. For example, the class, type, and unique identifier are added to an entity, while an entity group has a space-delimited list of unique identifiers that comprise the group.

Before becoming a document, named tuples are transformed into an internal representation. The internal representation adds flexibility to the system; so long as functions exist to convert between the internal representation and other forms, the system does not care about the source.

Clojure permits the annotation of data with metadata. Named tuples are returned as maps, with key-value pairs representing attributes and values for each tuple. Metadata may be associated with a named tuple that does not affect its key-value pairs.

The map of attributes and values of a named tuple is annotated with the function (`with-meta obj map`). The `obj` parameter is the named tuple, while `map` is a map of metadata as defined by the system. The keys of `map` for each type (value, entity, or entity group) are defined in Table 4.4.

The internal representation of the named tuple given in Table 2.4 on page 19 is given in Fig. 4.4 on the next page.

Once the internal representation is constructed, it may be transformed into a document. The

```

1 (with-meta
2   {:code    "CDPS 101"
3    :title   "Human-Mutant Relations"
4    :subject "CDPS"}
5   {:type :entity
6    :class :course
7    :id    "course|cdps_101"})

```

Figure 4.4: Internal representation of named tuple from Table 2.4 on page 19

Field	Value
code	cdps 101
title	human mutant relations
subject	cdps
__type__	entity
__class__	course
__id__	course cdps_101

Table 4.5: Document of internal representation from Fig. 4.4

mapping of a map to document is trivial; a field is created for every key in the map and the value corresponding to the key is the value of the field. Unfortunately Lucene does not facilitate the storage of metadata. Therefore the system must deal with metadata in a different way.

The system modifies each key of the metadata; two underscores are prepended and appended to the key name. This allows the system to differentiate between metadata and attributes. The transformation from internal representation to document is given in Table 4.5.

Indexing

With the named tuples transformed into documents, the index may be constructed. The first step is to open the index for writing. In Lucene, this is accomplished by creating an `IndexWriter` object on a `Directory` object that points to the index location. The `IndexWriter` object also expects an

analyzer to be used by default on documents it indexes. The system uses a `WhitespaceAnalyzer` by default.

For every named tuple, the transformed document is written to the index by the index writer object. In addition, the indexing document (Section 3.3) of every entity group is added.

4.2.3 Keyword Search in Document Space

With the entity graph encoded into the document model, users may begin issuing search queries. Every query follows the following pattern: users look up values (optionally using fuzzy search), these values are used to locate entities, and once two entities are selected, the system attempts to locate connections between the two.

Fuzzy Value Search

Recall that entity values are stored as their n -gram (Section 2.2.3). This allows users to make character substitutions, deletions, or additions, and still return values they may have intended on finding. Without the use of n -grams, a misspelling on the user's part would result in either no or irrelevant values being returned. Values that are approximately matched would be assigned a lower score than those which are fully matched, but they would at least appear in the results.

Rather than guessing the user's intention, the system presents the user with a list of values in order to give them the option of substituting their entry for an approximate match. This autosuggest feature is intended to improve the user experience by eliminating a source of frustration – irrelevant results as a result of a simple spelling error.

Flexible Keyword Search API

The system provides a simple – yet flexible – keyword search API. Recall the extended query (Example 23) is in the form

$$\text{query}(f, w)$$

The search API provides a function that accepts a symbol defining the field to search in, as well

```

1 (boolean-query
2   [[:and (query :subject 'MATH')]]
3   [[:and (query :text 'class')]])

```

Figure 4.5: Boolean query in Clojure

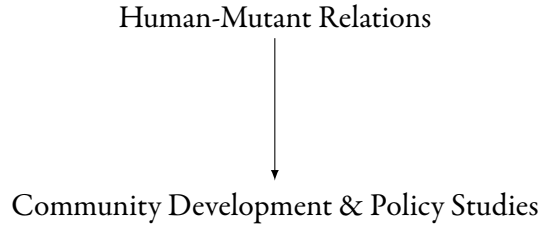


Figure 4.6: Human-Mutant Relations entity group

one or more words to search for. A phrase query comprised of every word is constructed.

$$\text{query}(f, w_1, w_2, \dots, w_n) = \bigcap_{w \in \{w_1, w_2, \dots, w_n\}} \text{query}(f, w)$$

Another function, `boolean-query`, accepts one or more query functions as well as a boolean operator for each and returns the result. The symbols `:and`, `:or`, and `:not` provide \cap , \cup , and \neg , respectively.

