# Towards a Concurrent Implementation of Keyword Search Over Relational Databases

## M.Sc. Thesis

### Richard Drake

richard.drake@uoit.ca

Supervisor: Dr. Ken Q. Pu

**University of Ontario Institute of Technology**
**Oshawa, Ontario, Canada**

9 June 2014

# Motivation

- Two important data models:
  - The relational model is rigid in structure and highly normalized
  - The document model is flexible and provides keyword search
- Choice between data integrity and accessibility
- Why can't we have both?

# Thesis Statement

*A system could be built that is capable of transforming data from the relational model to the document model. The transformation is reversible, allowing the original data model to be recovered. This system would use the keyword search capabilities, along with the relational information, to quickly discover related fragments of information.*

# Research Goals

▶ Define a formal framework for transforming data from the relational to document model

▶ Design a collection of expressive query operators for analyzing text from relational data sets

▶ Perform graph search over the document model

▶ Investigate implementation techniques to make the query operators performant on modern, multicore machines

# Relational Model

▶ A database is a collection of relations

▶ A relation is a set of named tuples

▶ Each named tuple consists of a set of attributes corresponding to values

▶ An entity group is a set of related tuples (joined by foreign key constraints)

# Document Model

- ▶ A document collection is a set of documents
- ▶ A document may have one or more fields
- ▶ A field is a bag of tokens
- ▶ A prescribed lexical analyzer converts bodies of text into a bag of tokens

# Framework Overview

▶ Indexing
  1. Iterate over all named tuples in relational database
  2. Convert each named tuple into a document
  3. Encode relational information into one or more indexing document(s)

▶ Search
  1. Fuzzy value search at the tuple-level
  2. Graph search for connections among entity groups

# Iterate

▶ Each entity and entity group defined in a configuration file

▶ Configuration specifies the SQL query to retrieve all entities from database

▶ This query is executed and each row is converted

```
(crawl
  [this db-conn idx-w]
  (let [sql (S :sql)]
    (execute-query db-conn sql (fn [row] ...))))
```

Figure: Code to iterate over every named tuple in entity group

# Convert

▶ Each attribute of a named tuple directly maps to a field in a document
▶ Analysis is performed on each attribute before storing in a field
▶ The framework supports multiple analyzers

$$\text{ATTR}[t] \xrightarrow{\text{analyzed}} \text{FIELD}[d] \qquad (1)$$

$$\alpha_1, \alpha_2, \ldots, \alpha_n \xrightarrow{\text{analyzed}} f_1, f_2, \ldots, f_n \qquad (2)$$

Where $\alpha_i$ is the value of an attribute, and $f_i$ are the resulting tokens in a document.

# Encode

▶ Relational information between related named tuples is encoded in the indexing document

▶ This document, $x$, is the concatenation of unique identifiers of every related document

$$x[\text{``entities''}] = \{\text{UID}[t] : t \in V(T)\} \tag{3}$$

Where $T$ is the entity group, and $V(T)$ is the set of named tuples in the entity group.

# Fuzzy Search at Tuple-Level

The encoding allows us to perform the following:

▶ Fuzzy search of relational attributes
▶ Search of named tuples in relations
▶ Find entity groups

# Graph Search Over Document Space

▶ Use the indexing document to discover adjacent nodes
▶ Issue search query to discover all entity groups containing UID[$u$]

```
(boolean-query [[(query :type :group) :and]
                [(query :entities id) :and]])
```

# Graph Search Algorithm and Implementation

▶ Chose Breadth-First Search (BFS) as search algorithm

▶ Adapted BFS to run concurrently without locking using Clojure's Software Transactional Memory (STM)

▶ Utilized Clojure's concurrency primitives (atoms, references) in the concurrent implementation of BFS

# Data Corpus

▶ Derived from MyCampus data[1]

| Relation | Attributes |
|----------|-----------|
| Course | <u>code</u>, title, subject |
| Section | <u>id</u>, actual, campus, capacity, credits, levels, registration_start, registration_end, semester, sec_code, sec_number, year, course |
| Schedule | <u>id</u>, date_start, date_end, day, schedtype, hour_start, hour_end, min_start, min_end, classtype, location, section_id |
| Instructor | <u>id</u>, name |
| Teaches | <u>id</u>, schedule_id, instructor_id, position |

Table: Subset of mycampus dataset schema

---

[1] http://uoit.ca/mycampus/

# Benchmarking Queries

▶ Keywords are sampled from the database to form search queries
  ▶ May not occur within the same entity group
  ▶ Sampled from instructor name, course code, section id, etc.

▶ Search for connections between these keywords
  ▶ Must utilize graph search to discover connections between keywords
  ▶ Monitor distance (hops) between keywords and measure performance

# Results



Figure: Growth of runtime of each implementation, by hops

# Results



Figure: Runtime of implementations, 1 Hop

# Results



Figure: Runtime of implementations, 4 Hops

# Results



Figure: Runtime of implementations, 6 Hops

# Contributions

▶ Provided a framework to transform data from the relational to document model

▶ Demonstrated the reversibility of this transformation

▶ Utilized the flexibility of the document model
  ▶ e.g. spelling correction, entity group search

▶ Performed graph search over document space

▶ Investigated the reduction in runtime from a concurrent graph search

# Lessons Learned

- Simple algorithms are easiest to parallelize
- Clojure's STM implementation is simple and effective
- Clojure is powerful and encouraged correct code that was easier to optimize later

# Publication

This work has been submitted to the 15th IEEE conference on Information Reuse and Integration.

# Future Work

- Simplify configuration
- Incremental indexing
- Generalize results by benchmarking standard datasets

# Towards a Concurrent Implementation of Keyword Search Over Relational Databases

## M.Sc. Thesis

### Richard Drake

richard.drake@uoit.ca

Supervisor: Dr. Ken Q. Pu

**University of Ontario Institute of Technology**
**Oshawa, Ontario, Canada**

9 June 2014