

Departmental Coversheet  
Hillary term 2022 Mini-project  
Paper title: Database System Implementation  
Candidate Number: 1058016  
Your degree: MSc Advanced Computer Science

To run the code, try make all. There will be an output file RDF.DB.out. Run the file and you can enter the database interface.

1. (a) The data structure in this paper is a six-column triple table. The subject, predicate, and object of a triple are encoded as integers in the first three columns,  $R_s$ ,  $R_p$ , and  $R_o$ , respectively. Conceptually, there are three linked lists, an  $sp$ -list that connects all triples with the same  $R_s$  grouped by  $R_p$ , an  $op$ -list that associates all triple with the same  $R_o$  grouped by  $R_p$ , and a  $p$ -list that relates all triples with the same  $R_p$  without any grouping. In the table, the last three columns,  $N_{sp}$ ,  $N_{op}$ , and  $N_p$ , store the next-pointes, which are going to be the row numbers in the triple table in the actual implementation.

Apart from the table, six index maps are also maintained.  $I_s$ ,  $I_p$ , and  $I_o$  store the head of the  $s$ -list,  $p$ -list, and  $o$ -list for each different  $s$ ,  $p$ , and  $o$ , respectively.  $I_{sp}$  maps to the head of the  $sp$ -list for each  $(s, p)$  tuple. So does  $I_{op}$ .  $I_{spo}$  stores the row number of each triple in the table and we use this to check existence of a triple in the database.

- (b) ADD(Triple  $t$ ) consists of two parts; first, append a new row to the end of the RDF-index table, and second, update all six index maps and alters the next-pointer columns,  $N_{sp}$ ,  $N_{op}$ , and  $N_p$  accordingly.

Since there is no need to worry about the concurrency in our setting, the RDF-index table is maintained as a fixed size vector of int list of size 6 with each position stands for different columns. Therefore, for each time we add a triple into the table, we simply append it to the last. However, if a triple is already in the  $I_{spo}$  map, we skip it. When the size reaches its limit, the entire table will resize.

---

**Algorithm 1** ADD ( $t$ )

---

```

if  $t$  in  $I_{spo}$  then                                     ▷  $t = (s, p, o)$  is a triple.
    return
end if
 $i = \#$  of elements in the triple-table
if  $i + 1 >$  the size of the triple-table then
    resize the triple table
end if
 $T_{new} = [t.s, t.p, t.o, -1, -1, -1]$                      ▷ The last three columns are left for update later.
triple-table[ $i$ ] =  $T_{new}$ 
 $I_{spo}[t] = i$ 
Update remaining indexes

```

---

The columns,  $N_{sp}$ ,  $N_{op}$ , and  $N_p$ , are maintained in a linked-list-like manner and are updated simultaneously with the index maps. The central idea is that we always append the new triple right after the head of the list, but there are two corner cases we need to discuss.

Take  $N_{op}$  for example, if a new triple,  $T_{new}$ , does not appear in  $I_o$  and  $I_{op}$ , it means that  $T_{new}$  is a triple with a brand new  $o$  and  $p$ . Therefore, we make  $T_{new}$  the head of the index map  $I_o$  and  $I_{op}$  (case1). If we found a  $T$  with  $T.o = T_{new}.o$  and  $T.p = T_{new}.p$ , then we insert  $T_{new}$  after  $T$ , and point the next of  $T_{new}$  to the original next of the  $T$  as described in the paper (case2). A special case is that when there is not a head of  $I_{op}$  for  $T$ , we will make  $T$  the new head of  $I_o$  and  $T.N_{op}$  point to the original start of  $I_{op}$  and  $I_o$ <sup>1</sup>, which is handled in case 3 of the Algorithm (2). The case of  $N_{sp}$  and  $N_p$  is similar.

---

<sup>1</sup>When encountering a new value, we always want to make the new value the head of the index map instead of appending it to the end, because appending to the end takes a lot of time in traversing the linked list.