Example 29 (`boolean-query`). For example, the query in Example 23 is given in Fig. 4.5.

Translation of Search Results to Relational Space

4.2.4 Graph Search in Document Space

With the ability for users to find relevant entities using fuzzy value search and the flexible keyword search API (Section 4.2.3), they must be able to find connections between two entities. As stated previously (Claim 1), the document encoding of relational data is a graph. This allows the system to search for links between entities by utilizing one or more graph search algorithms, such as `breadth-first search` (BFS).

- 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

In this chapter we evaluate our implementation of the system for transforming data in the relational model to the document model and vice versa described in Chapter 2. We cover the implementation details in Section 5.1, and the methodology and evaluation in Section 5.3.

5.1 Implementation

The system was implemented in Clojure, which “is a dynamic programming language that targets the JVM” [Hica]. Clojure was chosen due to its rich, immutable, and persistent data structures, excellent concurrency support, and seamless JVM interoperability. These features were discussed in detail in Section 4.1.1.

5.1.1 Code Base Statistics

The system consists of over 800 lines of Clojure, along with approximately 550 lines of Python. The Python code is used to construct the data set by crawling the course information site, as well as to aggregate the benchmark data produced by the system, producing graphs.

All development has occurred on GitHub [Dra]. ^{To do (7)}

5.2 The Data Corpus

The data corpus was derived from the University of Ontario Institute of Technology (UOIT) mycampus database. An HTML crawler was written in Python that scraped the information from the UOIT class schedule search page. This data was parsed, normalized, then placed in a SQLite database.

The data corpus consists of numerous classes of objects. These are campuses (Table B.1a), courses (Table B.1c), instructors (Table B.1b), locations (Table B.1d), schedules (Table B.1g), sections (Table B.1h), subjects (Table B.1e), and terms (Table B.1f). A graph representation of how these classes of objects are related can be found in Fig. 2.3. The data corpus is defined in Appendix B

The number of objects, as of the publication of this thesis, can be found in Table 5.1.

Class	Count
Campus	8
Course	1797
Instructor	560
Location	220
Schedule	20985
Section	6211
Subject	68
Term	32

Table 5.1: Number of objects in data corpus, grouped by class

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 issues, identified by Boyer [Boy08], 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 garbage collection after testing to estimate impact of cleanup on the timing results

This requires a much longer runtime as each function must be invoked numerous times.

During evaluation, Criterium collects performance metrics. Upon completion of the evaluation, it performs statistical analysis of these metrics using the bootstrap procedure developed by Efron [Efr87]. These metrics include mean, samples, variance, quartiles, outliers, and more.

Data Collection

The performance metrics computed by Criterium are returned as a Clojure map data structure. The evaluation process may take several hours to complete, necessitating a separation between data collection and post-processing. These metrics are stored offline for further processing.

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 the output is captured by the benchmark script. An example of this JSON output is given in Fig. 5.1.

¹<http://hugoduncan.org/criterium/>

```
[{  
  "max-hops": ...,  
  "method": ...,  
  "results": {  
    "execution-count": ...,  
    "final-gc-time": ...,  
    "lower-q": [...],  
    "mean": [...],  
    ...  
  }, ...]
```

Figure 5.1: Partial JSON output from Criterium.

5.3.2 Performance

Performance was measured for the various system components. An analysis of the metrics collected is presented in this section.

Indexing

The indexing process is computationally intensive but short lived. After the initial JVM warmup period, the time required to construct the index scales with the number of named tuples and relations between them.

Number of Groups	Elapsed Time (s)
1	16.895
2	21.624
3	24.436
4	24.528
5	25.334

Table 5.2: Indexing time growth by number of entity groups

We see in Table 5.2 the indexing time increases considerably between 1 and 2 groups. The number of entity groups also grew considerably, explaining the time increase. The number of entity groups grew at a slower pace from 2 to 3, and 3 to 5, causing the reduced increase in elapsed time.

Keyword Search

Graph Search

The worst-case performance of BFS is $\mathcal{O}(n^2)$. This is reflected in Fig. 5.3 which follows an exponential growth curve. In an attempt to mitigate the rapid increase in search time, concurrent variants of BFS were also implemented and benchmarked.

We see in Fig. 5.3 the decrease in rate of growth is not as high as expected. BFS is the slowest, closely followed by BFS with references and BFS with atoms.

