



Università degli studi di Messina Dipartimento di
Scienze Matematiche, Informatiche, Scienze
Fisiche e Scienze della Terra (MIFT)

Corso di Laurea Triennale in INFORMATICA

Progetto Basi di Dati NoSQL

PROBLEMATICA AFFRONTATA

Si calcola che in Italia vengano effettuate in media 160 milioni di telefonate al giorno. Una mole di dati esorbitanti, dunque, è immagazzinata nei tabulati telefonici, soprattutto se si pensa che le informazioni riguardanti ogni singola chiamata vengono cancellate solo dopo 24 mesi.

In ambito giudiziario l'acquisizione dei tabulati telefonici è una pratica che può rivelarsi fondamentale per lo svolgimento di un caso; tuttavia, ciò è possibile solamente entro le 48 ore successive all'attuazione della telefonata.

I database utilizzati per il confronto prestazionale sono MongoDB e Neo4j.

- MongoDB memorizza i dati in documenti flessibili JSON-like, il che significa che i campi possono variare da un documento all'altro ed è possibile modificare nel tempo la struttura dei dati.

Il modello di documento mappa gli oggetti del codice applicativo, semplificando il lavoro sui dati.

Le query ad hoc, l'indicizzazione e l'aggregazione in tempo reale offrono efficaci modalità di accesso e analisi dei dati.

MongoDB è concepito in origine come database distribuito; l'alta disponibilità, la scalabilità orizzontale e la distribuzione geografica sono quindi native e facili da usare.

- Neo4j è un DBMS di tipo NoSQL orientato ai grafi (graph-oriented), open-source, che fornisce un back-end transazionale conforme ad ACID, disponibile pubblicamente dal 2007.

È offerto come servizio gestito tramite AuraDB, ma è anche possibile eseguire personalmente Neo4j con la Community Edition o la Enterprise Edition.

Neo4j è scritto in Java e Scala e il codice sorgente è disponibile su GitHub.

Essendo un database a grafi nativo, implementa un vero modello grafico fino al livello di archiviazione; i dati non vengono archiviati come "astrazione dei grafi" in cima ad un'altra tecnologia.

SOLUZIONE DB CONSIDERATA

I database sono stati creati tramite l'utilizzo del linguaggio di programmazione JavaScript, eseguito tramite "node.js" (ovvero l'ambiente di runtime per l'esecuzione di JavaScript), e popolati con dati fittizi generati grazie alla libreria "Faker.js" e salvati in file JSON.

```
JS dati75.mjs •
JS dati75.mjs > _
1  import { faker } from '@faker-js/faker';
2  import fs from 'fs';
3
4  faker.locale = 'it';
5
6  var person = [];
7  var call = [];
8
9  for (var i=1; i<=5000; i++) {
10
11     var First_name = faker.name.firstName();
12     var Last_name = faker.name.lastName();
13     var Full_name = First_name + Last_name;
14     var Phone_number = faker.phone.number('### ### ####');
15
16     var First_name2 = faker.name.firstName();
17     var Last_name2 = faker.name.lastName();
18     var Full_name2 = First_name2 + Last_name2;
19     var Phone_number2 = faker.phone.number('### ### ####');
20
21     person.push({
22       "Full_name": Full_name,
23       "First_name": First_name,
24       "Last_name": Last_name,
25       "Phone_number": Phone_number,
26       "Calls_done": i
27     });
28     person.push({
29       "Full_name": Full_name2,
30       "First_name": First_name2,
31       "Last_name": Last_name2,
32       "Phone_number": Phone_number2,
33       "Calls_done": i
34     });
35
36     var Call_start = Math.floor(Math.random() * 31535399) + 1640991600;
37     var Duration = Math.floor(Math.random() * 600);
38     var Call_end = Call_start + Duration;
39
40     var City = faker.address.city();
41     var Address = faker.address.streetAddress(false);
42     var Cell_site = Math.floor(Math.random() * 999) + 1;
43
44     call.push({
45       "ID": i,
46       "Call_start": Call_start,
47       "Call_end": Call_end,
48       "Duration": Duration,
49       "City": City,
50       "Address": Address,
51       "Cell_site": Cell_site
52     });
53   }
54
55   let dataObj = person;
56
57   fs.writeFile('people75.json', JSON.stringify(dataObj, null, '\t'), function(err, result) {
58     if(err) console.log('error', err);
59     else console.log('File saved');
60   });
61
62   let dataObj2 = call;
63
64   fs.writeFile('calls75.json', JSON.stringify(dataObj2, null, '\t'), function(err, result) {
65     if(err) console.log('error', err);
66     else console.log('File saved');
67   });
```

PROGETTAZIONE

MongoDB

I dati all'interno di MongoDB sono stati gestiti con dei documenti incorporati (embedded documents), una funzionalità specifica del DBMS in questione che permette la gestione dei dati in gruppi e come singoli elementi.

Di seguito è presente un esempio di come si presentano i documenti gestiti da MongoDB

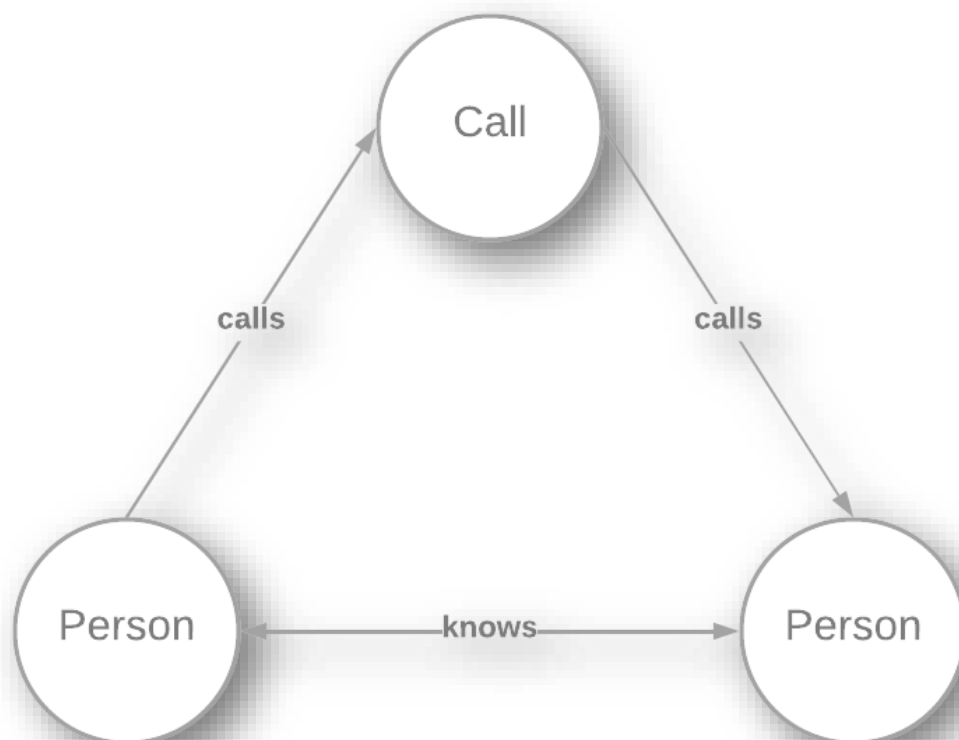
```
2      {
3          "Full_name": "LucianoMichelucci",
4          "First_name": "Luciano",
5          "Last_name": "Michelucci",
6          "Phone_number": "619 887 4624",
7          "Calls_done": [
8              {
9                  "_id": 1,
10                 "Call_start": 1653559846,
11                 "Call_end": 1653560382,
12                 "Duration": 536,
13                 "City": "Sesto Ghita",
14                 "Address": "Strada Ariberto 339",
15                 "Cell_site": 650
16             }
17         ]
18     },
```

