

Citation Network Proposal

Group 2

MITIU14005 - Nguyen Hoang Minh

MITIU14006 - Chiem Thach Phat

Ho Chi Minh City, International University

January 2, 2015

Table of Contents

- 1 Introduction
- 2 Implementation
- 3 Results
- 4 Lesson Learned
- 5 Future works

Table of Contents

- 1 Introduction
- 2 Implementation
- 3 Results
- 4 Lesson Learned
- 5 Future works

The world needs a citation network for patents to:

- Analyzing trending in field
- Looking up for the existence of similar patents
- Finding who invents most patents

Table of Contents

- 1 Introduction
- 2 Implementation**
- 3 Results
- 4 Lesson Learned
- 5 Future works

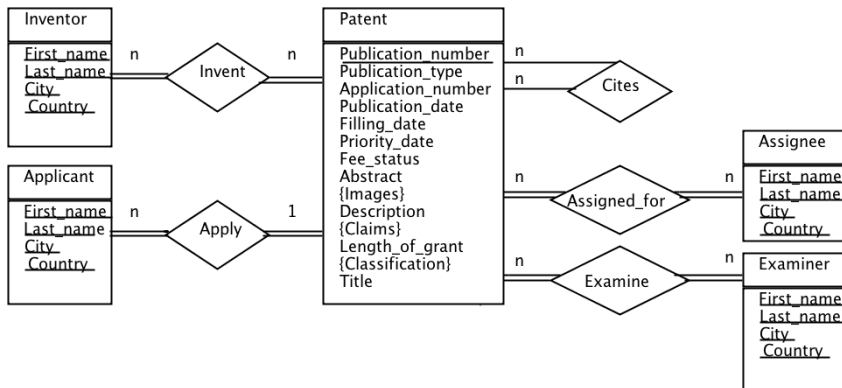
Analyzing the Scenario

We would like to build a citation network for patents with these constraints:

- Able to store and analyze patents from USPTO bulk data.
- A patent contains: an unique Publication Number, Publication Type, Application Number, Publication Date, Filing Date, Priority Date, Fee Status, Title, Length of Grant, Classification and some of its contents include Abstract, Images, Description and Claims
- A Patent may or may not cite/be cited other patents
- A Patent must be invented/applied/examined by or assigned for at least one Inventor, Applicant, Examiner or Assignee
- Inventor, Applicant, Examiner and Assignee are identified by first name, last name, city and country

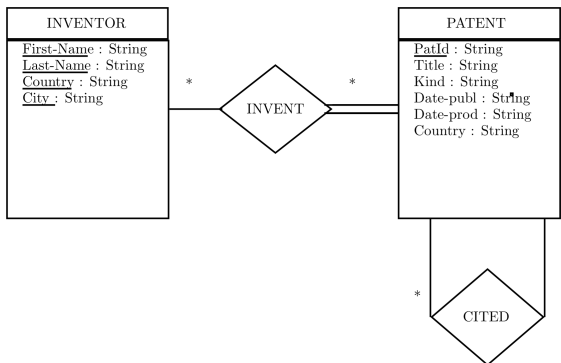
Conceptual Model

Based on the scenario, we have designed this conceptual model.



Implemented Conceptual Model

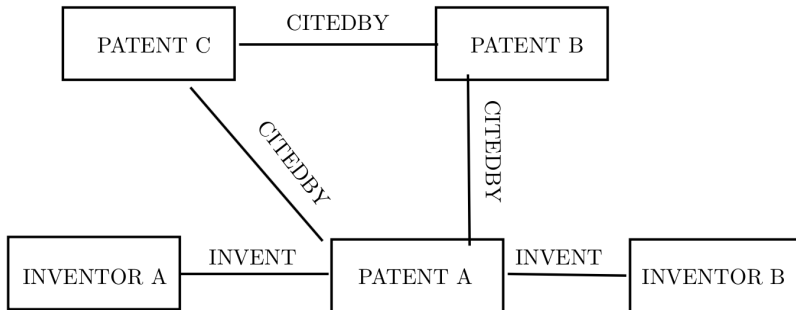
Due to time constraint, we use this conceptual model to implement the citation network.