---

**Algorithm 2** Update  $I_{op}(T_{new})$ 

---

```
 $T$  = the first triple with  $T.o = T_{new}.o$  and  $T.p = T_{new}.p$   
if  $T$  does not exist then ▷ Case 1  
    make  $T_{new}$  the head of  $I_o$  and  $I_{op}$   
end if  
if  $I_{op}$  has  $T$  then ▷ Case 2  
     $T_{new}.N_{op} = T_{next}$   
     $T.N_{op} = T_{new}$   
else if  $I_{op}$  does not have  $T$  then ▷ Case 3  
    make  $T_{new}$  the head of  $I_{op}$   
     $T_{new}.N_{op} = I_o.head$   
    make  $T_{new}$  the head of  $I_o$   
end if
```

---

- (c) Let  $t = \langle s, p, o \rangle$  be a triple pattern and  $X, Y$ , and  $Z$  be free variables. There are two inputs for EVALUATE;  $t$ , the matching pattern, and  $index$ , the previous row number. There are also two outputs,  $t'$ , the answer to the matched pattern,  $index$ , the current row number. The index are usually the pointer (row number) to the next triple, but there are two sentinels; *EndOfNode* will tell the caller there is no more next triple to search and *EndSearch* stands for an illegal search. The general framework is shown in Algorithm (3).

---

**Algorithm 3** Framework

---

```
while  $index \neq EndOfNode$  and  $index \neq EndSearch$  do  
     $index, t' \leftarrow Evaluate(t, index)$   
end while
```

---

In the while loop, it will return matched pattern one by one until *EndOfNode* or *EndSearch* is returned by Evaluate(.).

Inside EVALUATE, there are four categories, from no free variable to three variables.

When there is not a free variable, EVALUATE will check if  $t$  is in  $I_{spo}$  map, if it is return the triple and set the index to be *EndOfNode*. Otherwise, it will return nothing and set the index be *EndSearch*. The pseudo-code is shown in Algorithm (4).

---

**Algorithm 4** Evaluate  $\langle s, p, o \rangle$ 

---

```
Require:  $t = \langle s, p, o \rangle$   
if  $t$  in  $I_{spo}$  then  
     $index \leftarrow EndOfNode, t' \leftarrow t$   
else  
     $index \leftarrow EndSearch, t' \leftarrow Null$   
end if
```

---

When there is one free variable, there are three cases:  $\langle X, p, o \rangle$ ,  $\langle s, X, o \rangle$ , and  $\langle s, p, X \rangle$ . The first and the last cases are similar. For  $\langle X, p, o \rangle$ , in the first search, we use the index map,  $I_{op}$ , to locate the head of the triple with predicate and object the same as  $o$  and  $p$ , then we return the first match and the next pointer at  $N_{op}$ . For all subsequent evaluates, we traverse  $N_{op}$  list until there is not a match. The pseudo-code is suggested in Algorithm (5).  $\langle s, p, X \rangle$  is similar.

---

**Algorithm 5** Evaluate  $\langle X, p, o \rangle$ 

---

**Require:**  $t = \langle X, p, o \rangle$   
**if**  $\text{hash}(t.p, t.o)$  does not appear in  $I_{op}$  **then**  
     $\text{index} \leftarrow \text{EndSearch}$   
**end if**  
**if** it is the first search **then** ▷ locate the head  
     $t' \leftarrow$  the triple  $I_{op}$  points to,  $\text{index} \leftarrow t'.N_{op}$ .  
**else** ▷ traverse  $N_{op}$  list  
     $t' \leftarrow t.N_{op}$ ,  $\text{index} \leftarrow t'.N_{op}$   
**end if**  
**if**  $t'.o$  and  $t'.p$  are not we are looking for **then** ▷ the end of grouped by.  
     $\text{index} \leftarrow \text{EndSearch}$   
**end if**

---

Since we don't have an index on  $s$  and  $o$ ,  $\langle s, X, o \rangle$  needs some extra works. To match this pattern, we can traverse either  $sp$ -list or  $op$ -list and skip the triples does not match on  $s$  and  $o$ . The choose of the list depends on the size of  $I_s[s]$  and  $I_o[o]$  and we use the smaller one. The pseudo-code is in Algorithm (6).