e i rispettivi attributi:

- **"FULL_NAME"**: Nome e cognome della persona (stringa);
- **"FIRST_NAME"**: Nome della persona (stringa);
- **"LAST_NAME"**: Cognome della persona (stringa);
- **"PHONE_NUMBER"**: Numero di telefono della persona (stringa).
- **"CALLS_DONE"**: Insieme delle chiamate effettuate dalla persona (array di documenti), ciascuna con:
 - **"ID"**: L'identificativo unico per ciascuna chiamata (numero intero);
 - **"CALL_START"**: Data e ora dell'inizio della chiamata, rappresentata in UNIX timestamp (numero intero) (1/1/2022 0:00:00 – 31/12/2022 23:49:59);
 - **"DURATION"**: Durata della chiamata in secondi (numero intero);

- **"CALL_END"**: Data e ora della fine della chiamata, rappresentata in UNIX timestamp (numero intero);
- **"CITY"**: Città nella quale è avvenuta la chiamata (stringa);
- **"ADDRESS"**: Indirizzo in cui è avvenuta la chiamata (stringa);
- **"CELL_SITE"**: Cella telefonica della chiamata (numero intero);

Neo4j



Questo modello grafico rappresenta le relazioni che le persone ("Person") hanno con la telefonata ("Call") e tra di loro (chi telefona e chi riceve la telefonata).

Per ogni "Call" abbiamo sempre due "Persone" che, essendo state in contatto, sono l'una a conoscenza dell'esistenza dell'altra, questa relazione è fondamentale per avere una rete di contatti durante le indagini.

Ogni entità “Call” è composta dai campi:

- **“ID”**: L’identificativo unico per ciascuna chiamata (numero intero);
- **“CALL_START”**: Data e ora dell’inizio della chiamata, rappresentata in UNIX timestamp (numero intero) (1/1/2022 0:00:00 – 31/12/2022 23:49:59);
- **“DURATION”**: Durata della chiamata in secondi (numero intero);
- **“CALL_END”**: Data e ora della fine della chiamata, rappresentata in UNIX timestamp (numero intero);
- **“CITY”**: Città nella quale è avvenuta la chiamata (stringa);
- **“ADDRESS”**: Indirizzo in cui è avvenuta la chiamata (stringa);
- **“CELL_SITE”**: Cella telefonica della chiamata (numero intero);

Ogni entità “Person” è composta dai campi:

- **“FULL_NAME”**: Nome e cognome della persona (stringa);
- **“FIRST_NAME”**: Nome della persona (stringa);
- **“LAST_NAME”**: Cognome della persona (stringa);
- **“PHONE_NUMBER”**: Numero di telefono della persona (stringa).
- **“CALLS_DONE.ID”**: Insieme degli ID delle chiamate effettuate dalla persona (numero intero)

IMPLEMENTAZIONE

Dopo aver creato i file JSON contenenti i dati da inserire nei vari database, si effettua la connessione a MongoDB tramite JavaScript grazie alla libreria “MongoClient” e si inseriscono i dati letti dal file grazie al comando “.insertMany”.

```
JS inserimento75.mjs X
JS inserimento75.mjs > ...
1  import { MongoClient } from 'mongodb';
2  import fs from 'fs';
3
4  const dbName = 'CallsRecord75';
5  const client = new MongoClient('mongodb://127.0.0.1:27017');
6
7  var startTime = performance.now();
8
9  client.connect(function(err) {
10
11     console.log('Connected successfully to server');
12     const db = client.db(dbName);
13     const data = fs.readFileSync('calls75.json');
14     const docs = JSON.parse(data.toString());
15
16     db.collection('calls')
17     .insertMany(docs, function(err, result) {
18         if (err) throw err;
19         console.log('Inserted docs:', result.insertedCount);
20         client.close();
21
22         var endTime = performance.now();
23
24         let mia = `${endTime - startTime} msec\n`;
25
26         console.log(mia);
27
28         fs.appendFile('Inserimento75.txt', mia, (err) => {
29             if (err) throw err;
30         });
31     });
32 });
33
34 });
```

Per quanto riguarda Neo4j, per prima cosa sono stati creati due indici per velocizzare l'inserimento dei dati grazie ai seguenti comandi Cypher (query language di Neo4j):

```
neo4j$ CREATE INDEX Person FOR (nodeP:Person) ON (nodeP.Full_name)
neo4j$ CREATE INDEX Call FOR (nodeC:Call) ON (nodeC.ID)
```

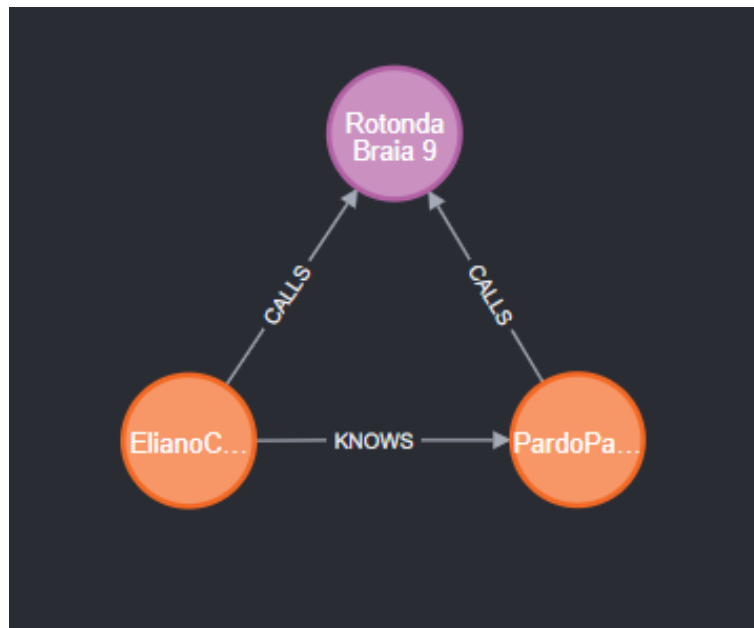
La connessione al database avviene grazie al driver ufficiale fornito da Neo4j per JavaScript “neo4j-driver”.