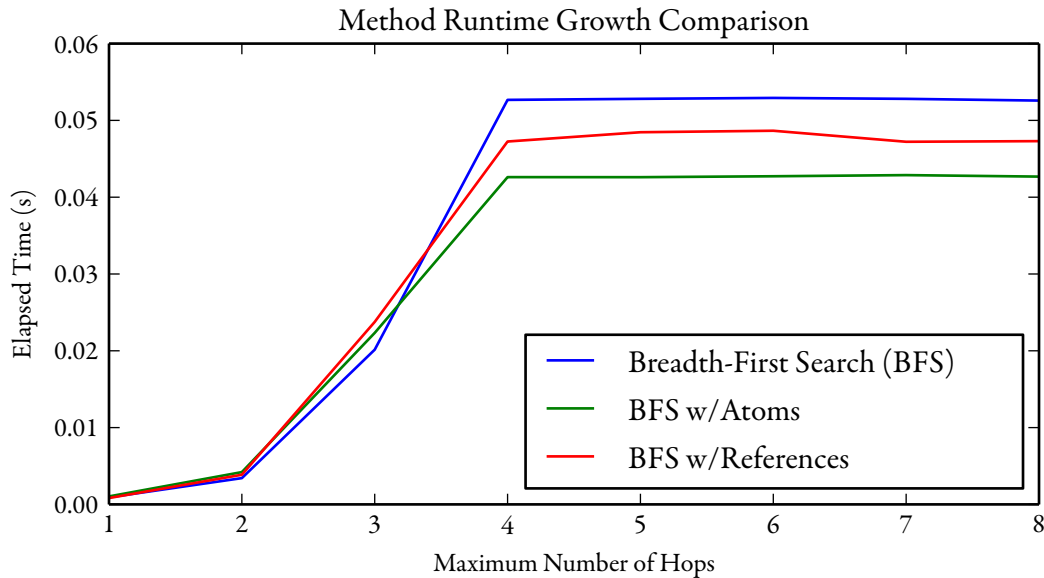


Figure 5.2: Comparison of time taken per hop between all methods (top- k entities: 25), target found

The most likely cause for the similarity is how the concurrency is implemented. In our implementation, we concurrently search for neighbours. These neighbours must be added to the frontier before proceeding. Therefore the algorithm must synchronize upon completion of the neighbour search in order to populate the new frontier.

Atoms are less expensive to use than references as they are handled in an uncoordinated manner in Clojure, reducing the amount of overhead [Nar12]. This explains the difference between the two concurrent implementations of BFS.

Also of note is the growth of Fig. 5.2 compared to that of Fig. 5.3. In the former figure, the target is located by the fourth hop. In the latter, the target is not found. This leads to hops 4 through 8 having approximately the same runtime as one another.

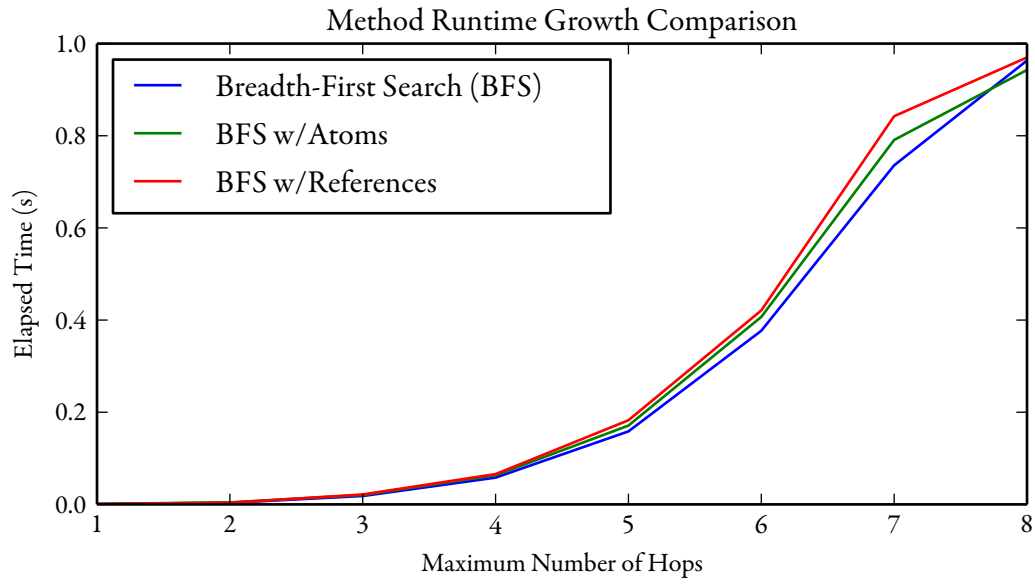


Figure 5.3: Comparison of time taken per hop between all methods (top- k entities: 25), target not found

Chapter 6

Conclusion

6.1 Summary

While the relational model is a powerful and well understood method of storing data, it is not without its shortcomings. The rigidity of the relational model comes at the cost of usability. A change to the data model may require a rewrite of queries to account for the different join paths, increasing the cognitive burden on users.