---

**Algorithm 6** Evaluate  $\langle s, X, o \rangle$ 

---

**Require:**  $t = \langle s, X, o \rangle$ , where  $X$  is a special code stands for a free variable.  
**if**  $|I_s[s]| < |I_o[o]|$  **then**  
    **do**  
        **if** it is the first search **then**  
             $t' \leftarrow$  the triple  $I_{sp}$  points to,  $\text{index} \leftarrow t'.N_{sp}$ .  
        **else**  
             $t' \leftarrow t.N_{sp}$ ,  $\text{index} \leftarrow t'.N_{sp}$   
        **end if**  
        **while**  $t'.s \neq s$  and  $t'.o \neq o$  and  $\text{index} \neq \text{EndOfNode}$   
    **else**  
        **do**  
            **if** it is the first search **then**  
                 $t' \leftarrow$  the triple  $I_{op}$  points to,  $\text{index} \leftarrow t'.N_{op}$ .  
            **else**  
                 $t' \leftarrow t.N_{op}$ ,  $\text{index} \leftarrow t'.N_{op}$   
            **end if**  
            **while**  $t'.s \neq s$  and  $t'.o \neq o$  and  $\text{index} \neq \text{EndOfNode}$   
        **end if**  
        **if**  $t'.s \neq s$  and  $t'.o \neq o$  **then** ▷ corner case  
             $\text{index} = \text{EndSearch}$   
        **end if**

---

When there are two free variables, there are also three cases:  $\langle X, Y, o \rangle$ ,  $\langle X, p, Z \rangle$ , and  $\langle s, Y, Z \rangle$ , where  $X$ ,  $Y$ , and  $Z$  can be equal. Their ideas are similar;  $\langle X, Y, o \rangle$ , for example, we first find  $I_o[t.o]$  and traverse the  $op$ -list till the end. If  $X = Y$ , then we skip those triples with  $p \neq o$ . The pseudo-code is given in Algorithm (7)<sup>2</sup>.

---

<sup>2</sup>Even I found myself a little confused by what I wrote in the do-while loop, so this footnote might be helpful. After a few minutes of thinking, I found only when we want  $X = Y$ , the condition  $t'.p \neq t'.s$  will make the loop start and try to locate the one with  $t.p = t.s$ . If  $X \neq Y$ , the loop will never start and return immediately.

---

**Algorithm 7** Evaluate  $\langle X, Y, o \rangle$ 

---

**Require:**  $t = \langle X, Y, o \rangle$ , where  $X$  and  $Y$  are special codes stand for variables.

```
do
  if it is the first search then
     $t' \leftarrow$  the triple  $I_o$  points to,  $index \leftarrow t'.N_{op}$ .
  else
     $t' \leftarrow t.N_{op}$ ,  $index \leftarrow t'.N_{op}$ 
  end if
  while  $X = Y$  and  $t'.p \neq t'.s$  and  $index \neq EndOfNode$ 
  if  $X = Y$  and  $t'.p \neq t'.s$  then  $\triangleright$  stop traversing when  $t'.p \neq t'.s$  when the pattern asks  $X = Y$ 
     $index = EndSearch$ 
  end if
```

---

The ideas for the other two patterns are similar; start from  $I_p$  or  $I_s$  map and traverse the  $p$ -list or  $sp$ -list.

When there are three free variables, the paper suggests to match, for example, the patterns like  $\langle X, Y, Z \rangle$ , we need to iterate over the triple table; if we want  $X = Y$ , we skip those  $X \neq Y$ . However, this is not efficient. Therefore, I modified a little bit and the idea and pseudo-code are discussed in 2.a, where I put everything that is different from the paper there.

2. (a) RDF indexing data structure that implements Add and Evaluate functions.