L'inserimento dei dati prelevati dai file avviene tramite comandi Cypher con il comando “apoc.load.json”.

Dopo aver importato tutte le entità sono state create le relazioni che le interessano usando sempre istruzioni Cypher.

```
JS inserimento75.mjs •
JS inserimento75.mjs > ...
1  import neo4j from 'neo4j-driver'
2  import fs from 'fs'
3
4  const uri = "neo4j://localhost:7687";
5  const user = "neo4j";
6  const password = "admin";
7
8  const driver = neo4j.driver(uri, neo4j.auth.basic(user, password));
9  const session = driver.session();
10
11 const startTime = performance.now();
12
13 await session.run(
14   'CALL apoc.load.json("file:///people75.json") YIELD value UNWIND value AS persone MERGE (p:Person{Full_name: persone.Full_name})',
15   'SET p.First_name = persone.First_name, p.Last_name = persone.Last_name, p.Phone_number = persone.Phone_number',
16 );
17 await session.run(
18   'CALL apoc.load.json("file:///calls75.json") YIELD value UNWIND value AS chiamate MERGE (c:Call{ID: chiamate.ID}) SET c.Call_start = chiamate.Call_start',
19   'c.Call_end = chiamate.Call_end, c.Duration = chiamate.Duration, c.City = chiamate.City, c.Address = chiamate.Address, c.Cell_site = chiamate.Cell_site',
20 );
21
22 const endTime = performance.now();
23
24 await session.run(
25   'CALL apoc.load.json("file:///people75.json") YIELD value UNWIND value AS persone MERGE (p:Person{Full_name: persone.Full_name}) WITH p',
26   'persone UNWIND persone.Calls_done as chiamate MERGE (c:Call{ ID: chiamate }) MERGE (p)-[:CALLS]->(c)',
27 );
28 await session.run(
29   'MATCH (a:Person)-[:CALLS]-(c:Call)-[:CALLS]-(b:Person) MERGE (a)-[:KNOWS]-(b)',
30 );
31
32 await session.close();
33
34
35 let mia = `${endTime - startTime} msec\n`;
36
37 console.log(mia);
38
39 fs.appendFile('Inserimento75.txt', mia, (err) => {
40   if (err) throw err
41 });
42
43 await driver.close();
```


Questo è un esempio di come si presentano le entità con le relative relazioni in Neo4j



In entrambi i programmi JavaScript sopra riportati, il tempo di inserimento dei dati nei database è stampato a video, oltre ad essere salvato in un file di testo, grazie a questa porzione di codice.

```
7   var startTime = performance.now();
8
9   // inserimento
10
11  var endTime = performance.now();
12
13  let mia = `${endTime - startTime} msec\n`;
14
15  console.log(mia);
16
17  fs.appendFile('Inserimento75.txt', mia, (err) => {
18    if (err) throw err;
19  });
```

Di seguito troviamo una tabella rappresentante i tempi d'inserimento dei vari dataset per entrambi i DBMS

	<i>MongoDB</i>	<i>Neo4j</i>
25%	102,22	162,80
50%	151,52	180,35
75%	396,90	865,73
100%	2.598,78	8.150,76

ESPERIMENTI

Per confrontare le prestazioni dei due DBMS in esame sono state eseguite 5 query ciascuno, calcolando il tempo di esecuzione.

Ogni query è stata eseguita 31 volte prendendo in considerazione:

- Il tempo della prima esecuzione (in millisecondi);
- Il tempo medio delle 30 esecuzioni successive (in millisecondi);
- La deviazione standard;
- L'intervallo di confidenza al 95%.

Tutti i dati sopra citati sono stati immagazzinati in delle tabelle e hanno permesso di creare degli istogrammi.

QUERY 1

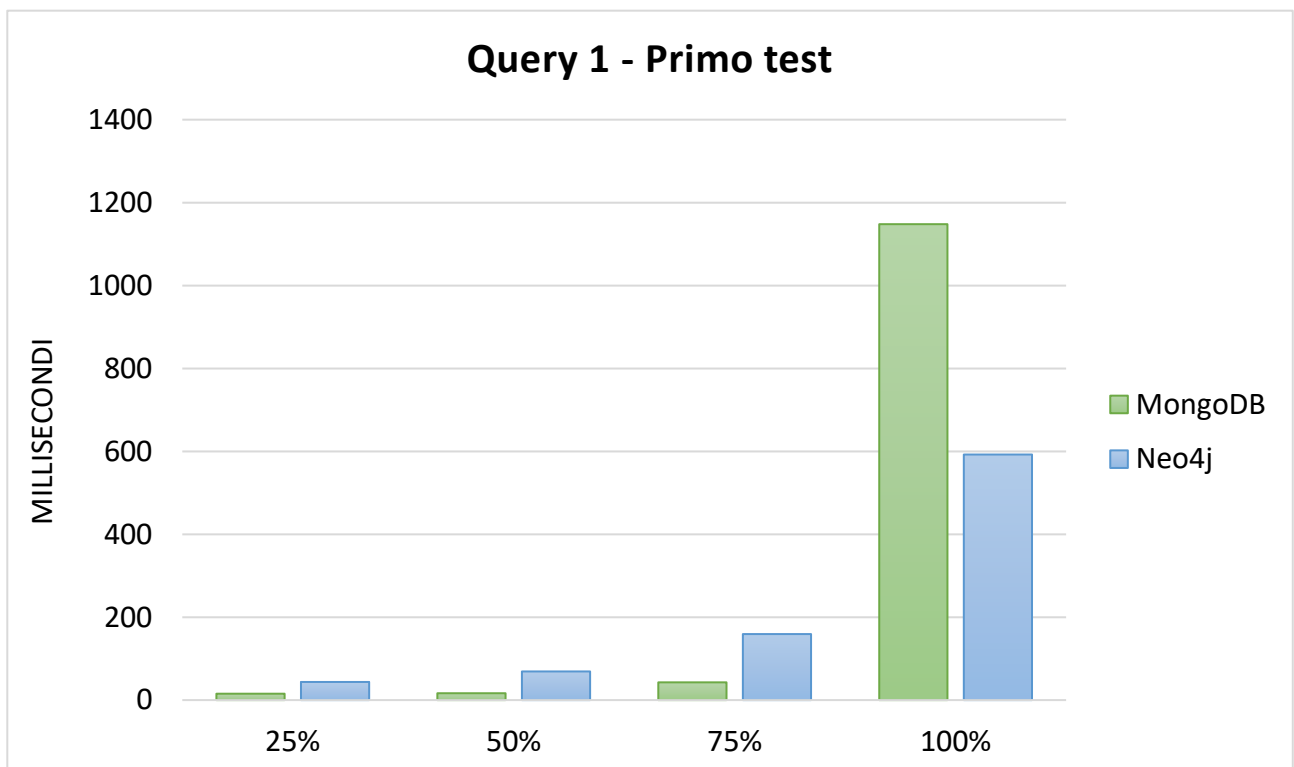
MongoDB:

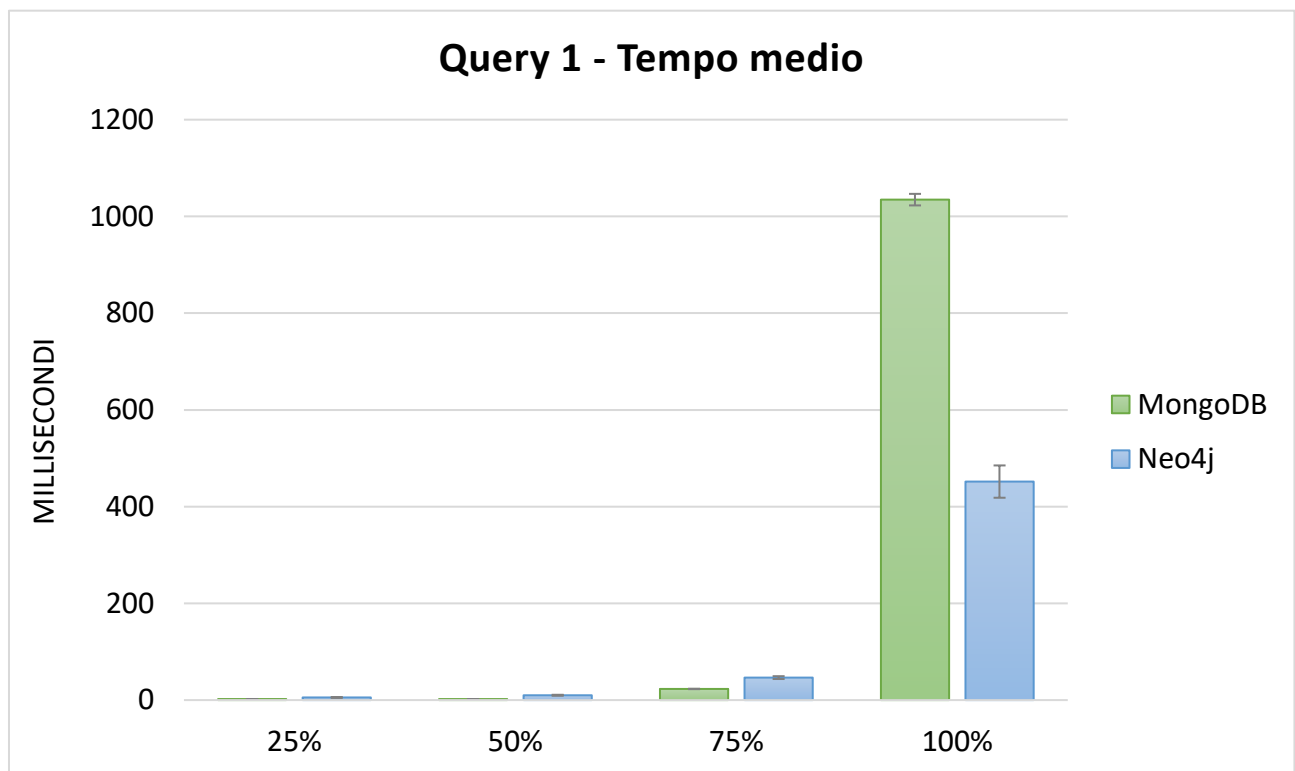
```
const query = { "Calls_done.Cell_site": { $gt: 100 } };
```

Neo4j:

```
await session.run('MATCH (n:Call) WHERE n.Cell_site > 100 RETURN n',  
);
```

	MongoDB 25%	Neo4j 25%	MongoDB 50%	Neo4j 50%	MongoDB 75%	Neo4j 75%	MongoDB 100%	Neo4j 100%
Primo test	15,42	43,74	16,54	69,10	42,81	159,26	1.147,91	592,33
Tempo medio	2,13	5,46	2,10	9,96	23,05	46,49	1.034,66	451,73
Dev. standard	0,44	0,98	0,34	2,82	1,23	8,15	33,31	93,34
Conf. 95%	0,16	0,35	0,12	1,01	0,44	2,92	11,92	33,40





QUERY 2

MongoDB:

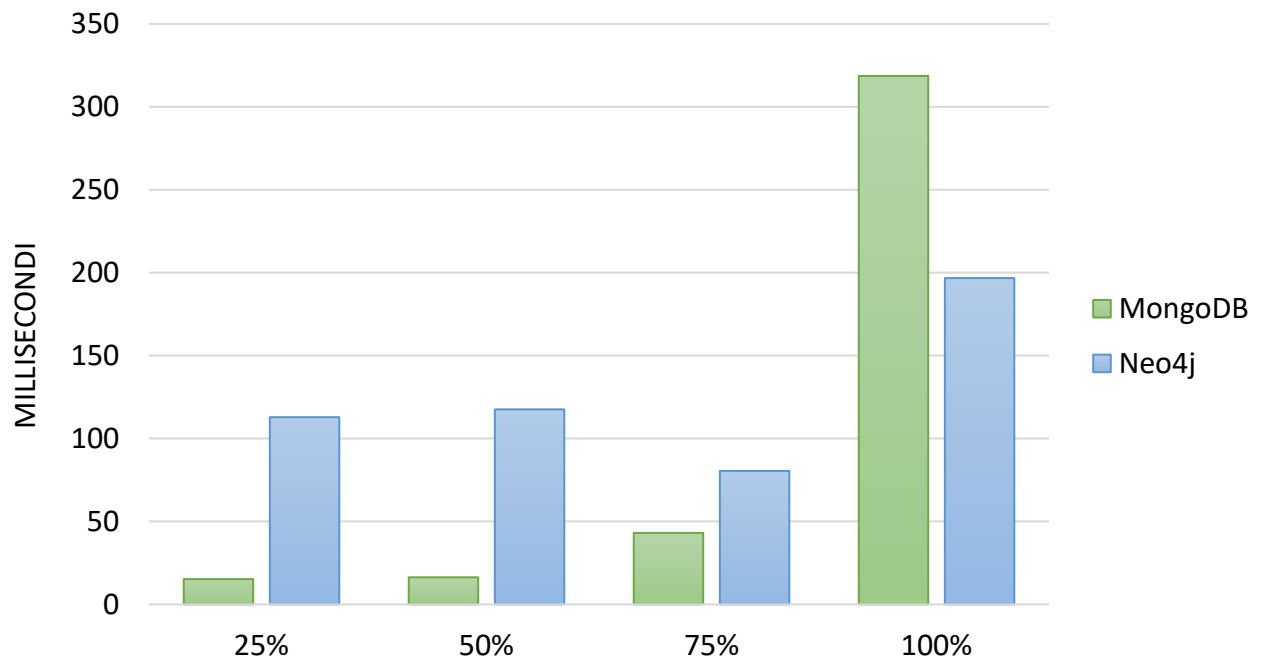
```
const query = { "Calls_done.Cell_site": { $gt: 100 }, "Calls_done.City": { $gt: "V" } };
```