Logical Model - Neo4j

Neo4j use the concept of **Nodes**, **Properties** and **Relationships** as shown in the picture below, where:

- Nodes: PATENT, INVENTOR
- Properties: Nodes and Relationships attributes, such as Patent_No or Cited_date
- Relationships: CITEDBY, INVENT



Logical Model - Neo4j (Cont.)

Pros:

- Can cope with complex relationship
- Offers more accurate ways to determine the relationship between each entity

Cons:

- Slow with complex search
- Built-in Web Server is not supported for visualizing more than 1000 rows

Logical Model - MongoDB

While Neo4j stores data in Nodes, Properties and Relationships, **MongoDB stores data as a BSON file.**

```
1 patent_schema = {
2     kind : string ,
3     date-published : string ,
4     number-of-claims : integer ,
5     country : string ,
6     patid : string ,
7     main-classification : string ,
8     title : string ,
9     app-number : string ,
10    citations : [ { patid : string } , {
11                  patid : string } ]
12    date-produced : string ,
13    inventors:[
14        { city : string , last-name :
15          string , country : string , first-name :
16          string } ,
17        { city: string , last-name : string ,
18          country : string , first-name : string }
19    ]
20 }
21
22 inventor_schema={
23     city : string ,
24     last-name : string ,
25     country : string ,
26     first-name : string
27 }
```

Logical Model - MongoDB (Cont.)

Pros:

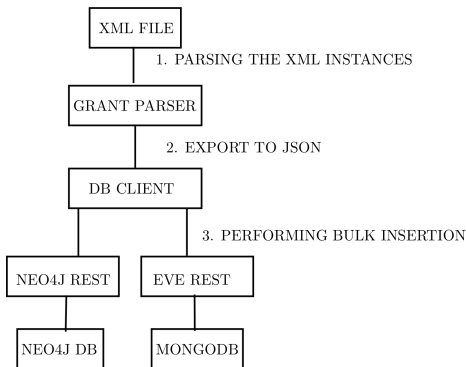
- Documents (i.e. objects) can be defined with native data types that can be seen in many programming languages. E.x: Date, String, Integer
- Embedded documents and arrays reduce need for expensive joins
- Dynamic schema supports fluent polymorphism

Cons:

- Does not support relationship

CINAS Structure

Our CINAS will follow the structure below, starts from the Patent XML File and ends at **Neo4J DB and MongoDB** as our data storages.



Grant Parser Module

Due to the enormous size of the Patent Bulk Download File, we will use **SAX method** instead of DOM method to **process the XML file**.

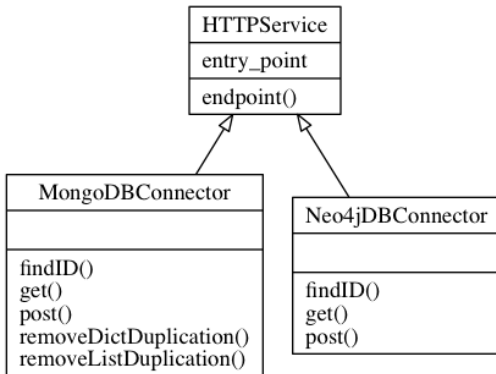
Pseudocode for **Grant Parser Module**:

- ① Load XML file and break it down into separated XML instances
- ② Parse each XML instance into patent handler. For each tag with nested tags, enable stack to store adjacent data
- ③ Tags that are explicitly defines in ignore list will be left out during parsing. Upon completion, reset the patent handler for new instance
- ④ Finally, an option to export the entire file to JSON format is available if chosen

DB Client Module

We chose to deploy REST services in communication with our backends instead of OM to provide more scalable solution as we just use HTTP response and request.