This component is included in `RDF_index.cpp`, in which ADD and EVALUATE are implemented as suggested in the problem 1. There are a few things that is slightly different from the paper.

1) I used XXHASH<sup>3</sup> by Facebook instead of Jenkins hashing, because it achieves state-of-the-art excellent performance on both long and small inputs. Jenkins hashing is implemented in `RDF_index.h` and but it is commented. I did some tests and observed that Jenkins is about 4% slower than XXHASH in average<sup>4</sup>.

2) Instead of open addressing, I eventually choose `std::unordered_map` for index maps  $I_{sp}$ ,  $I_{op}$  and  $I_{spo}$ . I had an open addressing hash implemented in `HashTable.cpp` and `HashTable.h` (attached in the submission), but it was much slower than the `unordered_map` by about 20% or more.

3) To match the patterns like  $\langle X, Y, Z \rangle$ , the paper suggests to iterate over the triple table; if we want  $X = Y$ , we skip those  $X \neq Y$ . I modify this a little bit to improve the efficiency. Again,  $\langle X, X, Z \rangle$ , for example, we first iterate  $I_s$ , and for each  $s$  in  $I_s$ , we find if  $I_{sp}$  includes `hash(s, s)`, if yes, we traverse over the triple table. As shown in the pseudo-code in Algorithm (8)<sup>5</sup>.

---

**Algorithm 8** Evaluate  $\langle X, X, Z \rangle$ 

---

```
for  $s$  in  $I_s$  do
   $i = \text{hash}(s, s)$ 
  if  $i$  in  $I_{sp}$  then
    Evaluate_SPZ( $s, s$ )
  end if
end for
```

---

$\langle X, Y, Y \rangle$  and  $\langle X, Y, X \rangle$  are similar. For  $\langle X, X, X \rangle$ , we only need to traverse  $s$  in  $I_s$  and find if  $\langle s, s, s \rangle$  is in  $I_{spo}$  as shown in Algorithm (9).

---

<sup>3</sup><https://cyan4973.github.io/xxHash/>

<sup>4</sup>For loading, XXHASH and Jenkins took 314.359, 4026.82, and 323.914, 4206.85, and 58582.4 ms on three datasets, respectively, achieving 2%, 4%, and 8% speed up on each.

<sup>5</sup>The actual implementation is slightly different. I put the check of if  $i$  in  $I_{sp}$  at the beginning of the Evaluate\_SPZ to reduce some code redundancy, but they are equivalent.

---

**Algorithm 9** Evaluate  $\langle X, X, X \rangle$ 

---

```
for  $s$  in  $I_s$  do  
    Evaluate_SPO( $s, s, s$ )  
end for
```

---

- (b) The engine for evaluating BGP SPARQL queries.

This part is in the file SPARQL\_engine.cpp and it is implemented strictly based on the model answer. A few additional lines are added for printing or improving performance.

- (c) The greedy join order optimization query planner.

This part can be sought in the file query\_planner.cpp. I see what the model answer is trying to say, but I believe there are some minor mistakes, so I made some slight changes, but it is still  $O(n^2)$ , where  $n$  is the number of triple patterns we try to fit. There were four things I modified.

Firstly, the last three lines should go to the outside while loop.

Secondly, I changed the criteria for updating the new triple patterns. In the model answer, the criteria is

$$t_{best} = \perp \text{ or } score < score_{best} \text{ and either } var(t) = \emptyset \text{ or } var(t) \cap B \neq \emptyset. \quad (1)$$