Neo4j:

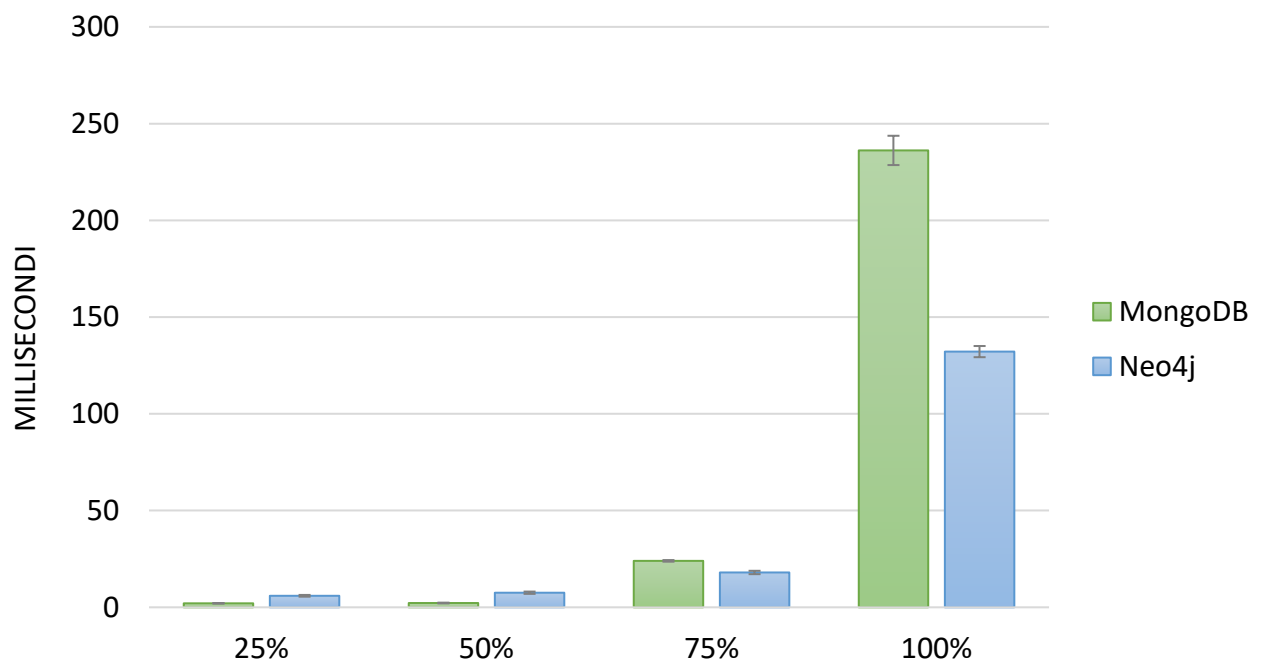
```
await session.run('MATCH (n:Call) WHERE n.Cell_site > 100 AND n.City > "V" RETURN n',
);
```

	MongoDB 25%	Neo4j 25%	MongoDB 50%	Neo4j 50%	MongoDB 75%	Neo4j 75%	MongoDB 100%	Neo4j 100%
Primo test	15,23	112,88	16,33	117,59	43,08	80,46	318,63	196,77
Tempo medio	2,06	5,95	2,22	7,54	24,00	18,01	236,15	132,17
Dev. standard	0,34	1,26	0,54	1,67	1,16	2,43	21,15	8,06
Conf. 95%	0,12	0,45	0,19	0,60	0,42	0,87	7,57	2,88

Query 2 - Primo test



Query 2 - Tempo medio



QUERY 3

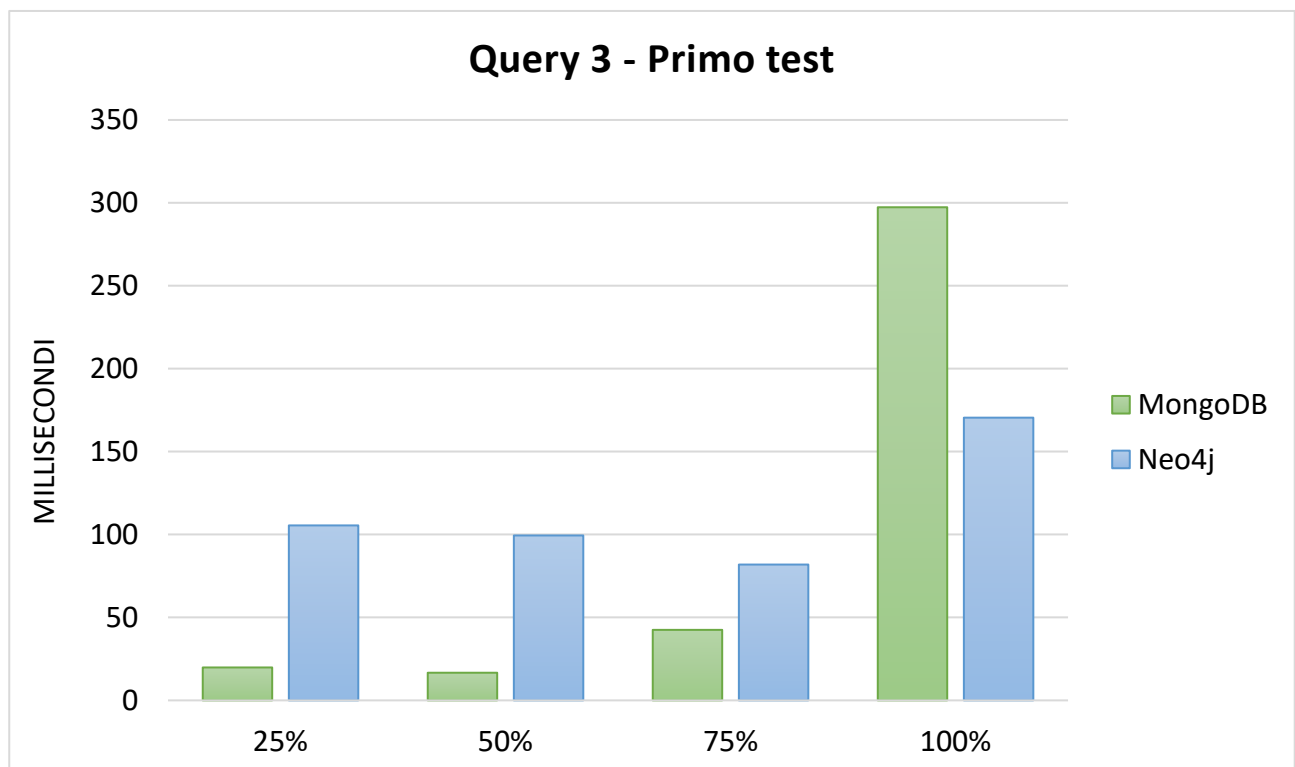
MongoDB:

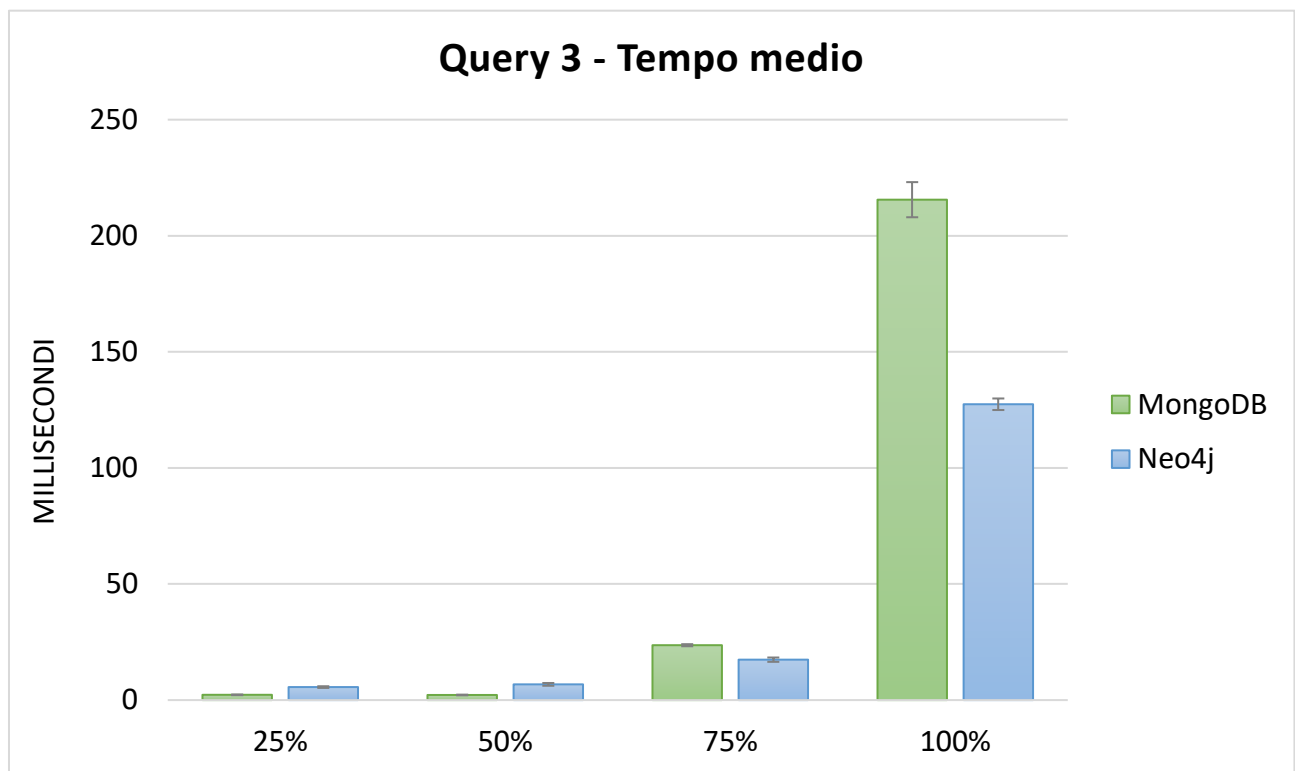
```
const query = { "Calls_done.Cell_site": { $gt: 100 }, "Calls_done.City": { $gt: "V" },  
                "Calls_done.Call_start": { $gt: 1660000000 } };
```

Neo4j:

```
await session.run('MATCH (n:Call) WHERE n.Cell_site > 100 AND n.City > "V",  
                  'AND n.Call_start > 1660000000 RETURN n',  
);
```

	MongoDB 25%	Neo4j 25%	MongoDB 50%	Neo4j 50%	MongoDB 75%	Neo4j 75%	MongoDB 100%	Neo4j 100%
Primo test	19,85	105,50	16,67	99,44	42,52	81,90	297,36	170,48
Tempo medio	2,28	5,61	2,20	6,78	23,64	17,41	215,53	127,43
Dev. standard	0,40	1,01	0,39	1,67	1,25	2,61	21,15	6,96
Conf. 95%	0,14	0,36	0,14	0,60	0,45	0,93	7,57	2,49





QUERY 4

MongoDB:

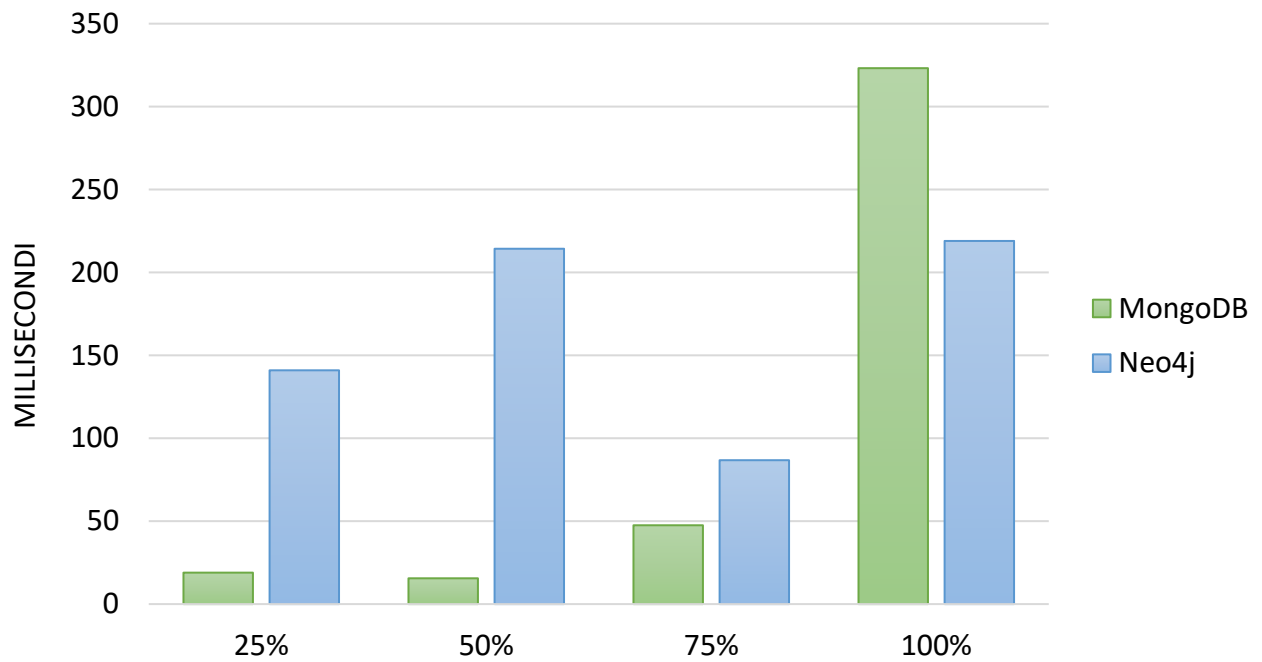
```
const query = { "Calls_done.Cell_site": { $gt: 100 }, "Calls_done.City": { $gt: "V" },
  "Calls_done.Call_start": { $gt: 1660000000 }, "Phone_number": { $gt: "700" } };
```

Neo4j:

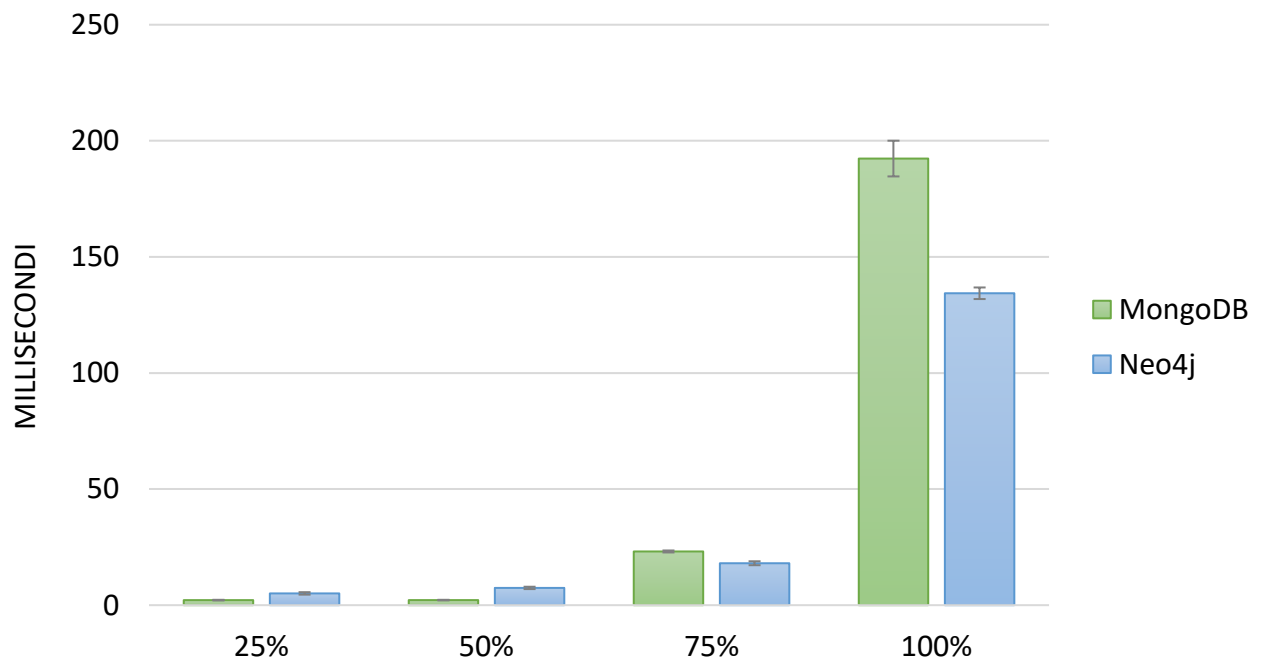
```
await session.run('MATCH (n:Call)-[r]-(p:Person) WHERE n.Cell_site > 100 AND n.City > "V",
  'AND n.Call_start > 1660000000 AND p.Phone_number > "700" RETURN n',
);
```

	MongoDB 25%	Neo4j 25%	MongoDB 50%	Neo4j 50%	MongoDB 75%	Neo4j 75%	MongoDB 100%	Neo4j 100%
Primo test	18,96	140,97	15,55	214,27	47,53	86,76	323,15	218,96
Tempo medio	2,17	5,06	2,18	7,43	23,13	18,06	192,36	134,34
Dev. standard	0,37	1,48	0,36	1,42	1,20	2,35	21,51	6,97
Conf. 95%	0,13	0,53	0,13	0,51	0,43	0,84	7,70	2,49

Query 4 - Primo test



Query 4 - Tempo medio



QUERY 5:

