# Dig Data Management Assignment 6

信科四 郑元嘉 1800920541

In this report, we've got to run 4 tasks individually implemented by 3 databases – Redis, MongoDB and Neo4j.

The 4 tasks are given below :

1. Given si，find all its P and O，<si, P, O>. In particular, we run this task 5 times to have a more reliable consumed time where si are 'Shamsuzzaman_Khan', 'Jake_Gyllenhaal', 'Rubn_Arocha', 'frBgnet', 'Cody_Kessler'.
2. Given oi, find all S and P，<S, P,oi>. In particular, we run this task 5 times to have a more reliable consumed time where oi are 'France', 'Belgium', 'Canada', 'Scotland', 'Spain'.
3. Given p1,p2, find all S which has both p1 and p2，<S, p1, *>, <S, p2, *>. In particular, the (p1, p2) pair is ('isLeaderOf', 'isCitizenOf').
4. Given oi, find the S which has the most oi. If there's more than one result, we only return one of them. In particular, we run this task 5 times to have a more reliable consumed time where oi are 'DurhamSud', 'Norway', 'Germany', 'France', 'China'

Plus, for every single task in some database, both insertion method and search method would be demonstrated below to show how we design the storage strategies and how we search based on it.

## Database 1 - Redis

- **Task 1**
  **Insertion Method**
  (Partial Source Code)

```
if task_index == 0:
    with open(file_path) as f:
        for line in f.readlines():
            data = line.strip().split(' ')
            if len(data) != 3:
                print('Datum error:', data)
            client.rpush(data[0], ' '.join([data[1], data[2]]))
```

Intuitively, we simply apply list key-value to store data with S as the key and (P, O) pair as the value element of the list.

**Search Method**

(Partial Source Code)

```python
1  def search_0(redis_host, redis_port, redis_db_index):
2      global timestamp
3      client = redis.Redis(host=redis_host, port=redis_port, db=redis_db_inde
4
5      timestamp.stamp('search')
6
7      si_list = ['Shamsuzzaman_Khan', 'Jake_Gyllenhaal', 'Rubn_Arocha', 'frBg
8
9      for si in si_list:
10         res = client.lrange(si, 0, -1)
11         pair_list = [set(r.split(' ')) for r in res]
12         print('Search: {} | Number of results: {} | {}'.format(si, len(res)
13
14     timestamp.stamp('search')
```

We search by the key equal to the si, retrieve all elements in its value as a list, and decompose the form of (P, O) pair.

**Consumed Time :**

```
File/hadoopHw/hw6  master ✗                                    2h14m ●
▶ python3 redis_search.py -m search -i 0
Search: Shamsuzzaman_Khan | Number of results: 1 | [{'isLeaderOf', 'Bangla_Academy'}]
Search: Jake_Gyllenhaal | Number of results: 6 | [{'isCitizenOf', 'Russia'}, {'isCitizenOf', 'Lithuania'}, {'isCitizenOf
', 'Sweden'}, {'isCitizenOf', 'Latvia'}, {'United_States', 'isCitizenOf'}, {'isCitizenOf', 'England'}]
Search: Rubn_Arocha | Number of results: 1 | [{'Venezuela', 'isCitizenOf'}]
Search: frBgnet | Number of results: 1 | [{'France', 'isCitizenOf'}]
Search: Cody_Kessler | Number of results: 1 | [{'United_States', 'isCitizenOf'}]
Total time: 0.0017004013061523438
Search time: 0.0016407966613769531
```

- **Task 2**

  **Insertion Method**

  (Partial Source Code)

```python
1  elif task_index == 1:
2      with open(file_path) as f:
3          for line in f.readlines():
4              data = line.strip().split(' ')
5              if len(data) != 3:
6                  print('Datum error:', data)
7              client.rpush(data[2], ' '.join([data[0], data[1]]))
```

It's highly logically similar to task 1. We use O as the key and (S, P) pair as the element in the list value this time.

**Search Method**

(Partial Source Code)

```
 1  def search_1(redis_host, redis_port, redis_db_index):
 2      global timestamp
 3      client = redis.Redis(host=redis_host, port=redis_port, db=redis_db_inde
 4
 5      timestamp.stamp('search')
 6
 7      oi_list = ['France', 'Belgium', 'Canada', 'Scotland', 'Spain']
 8
 9      for oi in oi_list:
10          res = client.lrange(oi, 0, -1)
11          pair_list = [set(r.split(' ')) for r in res]
12          print('Search: {} | Number of results: {}'.format(oi, len(res)))
13
14      timestamp.stamp('search')
```

It's almost the same as task 1 as well. Use oi as the key to retrieve the list value, and decompose the elements in it.