DB Client Class Diagram:



Pseudocode for **POST in MongoDBConnector**:

- 1 Load all data from JSON file
- 2 Construct list of citations from group of patent
- 3 Validate uniqueness within the list
- 4 Appending statements into query. Then, perform 1st POST these citations
- 5 If the uniqueness is violated, the response code from EVE server will determine the next step. If error response with duplicates found, remove it from the list and resend it again
- 6 Finally, perform an overwrite PUT for any individual patent. This step must be linearly execute to guarantee that all the dummy patents, which are created using citations, will be updated with information

POST in Neo4j is more simpler because it has an overwrite operation called **“MERGE”**, which **automatically detect duplication and merge together or create new one.**

Pseudocode for **POST in Neo4jDBConnector:**

- ① Load all data from JSON file
- ② Construct list of citations from group of patent
- ③ Appending statements into query. Then, perform POST these citations
- ④ Finally, we perform another POST for list of patents

Citation Network Visualization

For Citation Network Visualization, we wrote **generate_csv** to automate the process:

- 1 Assuming that Neo4J is populated using **DB Client**
- 2 Run **generate_csv** will results in 2 csv files : nodes.csv and links.csv
- 3 Use GEPHI to load these files and visualize the data

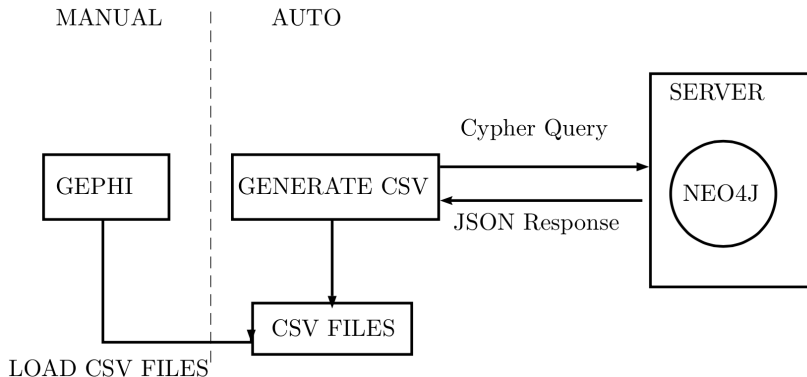


Table of Contents

- 1 Introduction
- 2 Implementation
- 3 Results**
- 4 Lesson Learned
- 5 Future works

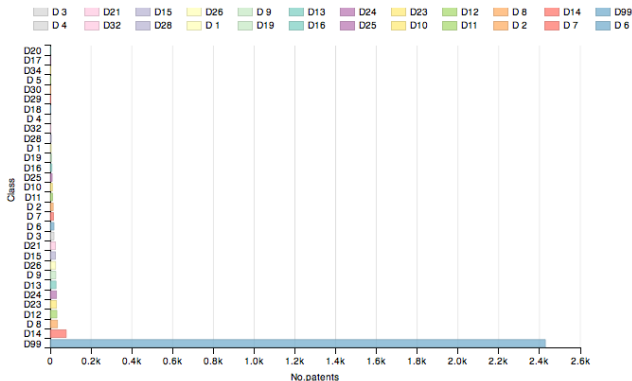
Processed patents:

	Number of Documents
Duplicated Citations	347
Patents	3000
Citations	106996
Total	110343

Results (Cont.)

The figure illustrates Patent Classification in October 2014.

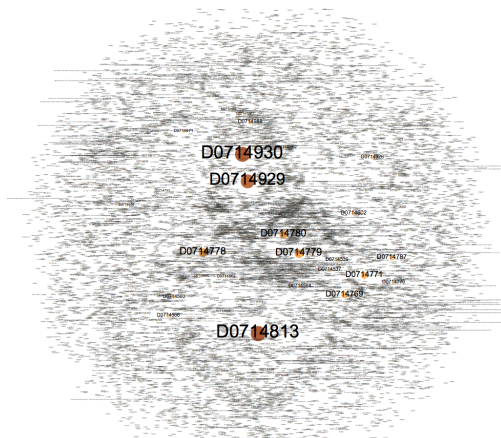
The MISCELLANEOUS (D99) industry were blossomed in October 2014



Citation Network

The figure shows a citation network with 500 patents in the database.