In contrast, the document model represents semi-structured and unstructured data. The queries issued against the document model are unstructured and flexible. This allows users with little or no prior domain knowledge to issue queries. Unfortunately this flexibility comes at the cost of foreign keys, data consistency, and aggregate queries.

In Chapter 2 we introduced the relational and document data models. We compared and contrasted the two data models, paving the way for our contribution. We introduced our contribution in Chapter 3, providing a formal definition of a framework for the translation between the relational and document data models. Our implementation was introduced in Chapter 4, which also covered Clojure. In Chapter 5, we evaluated our implementation, identifying performance characteristics.

6.2 Lessons Learned

The system evaluation in Chapter 5 yielded several important insights.

- We have found that the simplest algorithms are the easiest to parallelize. The reduced complexity, and thus state, reduces the amount of shared data that must be synchronized. This allows for higher concurrency.
- Clojure's `STM` implementation is simple to use and effective. A few simple functions provide a powerful concurrency model.
- Sometimes a simpler approach to concurrency is the most appropriate one. In our evaluation, atoms provided better performance than references. Atoms allow for finer granularity in concurrency, reducing the overhead associated with references. This is desirable in situations that do not require much shared state.
- Clojure is a powerful language that encouraged us to write correct code first, then optimize it later. The transition to a concurrent implementation of `BFS` was trivial. The switch from atoms to references was also trivial.

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

The core namespace is responsible for determining, provided a series of command line arguments, which action to take. This namespace is invoked as the main class when the Java archive (JAR) is executed.

```
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         (if (opts :debug)
63           (let [[marked dist prev] (f searcher
64                                       source

```

```
67                                     target
68                                     max-hops)]
69     (println marked)
70     (println dist)
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

This namespace contains helper functions for loading part of the system configuration. It also provides a protocol, `IConfig`, that is used to define the rest of the system configuration.

```
1 (ns molly.conf.config
2   (:require [propertea.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

This is a sample configuration. It defines the entities and relations in the mycampus dataset.

```

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

There are several datatypes used in the system.

A.3.1 molly.datatypes.database

A protocol and concrete datatype are defined which provide access to a relational database. Users creating an instance of this datatype are able to execute arbitrary queries, and must provide a function to apply to every tuple that is returned.

```
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

One of the most important namespaces represents entities. It includes functions to transform a named tuple from a database row into the internal representation as well as into documents. It also includes auxiliary functions to produce a unique identifier.

```

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))

```

[illegible]

```

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

The final datatype represents a schema. These schemas contain a function that is used to execute the necessary SQL statements to retrieve all data from the relational database and place it in the full-text search index.

Several of these schema datatypes are joined together in a configuration in order to produce a schema graph.

```

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

This namespace contains the function used to build the full-text search database. It takes advantage of the fact that each schema knows how to construct its own documents. The function simply iterates through every schema in the configuration, instructing them to index themselves.

```
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

The `q-gram` function computes the q -gram of a string. Optionally a value for q and the padding character can be specified.

```

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

This namespace contains various functions for interfacing with the Lucene library. These functions include opening, adding documents, searching, and closing indices.

```

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

Phrase queries are used as they require each term in the phrase to be in a specific order. This permits more accurate results as course titles and other items are in a specific order.

These queries may be combined, creating a boolean query.

```

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

This namespace contains functionality to expose the system functionality to clients over HyperText Transfer Protocol (HTTP).