**Consumed Time**

```
File/hadoopHw/hw6  master ✗
▶ python3 redis_search.py -m search -i 1
Search: France | Number of results: 58241
Search: Belgium | Number of results: 11275
Search: Canada | Number of results: 27381
Search: Scotland | Number of results: 13908
Search: Spain | Number of results: 20772
Total time: 0.454378604888916
Search time: 0.45215606689453125
```

- **Task 3**

  **Insertion Method**

  (Partial Source Code)

```
 1  elif task_index == 2:
 2      with open(file_path) as f:
 3          for line in f.readlines():
 4              data = line.strip().split(' ')
 5              if len(data) != 3:
 6                  print('Datum error:', data)
 7              client.rpush(data[1], data[0])
```

P is selected as the key this time with a list value carrying S as elements.

**Search Method**

(Partial Source Code)

```
 1  def search_2(redis_host, redis_port, redis_db_index):
 2      global timestamp
```

```
 3        client = redis.Redis(host=redis_host, port=redis_port, db=redis_db_ind

 4

 5        timestamp.stamp('search')

 6

 7        p_pair = ('isLeaderOf', 'isCitizenOf')

 8

 9        res_0 = client.lrange(p_pair[0], 0, -1)

10        res_1 = client.lrange(p_pair[1], 0, -1)

11

12        res_cross = set(res_0) & set(res_1)

13

14        print('Search: {} | Number of results: {} \n| Partial Samples: {}'.for

15

16        timestamp.stamp('search')
```

We retrieve two values as lists by the p1 and p2 as the keys individually. Then we turn the two lists into instances of `set` to see all elements in the cross set which is our goal.

**Consumed Time**

```
File/hadoopHw/hw6  master ✗                              1d ◓ ⊘
▶ python3 redis_search.py -m search -i 2
Search: ('isLeaderOf', 'isCitizenOf') | Number of results: 6083
| Partial Samples: ['David_A_Granger', 'Yin_Yicui', 'Alexander_Papagos', 'T_J_Harvey', 'Camille_Chamoun']
Total time: 2.8593850135803223
Search time: 2.843515396118164
```

- **Task 4**

  **Insertion Method**

  (Partial Source Code)

```
1  elif task_index == 3:
2      with open(file_path) as f:
3          for line in f.readlines():
4              data = line.strip().split(' ')
5              if len(data) != 3:
6                  print('Datum error:', data)
7              client.zincrby(data[2], 1, data[0])
```

In this task, we apply `zset` type to store the data. O is the key while S is the element in `zset` value. And the score in `zset` type is taken advantage of by us to do the ranking task by the number of occurrences.

  **Search Method**

  (Partial Source Code)

```
1  def search_3(redis_host, redis_port, redis_db_index):
2      global timestamp
3      client = redis.Redis(host=redis_host, port=redis_port, db=redis_db_inde
4
5      timestamp.stamp('search')
```

```
 6
 7        oi_list = ['DurhamSud', 'Norway', 'Germany', 'France', 'China']
 8
 9        for oi in oi_list:
10            res = client.zrevrange(oi, 0, 0)
11
12            print('Search: {} | Result: {}'.format(oi, res))
13
14        timestamp.stamp('search')
```

Thanks to the scoring property in `zset` type, we can obtain the element with highest score (number of occurrences) via mere function `zrevrange()`, and then the S which has the most oi is returned.

**Consumed Time**

```
File/hadoopHw/hw6  master ✗
▶ python3 redis_search.py -m search -i 3
Search: DurhamSud | Result: ['Michel_Nol']
Search: Norway | Result: ['Toril_Marie_ie']
Search: Germany | Result: ['fa_']
Search: France | Result: ['fa_']
Search: China | Result: ['Zhang_Yuan']
Total time: 0.0017142295837402344
Search time: 0.0016629695892333984
```

# Database 2 - MongoDB

- **Task 1**

  **Insertion Method**

  (Partial Source Code)

```
 1  if task_index == 0:
 2      collection.create_index('S')
 3      bulk = []
 4
 5      with open(file_path) as f:
 6          for line in f.readlines():
 7              doc = line.strip().split(' ')
 8              doc = dict([(k, v) for (k, v) in zip(['S', 'P', 'O'], doc)])
 9              bulk.append(doc)
10
11      collection.insert_many(bulk)
```

We purely save a datum as a document in MongoDB with keys `S`, `P` and `O`. To speed up the search task, we create index of `S` field. Here, we can use function `insert_many()` to insert all data at once so as to reduce the transmission.

