

# Graph Processing with Neo4j

## Lecture BigData Analytics

Julian M. Kunkel

julian.kunkel@gmail.com

University of Hamburg / German Climate Computing Center (DKRZ)

2017-12-08



*Disclaimer: Big Data software is constantly updated, code samples may be outdated.*

# Outline

- 1 Overview
- 2 Cypher Query Language (CQL)
- 3 Interfaces
- 4 Architecture
- 5 Summary

# Neo4j [31, 32]

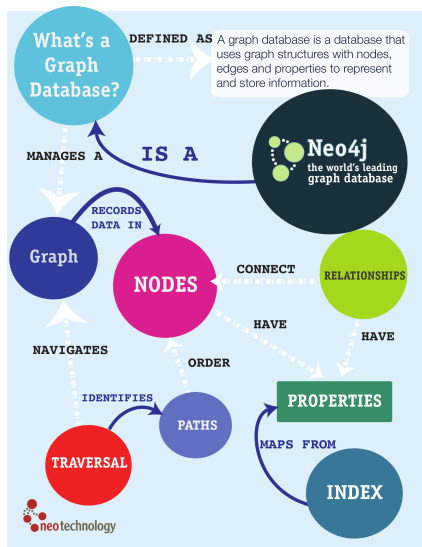
- Graph database written in Java
- Supports ACID transaction semantics
- One server scales to billions of nodes/relationships
  - Performance: Millions of node traversals/s
- High availability (and performance) through clustering
- Declarative query language Cypher (instead of SQL)
- Note: Very loose connection to Hadoop ecosystem
  - E.g., Prepare data in HBASE for batch import in Neo4j
  - Suboptimal import of Millions of nodes can take days
- Schema-optional: You can use a schema
  - To gain performance
  - To improve modeling, e.g., via constraints
- Many interfaces to the graph database

# Graph Data Model (Slight Changes for Neo4j)

- Nodes: Entity
- Edges: Relationship between two nodes
  - They have a direction and type
- Property: (key, value)
  - Attributes describe relationships/nodes
  - Key is string, the value has a datatype
- Label: Organize nodes into groups

## Definitions for queries

- Path: One or more nodes with connecting relationships
- Traversal: Navigates through a graph to find paths

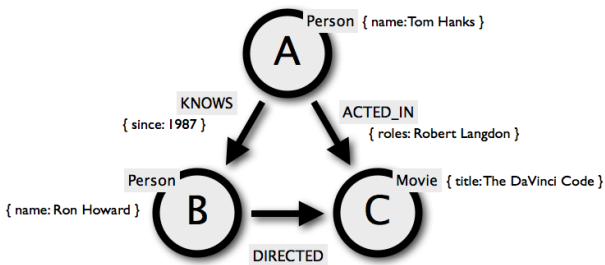


Source: What's a Graph Database [31]

# Example Graph Use-Cases

## Movie and actors data [31]

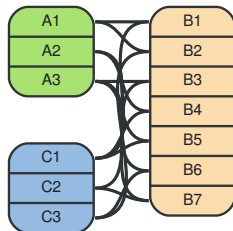
- Movies: label, title, released date, tagline
- People: label, name, born (optional date)
- Relationships
  - ACTED\_IN from actor to movie, roles (list of played chars)
  - DIRECTED from director to movie



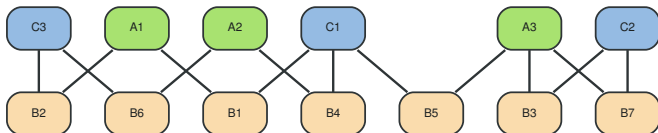
Source: Online Course: Introduction to Graph Databases and Neo4j [31]

# Converting RDBMS to Graphs

- Consider three tables A,B,C
- Relations between rows (foreign keys) become edges

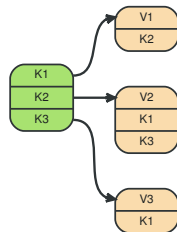


Source: RDBMS. The Neo4j Manual v2.2.5 [33]