It is troublesome because it will always take the first unprocessed triple pattern as the best pattern and compare this with the remaining. However, it might cause some issues in many cases. For example, the following triple patterns

$$\begin{aligned} \langle X, 1, 2 \rangle \\ \langle Y, 2, 3 \rangle \\ \langle X, 4, Y \rangle \end{aligned}$$

will produce the plan  $\langle X, 1, 2 \rangle \mapsto \langle Y, 2, 3 \rangle \mapsto \langle X, 4, Y \rangle$  because even after we replace the  $X$  in  $\langle X, 4, Y \rangle$ , this triple pattern still has a higher selectivity than  $\langle Y, 2, 3 \rangle$ . Therefore, this plan will create an unnecessary cross product, and I believe the criteria (1) meant choosing the triple pattern that has the lowest selectivity and includes some variables that are in the processed patterns. Thus, I made some modifications to the algorithm, and it will 1) as long as an unprocessed pattern has a variable that is processed, we will always consider this pattern and 2) if this pattern has lower selectivity than the previous one, we choose this pattern as the best triple. With this modification, the plan becomes  $\langle X, 1, 2 \rangle \mapsto \langle X, 4, Y \rangle \mapsto \langle Y, 2, 3 \rangle$ .

Thirdly, there are some cases that  $|\langle X, 1, 2 \rangle| = 10,000$  but  $|\langle Y, 2, 3 \rangle| = 10$ . If we ignore the size of the pattern, we might produce a lot of unnecessary joins, but we can store this information while we are updating  $I_{op}$ . Therefore, I also record the size of  $I_{op}[op]$  and  $I_{sp}[sp]$ , and we choose the best triple pattern based on both selectivity and size. Our final query plan is  $\langle Y, 2, 3 \rangle \mapsto \langle X, 4, Y \rangle \mapsto \langle X, 1, 2 \rangle$ , which is the most optimal plan. Notice that only  $I_{op}$  and  $I_{sp}$  are considered because they are the most common cases. Moreover, since the size of  $I_{op}[t]$  and  $I_{sp}[t]$  can be recored while we are inserting triple, the cost of maintaining the size is negligible<sup>6</sup>.

Lastly, the original query plan was based on the following precedence relation  $\prec$  on selectivity of triple patterns:

$$(s, p, o) \prec (s, ?, o) \prec (?, p, o) \prec (s, p, ?) \prec (?, ?, o) \prec (s, ?, ?) \prec (?, p, ?) \prec (?, ?, ?).$$

Now, since we have maintained the size of  $(?, p, o)$  and  $(s, p, ?)$ , a new order of the selectivity should be used:

$$(s, p, o) \prec (s, ?, o) \prec (?, p, o) = (s, p, ?) \prec (?, ?, o) \prec (s, ?, ?) \prec (?, p, ?) \prec (?, ?, ?).$$

---

<sup>6</sup>There are many tricks to save memory. For example, let's say  $I_{sp}$  is an unordered\_map<int, int>, where the first int is the hash key and the second int is the index of row number. For a table size smaller than  $2^{24}$  rows, we can use the higher 7 bits to store the size and lower 24 bits to store the index. Other trick might apply also.

Therefore, when we are precessing these two patterns, instead of comparing the selectivity, we can directly compare the actual size of the triple patterns and select the smaller one. For example, for the following triple patterns,

$$\begin{aligned} &\langle X, 1, 2 \rangle \\ &\langle 2, 3, Y \rangle \end{aligned}$$

with  $|\langle 2, 3, Y \rangle| > |\langle X, 1, 2 \rangle|$  will produce the following query plan,

$$\begin{aligned} &\langle 2, 3, Y \rangle \\ &\langle X, 1, 2 \rangle \end{aligned}$$

The pseudo-code goes as the following.