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

Rather than handle serialization over HTTP manually, the system uses the Shoreleave library¹. It permits ClojureScript clients to transparently call functions exposed on the server. The `defremote` macro is used to expose these functions.

```
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)))
```

¹<https://github.com/shoreleave/shoreleave-remote>

A.7.3 molly.server.search

This namespace provides the “glue” between the system and HTTP interface.

```

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.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 t]
28   {:Q (conj (clojure.lang.PersistentQueue/EMPTY) s)
29    :marked #{s}
30    :dist {s 0}
31    :prev {s nil}
32    :done false
33    :target t})
34
35 (defn update-state
36   [state u v max-hops]
37   (let [Q (state :Q)
38         marked (state :marked)
39         dist (state :dist)
40         prev (state :prev)]

```



```
41     target (state :target)
42     done   (or (>= (dist u) max-hops)
43              (= v target))
44 (if done
45     (assoc state
46          :marked (conj marked v)
47          :dist   (assoc dist v (inc (dist u)))
48          :prev   (assoc prev v u)
49          :done   done)
50     (assoc state
51          :Q       (conj Q v)
52          :marked  (conj marked v)
53          :dist    (assoc dist v (inc (dist u)))
54          :prev    (assoc prev v u))))
55
56 (defn deref-future
57   [dfd]
58   (if (future? dfd)
59       (deref dfd)
60       dfd))
```