**Search Method**

(Partial Source Code)

```python
1   def search_0(db_host, db_port, db_name, db_col_name, task_index):
2       global timestamp
3       client = MongoClient(db_host, db_port)
4       db = client[db_name]
5       collection = db[db_col_name]
6
7       si_list = ['Shamsuzzaman_Khan', 'Jake_Gyllenhaal', 'Rubn_Arocha', 'frBg
8
9       timestamp.stamp('search')
10
11      for si in si_list:
12          cursor = collection.find({'S': si})
13          res = [(_doc['P'], _doc['O']) for _doc in cursor]
14          print('Search: {} | Number of results: {} | {}'.format(si, len(res)
15
16      timestamp.stamp('search')
```

We use a simple function `find()` to search for our goal.

**Consumed Time**

```
File/hadoopHw/hw6  master ✗                    2d ●
▶ python3 mongo_search.py -m search -i 0
Search: Shamsuzzaman_Khan | Number of results: 1 | [('isLeaderOf', 'Bangla_Academy')]
Search: Jake_Gyllenhaal | Number of results: 6 | [('isCitizenOf', 'Russia'), ('isCitizenOf', 'Lithuania'), ('isCitizenOf', 'S
weden'), ('isCitizenOf', 'Latvia'), ('isCitizenOf', 'United_States'), ('isCitizenOf', 'England')]
Search: Rubn_Arocha | Number of results: 1 | [('isCitizenOf', 'Venezuela')]
Search: frBgnet | Number of results: 1 | [('isCitizenOf', 'France')]
Search: Cody_Kessler | Number of results: 1 | [('isCitizenOf', 'United_States')]
Total time: 0.003747701644897461
Search time: 0.0028231143951416016
```

- **Task 2**

**Insertion Method**

(Partial Source Code)

```python
1   elif task_index == 1:
2       collection.create_index('O')
3       bulk = []
4
5       with open(file_path) as f:
6           for line in f.readlines():
7               doc = line.strip().split(' ')
8               doc = dict([(k, v) for (k, v) in zip(['S', 'P', 'O'], doc)])
9               bulk.append(doc)
10
11      collection.insert_many(bulk)
```

It's almost similar to Task 1 while we set field `O` as index here.

**Search Method**

(Partial Source Code)

```
1   def search_1(db_host, db_port, db_name, db_col_name, task_index):
2       global timestamp
3       client = MongoClient(db_host, db_port)
4       db = client[db_name]
5       collection = db[db_col_name]
6
7       oi_list = ['France', 'Belgium', 'Canada', 'Scotland', 'Spain']
8
9       timestamp.stamp('search')
10
11      for oi in oi_list:
12          cursor = collection.find({'O': oi})
13          res = [(_doc['S'], _doc['P']) for _doc in cursor]
14          print('Search: {} | Number of results: {}'.format(oi, len(res)))
15
16      timestamp.stamp('search')
```

Doing the same thing as the one in Task 1, we use function `find()` and use `O` as our condition,

**Consumed Time**

```
File/hadoopHw/hw6   master ✗
▶ python3 mongo_search.py -m search -i 1
Search: France | Number of results: 58241
Search: Belgium | Number of results: 11275
Search: Canada | Number of results: 27381
Search: Scotland | Number of results: 13908
Search: Spain | Number of results: 20772
Total time: 0.4867572784423828
Search time: 0.4843409061431885
```

- **Task 3**

**Insertion Method**

(Partial Source Code)

```
1   elif task_index == 2:
2       collection.create_index('S', unique=True)
3       collection.create_index('P_list')
4       bulk = []
5
6       with open(file_path) as f:
7           for line in f.readlines():
8               doc = line.strip().split(' ')
9               bulk.append(UpdateMany({'S': doc[0]}, {'$addToSet': {'P_list':
```

```
10
11        collection.bulk_write(bulk)
```

Here's something different. We save data as a document with fields `S` and `P_list` where all Ps would be saved into a list in the same document by operation `$addToSet` and unique index `S`. To speed up update work, we have class `UpdateMany` and function `bulk_write()` to give us a favor, which allows us to send mere one request to MongoDB and run multiple operations.