---

**Algorithm 10** New-Plan-Query( $U$ )

---

```

 $P \leftarrow [], B \leftarrow \emptyset$ 
while  $U \neq \emptyset$  do
   $t_{best} \leftarrow \perp, score_{best} \leftarrow 100, intersected \leftarrow \text{false}$ 
  for each unprocessed tripple pattern  $t \in U$  do
     $score \leftarrow$  the position of  $t$  in  $\prec$  where the variables in  $B$  are considered bounded
    if  $t_{best} = \perp$  then
       $intersected = (var(t) \cap B \neq \emptyset)$ 
       $t_{best} \leftarrow t, score_{best} \leftarrow score$ 
    end if
    if  $t$  and  $t_{best}$  are of the form  $\langle X, p, o \rangle$  or  $\langle s, p, X \rangle$  then
       $t_{best} \leftarrow$  the one with a smaller size
    end if
    if  $intersected$  then
      if  $score < score_{best}$  and either  $var(t) = \emptyset$  or  $var(t) \cap B \neq \emptyset$  then
         $t_{best} \leftarrow t, score_{best} \leftarrow score$ 
      end if
    else
      if  $var(t) = \emptyset$  or  $var(t) \cap B \neq \emptyset$  then  $\triangleright$  make sure that the next pattern share some variables.
         $t_{best} \leftarrow t, score_{best} \leftarrow score, intersected \leftarrow \text{true}$ 
      end if
    end if
  end for
end while

```

---

From Table 1, we see a significant improvement in all queries. The testing protocol is the same as the one in question 3. To do the test, uncomment lines 20-25 in query\_planner.cpp and comment lines 26-53 and MAKE TEST. The analysis is included in 3.c.

- (d) The component for parsing and importing Turtle files.  
This part is included in Turtle\_handler.cpp. Some corner cases are handled and it achieves what the question asks.
- (e) The parser for SPARQL queries  
This part is included in query\_parser.cpp. Some corner cases are handled and it achieves what the question asks.
- (f) The component implementing the command line  
This part is included in interface.cpp.

	Model Answer (ms)	New Query Plan (ms)	Improvement (%)
q1	0.2269	0.0038	5971.05263
q2	370.891	0.8707	42596.8761
q3	1.289	0.0022	58590.9091
q4	0.1669	0.0241	692.53112
q5	2.2444	0.06	3740.66667
q6	0.3	0.2121	141.442716
q7	2716.41	0.0106	25626509.4
q8	28.3936	1.6634	1706.96164
q9	1033.41	1.1611	89002.6699
q10	2.3399	0.0018	129994.444
q11	0.0617	0.0205	300.97561
q12	0.0517	0.0105	492.380952
q13	1.3423	0.0018	74572.2222
q14	0.2384	0.1947	122.444787

Table 1: Time took to process and evaluate the query by the model answer and the new query plan in ms on the dataset LUBM-001-mat.ttl. The new plan achieves more than 1,859,602.501% speed up in average on the 14 queries.

3.a (a) Hardware and Software Configuration:

- i. Model Identifier: MacBookPro14,3
- ii. Processor Name: Apple M1 Pro
- iii. Processor Speed: 3.2 GHz
- iv. Number of Processors: 1
- v. Total Number of Cores: 10
- vi. Memory: 32 GB
- vii. Operating System: macOS Monterey, Version 12.3
- viii. Compiler Version: Apple clang version 13.1.6 (clang-1316.0.21.2)

- (b) We set a timeout if one instruction takes more than 3 minutes. The following protocol is implemented in experiment.cpp. If you want to run the load test, fill in the correct path to the data. It will take about 5-10 minutes to run.

Loading Data Test Protocol:

- i. \*Restart computer\* and turn off all irrelevant applications.
- ii. Start a clean terminal.
- iii. Load dataset 001 and markdown the time. If timeout, mark TO.
- iv. Delete all objects.
- v. Load dataset 010 and markdown the time. If timeout, mark TO.
- vi. Delete all objects.
- vii. Load dataset 100 and markdown the time. If timeout, mark TO.
- viii. Delete all objects.
- ix. Go back to step ii repeat this process for 10 times.
- x. Take the average.

To eliminate the cache effect, I load 3 datasets one by one for 10 times instead of running a single dataset for 10 times. This will flush most of caches. Also, I should play my phone and don't touch my laptop while running. The protocol is implemented the function load\_test in experiment.cpp.



- (c) We choose COUNT instead of SELECT and for each query. We set a timeout if one instruction takes more than 3 minutes.