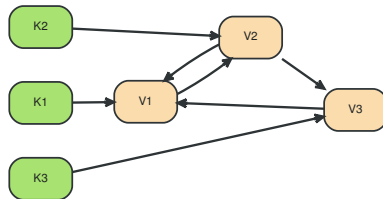


Source: Graph Database as RDBMS. The Neo4j Manual v2.2.5 [33]

# Converting Key-Value Store Models to Graphs

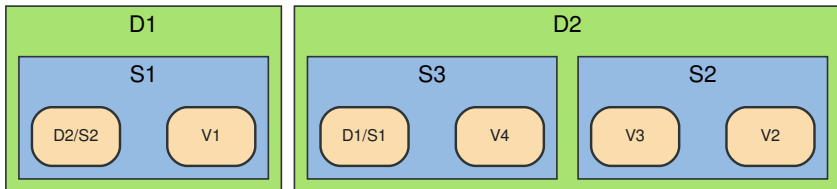


Source: Key-Value Store. The Neo4j Manual v2.2.5 [33]



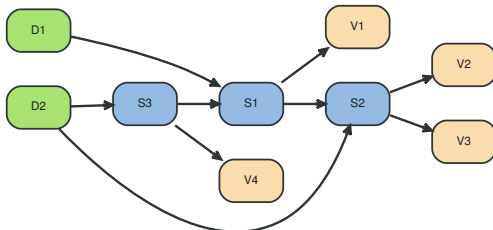
Source: Graph Database as Key-Value Store. The Neo4j Manual [33]

# Converting the Document Store Model to Graphs



Source: Document Store. The Neo4j Manual v2.2.5 [33]

D=Document, S=Subdocument, V=Value, X/Y=reference to a subdocument in another document



Source: Graph Database as Document Store. The Neo4j Manual v2.2.5 [33]



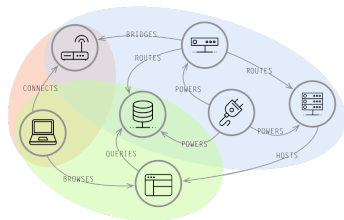
# Neo4j Case Success Studies [31]

## For the logistics company Accenture

- Use case: Dynamic parcel routing (5 million parcels/day)
- With Neo4j: Routing of packets online, i.e., where to load a parcel

## For the communication company SFR

- Use case: Prioritize hardware replacement to minimize downtime
  - Run automated “what if” analysis to ensure resilience
- With Neo4j: Loading data from > 30 systems works; easier analysis model



Source: [36]

1 Overview

**2 Cypher Query Language (CQL)**

3 Interfaces

4 Architecture

5 Summary

# Cypher Query Language Basics [31]

- Declarative query language for formulating graph queries
  - Define matches based on structural patterns
- Allows query and/or update of the graph
  - Each part of a query must be read-only or write-only
  - A query consists of multiple clauses
- Transactions can span multiple queries
- Supports: variables, expressions<sup>1</sup>, operators, comments
- Supports collections (list, dictionary)
- Provides functions for aggregation, collections, strings, math

---

<sup>1</sup>Handling missing values with NULL is possible, see <http://neo4j.com/docs/stable/cypher-working-with-null.html>

# Cypher Query Language [33]

## Syntax: specifying graph structures via patterns

### ■ Node

- Anonymous node:  $()$
- Named node:  $(x)$ , the variable  $x$  is used to refer to it
- Node with a specific label/class:  $(x : label)$

### ■ Relationship

- Named relationship:  $-[r]- >$
- Typed relationship:  $-[r : t]- >$
- Two nodes with a relationship:  $(a) - [r] - > (b)$