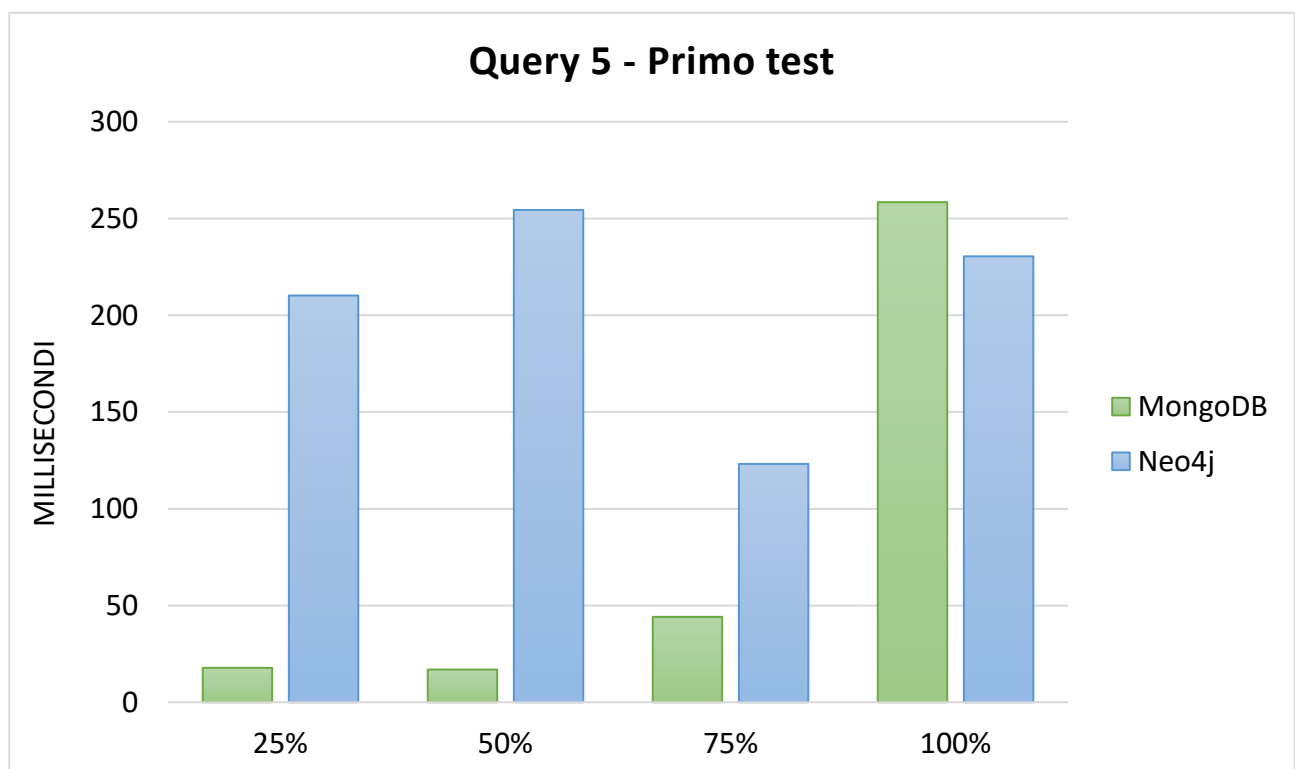
MongoDB:

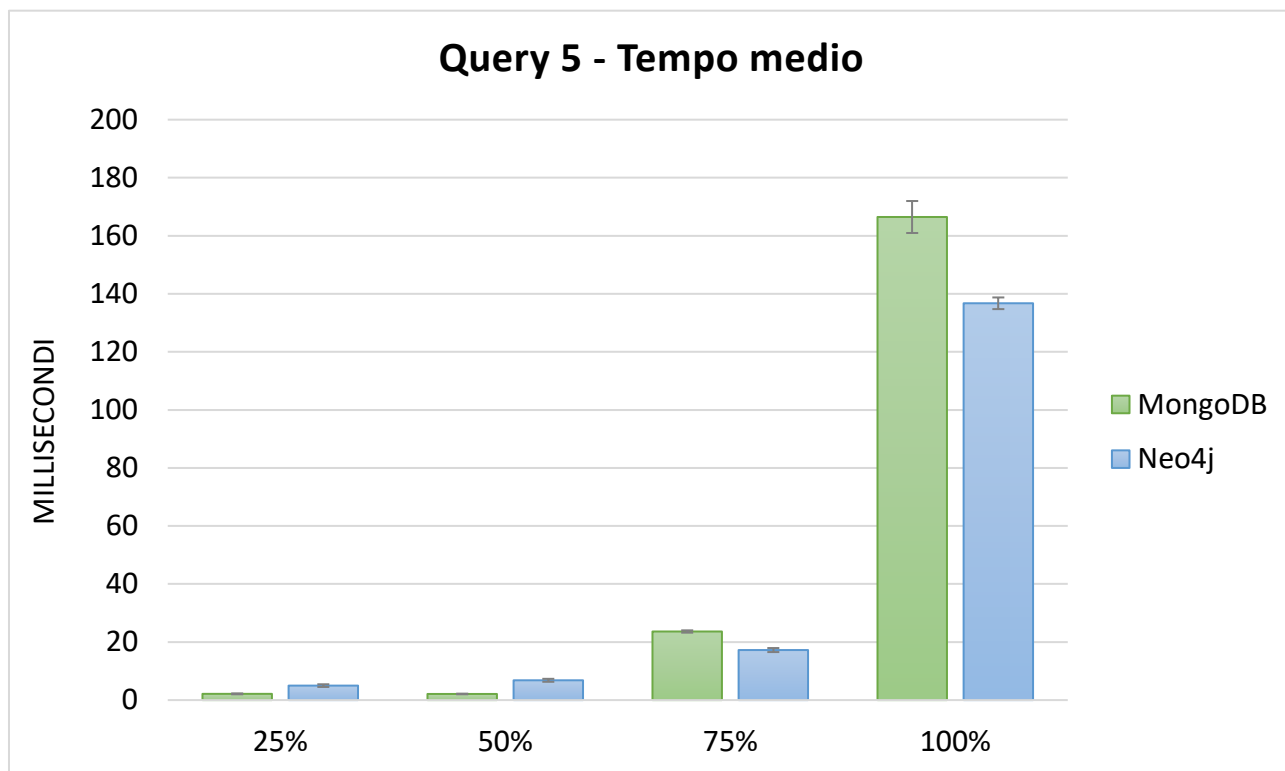
```
const query = { "Calls_done.Cell_site": { $gt: 100 }, "Calls_done.City": { $gt: "V" }, "Calls_done.Call_start": { $gt: 1660000000 },  
$or: [ { "Phone_number": { $gt: "700" } }, { "Last_name": { $gt: "S" } } ] };
```

Neo4j:

```
await session.run('MATCH (n:Call)-[r]-(p:Person) WHERE n.Cell_site > 100 AND n.City > "V" AND n.Call_start > 1660000000',  
'AND (p.Phone_number > "700" OR p.Last_name > "S") RETURN n',  
);
```

	MongoDB 25%	Neo4j 25%	MongoDB 50%	Neo4j 50%	MongoDB 75%	Neo4j 75%	MongoDB 100%	Neo4j 100%
Primo test	17,78	210,23	16,89	254,43	44,13	123,18	258,48	230,46
Tempo medio	2,15	4,99	2,11	6,82	23,63	17,24	166,46	136,73
Dev. standard	0,41	1,30	0,29	1,43	1,10	1,95	15,40	5,62
Conf. 95%	0,15	0,47	0,11	0,51	0,39	0,70	5,51	2,01





CONCLUSIONI

Dai vari test effettuati si può evincere una maggiore efficienza da parte di MongoDB per quanto riguarda la gestione e l'interrogazione di database con piccole quantità di dati; al contrario, in caso di dataset di grandi dimensioni i tempi d'esecuzione delle interrogazioni salgono notevolmente.

Neo4j è più lento nella gestione di database di piccole dimensioni e risulta, grazie alla struttura a grafo che ottimizza le interrogazioni in presenza di dataset di grandi dimensioni, più efficace con grandi quantità di dati.