A.8.2 molly.algo.bfs

```

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

Appendix B

Data Corpus

Tables representing various object classes of data corpus used in thesis implementation and evaluation.

Property	Data Type
id	INT
name	VARCHAR

(a) Campus Data Structure

Property	Data Type
code	VARCHAR
title	VARCHAR
subject_id	VARCHAR

(c) Course Data Structure

Property	Data Type
id	INT
name	VARCHAR

(b) Instructor Data Structure

Property	Data Type
id	INT
name	VARCHAR
campus_id	VARCHAR

(d) Location Data Structure

Property	Data Type	Property	Data Type
id	VARCHAR	id	VARCHAR
name	VARCHAR	name	VARCHAR
(e) Subject Data Structure		(f) Term Data Structure	
Property	Data Type	Property	Data Type
id	INT	crn	VARCHAR
days	VARCHAR	reg_start	DATE
sch_type	VARCHAR	reg_end	DATE
date_start	DATE	credits	FLOAT
date_end	DATE	section_num	VARCHAR
week	VARCHAR	term_id	VARCHAR
time_start	TIME	course_code	VARCHAR
time_end	TIME	levels	VARCHAR
location_id	INT		
section_id	INT		
instructor_id	INT		
(g) Schedule Data Structure		(h) Section Data Structure	

Table B.1: Structure of data corpus

Bibliography

- [] What is jaql? Technical report. URL: <http://www-01.ibm.com/software/data/infosphere/hadoop/jaql/>.
- [10] Xquery 1.0: an xml query language. XML Query WG, December 2010. URL: <http://www.w3.org/TR/xquery/> (visited on 12/14/2010).
- [BCNS13] Véronique Benzaken, Giuseppe Castagna, Kim Nguyen, and Jérôme Siméon. Static and dynamic semantics of nosql languages. *Sigplan not.*, 48(1):101–114, January 2013. ISSN: 0362-1340. DOI: 10.1145/2480359.2429083. URL: <http://dx.doi.org/10.1145/2480359.2429083>.
- [BHNCS02] G. Bhalotia, A. Hulgeri, C. Nakhe, S. Chakrabarti, and S. Sudarshan. Keyword searching and browsing in databases using banks. In *Data engineering, 2002. proceedings. 18th international conference on*, 2002, pages 431–440. DOI: 10.1109/ICDE.2002.994756. URL: <http://dx.doi.org/10.1109/ICDE.2002.994756>.
- [Boy08] Brent Boyer. Robust java benchmarking, part 1: issues. June 2008. URL: <http://www.ibm.com/developerworks/java/library/j-benchmark1/index.html>.
- [CKKS02] W. F. Cody, J.T. Kreulen, V. Krishna, and W. S. Spangler. The integration of business intelligence and knowledge management. *Ibm systems journal*, 41(4):697–713, 2002. ISSN: 0018-8670. DOI: 10.1147/sj.414.0697. URL: <http://dx.doi.org/10.1147/sj.414.0697>.

- [Cod79] E. F. Codd. Extending the database relational model to capture more meaning. *Acm trans. database syst.*, 4(4):397–434, December 1979. ISSN: 0362-5915. DOI: 10.1145/320107.320109. URL: <http://dx.doi.org/10.1145/320107.320109>.
- [Col96] George Colliat. Olap, relational, and multidimensional database systems. *Sigmod record*, 25(3):64–69, 1996. DOI: 10.1145/234889.234901. URL: <http://dx.doi.org/10.1145/234889.234901>.
- [DFLDH88] S. T. Dumais, G. W. Furnas, T. K. Landauer, S. Deerwester, and R. Harshman. Using latent semantic analysis to improve access to textual information. In *Proceedings of the sigchi conference on human factors in computing systems*. In CHI '88. ACM, Washington, D.C., USA, 1988, pages 281–285. ISBN: 0-201-14237-6. DOI: 10.1145/57167.57214. URL: <http://dx.doi.org/10.1145/57167.57214>.
- [Dra] Richard Drake. Rdrake/molly. URL: <https://github.com/rdrake/Molly>.
- [ECM13] ECMA International. *Standard ecma-404 - the json data interchange format*, 1st edition, October 2013. URL: <http://www.ecma-international.org/publications/standards/Ecma-404.htm>.
- [Efr87] Bradley Efron. Better bootstrap confidence intervals. *Journal of the american statistical association*, 82(397):171–185, 1987. DOI: 10.1080/01621459.1987.10478410. URL: <http://dx.doi.org/10.2307/2289144>.
- [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.
- [HGP03] Vagelis Hristidis, Luis Gravano, and Yannis Papakonstantinou. Efficient ir-style keyword search over relational databases. In *Proceedings of the 29th international conference on very large data bases - volume 29*. In VLDB '03. VLDB Endowment, Berlin,

- Germany, 2003, pages 850–861. ISBN: 0-12-722442-4. URL: <http://dl.acm.org/citation.cfm?id=1315451.1315524>.
- [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. DOI: 10.1093/comjnl/32.2.98. URL: <http://dx.doi.org/10.1093/comjnl/32.2.98>.
- [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.
- [Jon72] Karen Spärck Jones. A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 28:11–21, 1972. DOI: 10.1108/eb026526. URL: <http://dx.doi.org/10.1108/eb026526>.
- [LJLF11] Guoliang Li, Shengyue Ji, Chen Li, and Jianhua Feng. Efficient fuzzy full-text type-ahead search. *The vldb journal*, 20(4):617–640, August 2011. ISSN: 1066-8888. DOI: 10.1007/s00778-011-0218-x. URL: <http://dx.doi.org/10.1007/s00778-011-0218-x>.
- [MRS08] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to information retrieval*. Cambridge University Press, New York, NY, USA, 2008. ISBN: 0521865719, 9780521865715.
- [Nar12] Ronen Narkis. Clojure concurrency. Slide 17. 2012. URL: <http://narkisr.github.io/clojure-concurrency/>.

- [RZ09] Stephen Robertson and Hugo Zaragoza. The probabilistic relevance framework: bm25 and beyond. *Found. trends inf. retr.*, 3(4):333–389, April 2009. ISSN: 1554-0669. DOI: 10.1561/1500000019. URL: <http://dx.doi.org/10.1561/1500000019>.
- [SB88] Gerard Salton and Christopher Buckley. Term-weighting approaches in automatic text retrieval. In *Information processing and management*, 1988, pages 513–523. DOI: 10.1016/0306-4573(88)90021-0. URL: [http://dx.doi.org/10.1016/0306-4573\(88\)90021-0](http://dx.doi.org/10.1016/0306-4573(88)90021-0).
- [Sch11] Silvia Schreier. Modeling restful applications. In *Proceedings of the second international workshop on restful design*. In WS-REST '11. ACM, Hyderabad, India, 2011, pages 15–21. ISBN: 978-1-4503-0623-2. DOI: 10.1145/1967428.1967434. URL: <http://dx.doi.org/10.1145/1967428.1967434>.
- [Suc98] Dan Suciu. An overview of semistructured data. *Sigact news*, 29(4):28–38, December 1998. ISSN: 0163-5700. DOI: 10.1145/306198.306204. URL: <http://dx.doi.org/10.1145/306198.306204>.
- [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. URL: <http://dx.doi.org/10.1145/286385.286387>.
- [Yeg08] Steven Yegge. The universal design pattern. October 2008. URL: <http://steve-yegge.blogspot.ca/2008/10/universal-design-pattern.html>.

To do...

- ☐ 1 (p. 4): last sentence doesn't flow
- ☐ 2 (p. 24): Above proof could use some love
- ☐ 3 (p. 26): Finish example
- ☐ 4 (p. 27): Talk about how great Clojure is.
- ☐ 5 (p. 30): More concurrency
- ☐ 6 (p. 30): Compare to other functional languages with decent library support
- ☐ 7 (p. 38): Expand on this. Why is GitHub significant?
- ☐ 8 (p. 81): Fix bibliography entries with short or no keys