Loading Data Test Protocol:

- i. \*Restart computer\* and turn off all irrelevant applications.
- ii. Start a clean terminal.
- iii. Load one dataset.
- iv. Run 14 queries one by one and markdown the time.
- v. Go back to the step v and repeat 10 times.
- vi. Take the average.

To avoid the cache effect, I run 14 queries one by one for 10 times instead of running a single queries for 10 times. The protocol is implemented the function `query_test` in `experiment.cpp`.

- 3.b The time needed to load and index the data and the time needed to produce all query answers for each RDF graph and query are showed in Table 2 and Table 3, respectively<sup>7</sup>.

	LUBM-001-mat (ms)	LUBM-010-mat (ms)	LUBM-100-mat (ms)
Load & Index Time	114.651	2058.07	25624.2

Table 2: Time needed to load and index the data in ms.

	LUBM-001-mat (ms)	LUBM-010-mat (ms)	LUBM-100-mat (ms)
q1	0.0038	0.0063	0.0371
q2	0.8707	31.8538	348.732
q3	0.0022	0.008	0.012
q4	0.0241	0.0599	0.0932
q5	0.06	0.1908	0.2434
q6	0.2121	11.3648	253.683
q7	0.0106	0.0314	0.1369
q8	1.6634	4.153	9.8969
q9	1.1611	42.0941	981.896
q10	0.0018	0.0058	0.0163
q11	0.0205	0.0599	0.1568
q12	0.0105	0.2159	0.4706
q13	0.0018	0.01	0.2973
q14	0.1947	8.8916	272.475

Table 3: Time needed to produce all query answers in ms.

- 3.c q1: This query consists of two triple patterns, and they are both in the form of  $\langle X, p, o \rangle$ . My query plan picks the second one as a start because it is considerably smaller (1874 v.s. 4). Otherwise, it will be very slow, as Table 1 suggested.
- q2: This query contains two groups of the forms  $\langle X, p, o \rangle$ , and  $\langle X, p, Z \rangle$ , and each group has three triples. The plan is to find a  $\langle X, p, o \rangle$ , then a  $\langle X, p, Z \rangle$  that matches the first pick, and pick whatever matches the variable,  $Z$ , in the remaining unprocessed patterns. It picks not only the least selectivity pattern but also the smallest pattern in each step, resulting in a very short running time (compared to the model solution in Table 1).

<sup>7</sup>Since I brought a new Mac while I was working on this, I attached the results from my old Mac and the school server in Appendix 1 and 2. I was astonished by my new Mac.

- q3: This one is similar to q1, but they have different sizes in each pattern. In LUBM-001-mat.ttl, q3 has 5999 possible matches in the first pattern and 6 in the second; q1 has 1874 in the first and 4 in the second. Therefore, q3 is slightly slower than q1.
- q4: Though this query looks terrifying, the speed is fast because all patterns share the same  $X$  in  $p$ , and besides this, there is at most one free variable in each. Therefore, once  $X$  with the smallest size is picked, the remaining is just to check if a  $Y$  can go with the  $X$ .
- q5: This one is similar to q1 and q3 but is much slower because the search space is large. This one has 8330 possibilities for the first pattern and 719 for the second, so even starting from the second one would slow down by a lot.
- q6: This query simply matches just one pattern. No plan is needed. From the smallest to the largest dataset, the result sizes are 7790, 99566, and 1048532. So the time increases as we expected.
- q7: There are two triples like  $\langle X, p, o \rangle$ . Since the first one has the size 7790 and the second one is 1627, starting at the second one, the performance is much better (see Table 1).
- q8: All three datasets produce 7790 results, so the time is about the same. It also illustrates that we have an excellent query plan because pattern 1 has 7790 matches in the smallest dataset but 1048532 matches in the largest dataset. However, pattern 3 has 239 possible matches among all datasets. Therefore, by choosing pattern 3 as a start, we ensure that the time does not change as the size of pattern 1 changes.
- q9: It is like q8, but the sizes of the first three patterns increase by 12x from the smallest dataset to the largest. Therefore, even if we pick the smallest pattern to start with, the time will still increase.
- q10: Same as q1 and q3.
- q11: Like q10, the performance of q1, q3, q10, and q11 are good because they choose the pattern with the smallest size to start.
- q12: It is almost like q11. It first finds the smallest of the form  $\langle X, p, o \rangle$  and matches  $\langle X, p, Y \rangle$ . Then all others are just  $\langle s, p, o \rangle$  and can be verified very quickly.
- q13: There are 8330 choices for the first pattern but just one for the second. If we do it in order, it will be slow. However, the plan first takes  $\langle s, p, X \rangle$  and then  $\langle X, p, o \rangle$  because they have the same selectivity but different in sizes.
- q14: Same as q6. The running time increases as the size increases.