**Search Method**

(Partial Source Code)

```python
1   def search_2(db_host, db_port, db_name, db_col_name, task_index):
2       global timestamp
3       client = MongoClient(db_host, db_port)
4       db = client[db_name]
5       collection = db[db_col_name]
6
7       p_pair = ('isLeaderOf', 'isCitizenOf')
8
9       timestamp.stamp('search')
10
11      cursor = collection.find({'P_list': {'$all': list(p_pair)}})
12      res = [_doc['S'] for _doc in cursor]
13      print('Search: {} | Number of results: {} \n| Partial Samples: {}'.form
14
15      timestamp.stamp('search')
```

We use MongoDB list operation `$all$` to check whether the elements in `P_list` can match all of ones in our query list. By so, we can retrieve all matched documents with ease.

**Consumed Time**

```
File/hadoopHw/hw6  master x                           2d ●
▶ python3 mongo_search.py -m search -i 2
Search: ('isLeaderOf', 'isCitizenOf') | Number of results: 6083
| Partial Samples: ['Elizabeth_II', 'William_H_Seward_Jr', 'Andranik', 'Ramasamy_Palanisamy', 'Leonard_Leo']
Total time: 1.0357365608215332
Search time: 1.0338950157165527
```

- **Task 4**

**Insertion Method**

(Partial Source Code)

```python
1   elif task_index == 3:
2       collection.create_index('O')
3       # create compound index
4       collection.create_index([('O', pymongo.ASCENDING), ('S', pymongo.ASCEND
5       bulk = []
6
```

```
 7        with open(file_path) as f:
 8            for line in f.readlines():
 9                doc = line.strip().split(' ')
10                bulk.append(UpdateOne({'O': doc[2], 'S': doc[0]}, {'$inc': {'co
11
12        collection.bulk_write(bulk)
13
14   timestamp.stamp('insert')
```

This time, we save a datum as a document with fields `O` , `S` and `count` . Fields `O` and `S` work as the index for searching with field `count` denotes the number of occurrence of the ( `O` , `S` ) pair. We create a unique compound index consisting of `O` and `S` to ensure the uniqueness and to enable us to use class `UpdateOne` rather than `UpdateMany` furthermore. Also, we use function `bulk_write()` to speed up.

**Search Method**

(Partial Source Code)

```
 1  def search_3(db_host, db_port, db_name, db_col_name, task_index):
 2    global timestamp
 3    client = MongoClient(db_host, db_port)
 4    db = client[db_name]
 5    collection = db[db_col_name]
 6
 7    oi_list = ['DurhamSud', 'Norway', 'Germany', 'France', 'China']
 8
 9    timestamp.stamp('search')
10
11    for oi in oi_list:
12        cursor = collection.aggregate([
13            {'$match': {'O': oi}},
14            {'$sort': {'count':-1}},
15            {'$limit': 1},
16        ])
17
18        res = list(cursor)[0]['S']
19        print('Search: {} | Result: {}'.format(oi, res))
20
21    timestamp.stamp('search')
```

To find the document with maximal value of `count` , we adapt an aggregation pipeline. The first step `$match` is simply to filter and get the candidates, and the second step `$sort` as well as third step `$limit` help us to obtain the first element in a list in descendent order by `count` .

```
File/hadoopHw/hw6   master ✗
▶ python3 mongo_search.py -m search -i 3
Search: DurhamSud | Result: Michel_Nol
Search: Norway | Result: Harald_V_of_Norway
Search: Germany | Result: fa_
Search: France | Result: fa_
Search: China | Result: Duke_Shn_of_Chen
Total time: 0.45873594284057617
Search time: 0.45747876167297363
```

# Database 3 - Neo4j

- **Task 1**

**Insertion Method**

(Partial Source Code)

```python
1  if task_index == 0:
2      statement = 'CREATE (:set1{s: {s}, p: {p}, o: {o}})'
3      with driver.session() as session:
4          with open(file_path) as f:
5              for i, line in enumerate(f.readlines()):
6                  if i % 100 == 0: print('{} rounds'.format(i))
7                  doc = line.strip().split(' ')
8                  session.run(statement, s=doc[0], p=doc[1], o=doc[2])
```

Because the speed of creating nodes in Neo4j is dramatically slow, we apply this simplest way to all 4 tasks. We create one node with 3 fields `S`, `P` and `O`. Therefore, there's no distinct insertion method going to be demonstrated below.

**Search Method**

(Partial Source Code)

```python
1  def search_0(file_path, uri, acc, pwd, task_index):
2      driver = GraphDatabase.driver(uri, auth=(acc, pwd))
3
4      global timestamp
5      timestamp.stamp('search')
6
7      si_list = ['Shamsuzzaman_Khan', 'Spencer_Chandra_Herbert', 'Minnesota_P
8      statement = 'MATCH (set:set1{s: {s}}) RETURN set'
9
10     for si in si_list:
11         with driver.session() as session:
12             res = session.run(statement, s=si)
```

```
13                data = list(res.records())
14                print('Search task returns {} matched items.'.format(len(data))
15
16        timestamp.stamp('search')
```

It's a simple search task. We only apply the condition to retrieve our wanted data.