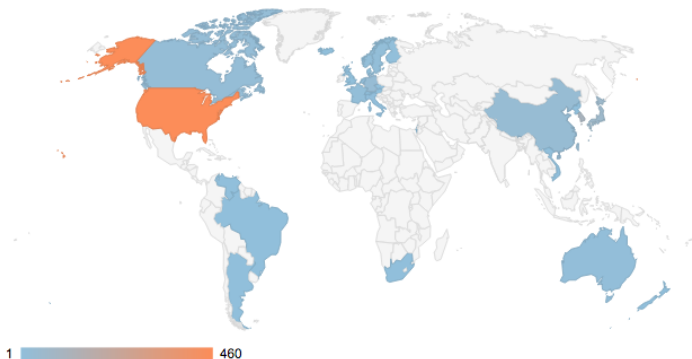
D0714930, D0714929 and D0714813 were the most popular in rank based on connected components.



Inventors Maps

To measure the activity of inventors community, we classify them based on their regions.

Most of inventors come from North America. Furthermore, US patents industry also attracts the numerous inventors from different parts of the world, including Asian and Australian.



Benchmark Neo4j and MongoDB

We perform the time measurements on 500 patents with total of 12016 citations.

Results:

- **DOM** parser is **faster** than **SAX**
- **Object-Mapping** is **faster** than **REST**, but **REST** is **easier to extend** than Object-Mapping
- Populate data in **MongoDB** is **faster** than **Neo4j**, but the drawback is MongoDB client has to implement references between documents

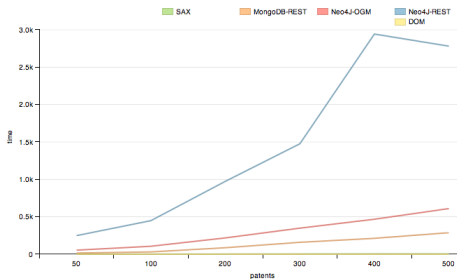


Table of Contents

- 1 Introduction
- 2 Implementation
- 3 Results
- 4 Lesson Learned**
- 5 Future works

- Choice of technology
 - Analyzing citation network with Neo4J is much simpler than MongoDB & RDBMS.
 - Neo4J web server only visualizes up to 1000 rows.
 - Object-Mapping is faster and easier to use than REST. However, REST is more suitable for scalable system that involves different languages & platforms that communicate via JSON.
 - Drawbacks of Normalization in RDBMS since it can't fully represent the natural language . Hence, we chose MongoDB to store data.
- Patent
 - It's inevitable that patent can help identify the trend of inventions in one country's industry.
 - Further analysis may take into account the backward citation and forward citation to identify the trend on time basis.

Table of Contents

- 1 Introduction
- 2 Implementation
- 3 Results
- 4 Lesson Learned
- 5 Future works**

- Divide patent data into different databases instead of storing the whole file into one particular back-end
- Propose a refined structure for storing and querying patent data across multiple database systems
- Provide an end-user services that would transform natural language into our CINAS's language, process the queries, determine the keywords and then pass it to the suitable database system
- Final goals:
 - ① Take advantages of multiple database systems
 - ② Reduce the query costs
 - ③ Enhance end-users experience with patent data

- [1] “Transitioning from Relational Databases to MongoDB - Data Models” B. Reiner, <http://blog.mongodb.org/post/72874267152/transitioning-from-relational-databases-to-mongodb>, Jan 2014
- [2] “USPTO Patents Database Construction” I. Torvik, D. Dubin and Q.Liu, <http://abel.lis.illinois.edu/UPDC/USPTOPatentsDatabaseConstruction.pdf>, August 2012
- [3] “Construction of Japanese Patent Database for Research on Japanese patenting activities” A. Goto, K. Motohashi, <http://www.iip.or.jp/e/images/paper.pdf>

Thank You

Thank you for your listening!