## Appendix 1: result from my old Mac

### 1. Hardware and Software Configuration:

- (a) Model Identifier: MacBookPro14,3
- (b) Processor Name: Quad-Core Intel Core i7
- (c) Processor Speed: 2.9 GHz
- (d) Number of Processors: 1
- (e) Total Number of Cores: 4
- (f) L2 Cache (per Core): 256 KB
- (g) L3 Cache: 8 MB
- (h) Memory: 16 GB
- (i) Operating System: macOS Big Surf, Version 11.6 (20G165)
- (j) Compiler Version: Apple clang version 11.0.0 (clang-1100.0.33.8)

	LUBM-001-mat (ms)	LUBM-010-mat (ms)	LUBM-100-mat (ms)
Load & Index Time	314.359	4026.82	53525.3

Table 4: Time needed to load and index the data in ms.

	LUBM-001-mat (ms)	LUBM-010-mat (ms)	LUBM-100-mat (ms)
q1	0.0152	0.0143	0.0248
q2	4.4338	72.8246	708.436
q3	0.0162	0.0239	0.0274
q4	0.1119	0.1494	0.1827
q5	0.3815	0.5832	0.7189
q6	0.974	19.8416	186.464
q7	0.0674	0.0898	0.1306
q8	8.3644	11.2896	14.8031
q9	6.1915	100.174	1376.28
q10	0.0106	0.0202	0.0278
q11	0.1376	0.2023	0.2621
q12	0.0409	0.4233	0.5893
q13	0.0075	0.0388	0.5113
q14	0.7846	15.3656	141.389

Table 5: Time needed to produce all query answers in ms.

## Appendix 2: result from the school server

### 1. Hardware and Software Configuration:

- (a) Processor Name: Intel(R) Xeon(R) CPU E5-2660 v3 @ 2.60GHz
- (b) Processor Speed: 2.6 GHz
- (c) Number of Processors: 1
- (d) Total Number of Cores: 4
- (e) L2 Cache (per Core): 1 MiB
- (f) L3 Cache: 100 MiB
- (g) Memory: 32 GB
- (h) Operating System: Fedora 34 (Workstation Edition)
- (i) Compiler Version: g++ (GCC) 11.2.1 20220127 (Red Hat 11.2.1-9)

	LUBM-001-mat (ms)	LUBM-010-mat (ms)	LUBM-100-mat (ms)
Load & Index Time	307.796	4161.18	52040.1

Table 6: Time needed to load and index the data in ms.

	LUBM-001-mat (ms)	LUBM-010-mat (ms)	LUBM-100-mat (ms)
q1	0.0098	0.0156	0.0225
q2	3.1313	62.3343	643.443
q3	0.0103	0.0175	0.0233
q4	0.0845	0.1269	0.1621
q5	0.2518	0.4658	0.5873
q6	0.6214	15.222	165.524
q7	0.0419	0.0622	0.0907
q8	6.0605	8.9462	11.8175
q9	4.1269	78.3835	1003.39
q10	0.0071	0.0139	0.0183
q11	0.0753	0.1217	0.147
q12	0.0273	0.3186	0.4493
q13	0.0038	0.0246	0.3449
q14	0.5543	9.7041	103.243

Table 7: Time needed to produce all query answers in ms.