**Consumed Time**

```
File/hadoopHw/hw6   master ✗
▶ python3 neo4j_search.py -m search -i 0
Search task returns 1 matched items.
Search task returns 2 matched items.
Search task returns 40 matched items.
Search task returns 1 matched items.
Search task returns 1 matched items.
Total time: 0.47410130500793457
Search time: 0.440701961517334
```

- **Task 2**

  **Search Method**

  (Partial Source Code)

```
 1  def search_1(file_path, uri, acc, pwd, task_index):
 2      driver = GraphDatabase.driver(uri, auth=(acc, pwd))
 3
 4      global timestamp
 5      timestamp.stamp('search')
 6
 7      oi_list = oi_list = ['France', 'Belgium', 'Canada', 'Scotland', 'Spain'
 8      statement = 'MATCH (set:set1{o: {o}}) RETURN set'
 9
10      for oi in oi_list:
11          with driver.session() as session:
12              res = session.run(statement, o=oi)
13              data = list(res.records())
14              print('Search task returns {} matched items.'.format(len(data))
15
16      timestamp.stamp('search')
```

It's the same as Task 1.

**Consumed Time**

```
File/hadoopHw/hw6  master ✗
▶ python3 neo4j_search.py -m search -i 1
Search task returns 4833 matched items.
Search task returns 992 matched items.
Search task returns 2196 matched items.
Search task returns 1189 matched items.
Search task returns 1770 matched items.
Total time: 0.8494319915771484
Search time: 0.8033356666564941
```

- **Task 3**

  **Search Method**

  (Partial Source Code)

```
 1  def search_2(file_path, uri, acc, pwd, task_index):
 2      driver = GraphDatabase.driver(uri, auth=(acc, pwd))
 3
 4      global timestamp
 5      timestamp.stamp('search')
 6
 7      p_pair = ('isLeaderOf', 'isCitizenOf')
 8      statement = 'MATCH (set:set1{p: {p}}) RETURN set'
 9
10      with driver.session() as session:
11          res_0 = session.run(statement, p=p_pair[0])
12          res_1 = session.run(statement, p=p_pair[1])
13          data_0 = set([item['set']['s'] for item in list(res_0.records())])
14          data_1 = set([item['set']['s'] for item in list(res_1.records())])
15          print('Search task returns {} matched items.'.format(len(data_0 & d
16
17      timestamp.stamp('search')
```

  Here, we use the simple way to get two set according to p1 and p2, and then we extract the cross get by manipulating class `set` to get the final result.

  **Consumed Time**

```
File/hadoopHw/hw6  master ✗
▶ python3 neo4j_search.py -m search -i 2
Search task returns 575 matched items.
Total time: 4.998905420303345
Search time: 4.96355000789795
```

- **Task 4**

  **Search Method**

  (Partial Source Code)

```
 1  def search_3(file_path, uri, acc, pwd, task_index):
 2      driver = GraphDatabase.driver(uri, auth=(acc, pwd))
 3
 4      global timestamp
```

```
5        timestamp.stamp('search')
6
7        oi_list = ['DurhamSud', 'Norway', 'Germany', 'France', 'China']
8        statement = 'MATCH (set:set1{o: {o}}) RETURN set'
9
10       for oi in oi_list:
11           with driver.session() as session:
12               res = session.run(statement, o=oi)
13               data = [item['set']['s'] for item in list(res.records())]
14               counter = Counter(data)
15               max_ele = sorted(counter.items(), key=lambda x: x[1], reverse=
16               print('Search: {} | Result: {}'.format(oi, max_ele))
17
18       timestamp.stamp('search')
```

We get the matched data back, and do the sorting task by number of concurrence in aid of class `Counter` and function `sort()` to achieve our goal.

**Consumed Time**

```
File/hadoopHw/hw6  master ✗
▶ python3 neo4j_search.py -m search -i 3
Search: DurhamSud | Result: ('Michel_Nol', 1)
Search: Norway | Result: ('Toril_Marie_ie', 2)
Search: Germany | Result: ('FrankWalter_Steinmeier', 1)
Search: France | Result: ('Paul_Tissandier', 1)
Search: China | Result: ('Xi_Jinping', 1)
Total time: 1.0068144798278809
Search time: 0.9719753265380859
```