### ■ Properties can be specified in $\{\}$ , e.g., $(x \{name:"Hans"\})$

### ■ A pattern combines several nodes/relations

# Cypher Query Language Read Clauses [33]

- **LOAD CSV:** read data from a CSV file, can be used for importing
- **MATCH:** search for something (returns a relational table)
  - **DISTINCT** keyword: Avoid replicates (e.g., returning a node twice)
  - **OPTIONAL MATCH:** optional relationship like SQL outer join
- **WHERE:** Filtering based on properties and pattern
  - Supports regex matching of strings
  - Pattern predicates restrict the graph's shape
- **Aggregation functions**
  - **Automatic grouping on all non-aggregated columns**
  - `sum`, `avg`, `percentileDisc`, `count`
    - e.g., `count(*)`, `count(DISTINCT X)`
  - `collect(x)`: creates a list of all values

# Cypher Query Language Write Clauses [33]

- CREATE: an element or relation
- MERGE: Create or lookup (CREATE + MATCH)
- SET: Modify/Add data/labels
- REMOVE: remove labels and properties
- DELETE: remove graph elements

# Cypher Query Language: Interactive Session

```

1 # Create a star graph
2 $ CREATE (c) FOREACH (x IN range(1,6) | CREATE (l),(c)-[:X]->(l)) RETURN id(c);
3 id
4 0
5 Updated the graph - created 7 nodes and 6 relationships
6
7 # Count the number of nodes
8 $ MATCH (n) RETURN count(n); # since we have not defined any restriction, all nodes
9 count(n)
10 7
11
12 # Count relationships based on their relationship type
13 $ MATCH ()-[r]->() RETURN type(r), count(*);
14 type(r) count(*)
15 X 6
16
17 # Set the center node's name property to CENTER
18 $ MATCH (n) WHERE id(n) = 184 SET n.name = "CENTER";
19
20 # Clean the database
21 $ MATCH (n) OPTIONAL MATCH (n)-[r]-() DELETE n, r;

```

# Cypher Query Language General Clauses [33]

- FOREACH(< col >|< op >): iterates through a collection, apply op
- RETURN: return the subgraph/table
  - Usually you can convert those into a response table
- AS x: rename column to x
- ORDER BY x (ASC|DESC): sorting
- SKIP, LIMIT X: paginate – organize output
- UNION: compose statements
- WITH: a barrier for a pipeline of multiple statements
  - Example: retrieve the top entries by a criteria and join it with other data
  - Allows also to combine read-only and write-only parts
  - Aggregated results must pass through a WITH clause
- UNWIND: expand a collection into a sequence of rows
- USING: instruction to use/avoid indexes



# Cypher Query Language [33]: Selection of Functions

- `id()`: the node id
- `timestamp()`: a timestamp
- `label()`: the node label
- `upper()`, `lower()`: change case
- `range(l,u)`: return a collection with numbers from l to u
- `length(x)`: size of a collection
- `keys(x)`: keys of a dictionary
- `coalesce(x, y)`: use property x if available, else y
- `nodes(path)`, `rels(path)`, `length(path)`

# Cypher Query Language: Examples [33]

```

1 # Return a collection
2 $ RETURN [1, 2, 3]
3
4 # Return a string with a row name of X
5 $ RETURN "BigData" as X
6
7 # Return a dictionary
8 $ RETURN {key1 : 2, key2 : "test"}
9
10 # Return a list of x^3 where x is an even number
11 $ RETURN [x IN range(1,10) WHERE x % 2 = 0 | x^3] AS result
12
13 # populate a table
14 $ CREATE (matrix1:Movie { title : 'The Matrix', year : '1999-03-31' })
15 $ CREATE (keanu:Actor { name:'Keanu Reeves' })
16 $ CREATE (keanu)-[:ACTS_IN { role : 'Neo' }]->(matrix1)
17
18 # Create actor keanu if he does not exist
19 $ MERGE (keanu:Actor { name:'Keanu Reeves' })
20
21 # Eliminate duplicates from a collection
22 $ WITH [1,1,2,2] AS coll UNWIND coll AS x WITH DISTINCT x RETURN collect(x) AS SET
23 # [1,2]

```

# Cypher Query Language: Examples [33]

```

1 # Read a table from a (large) CSV
2 USING PERIODIC COMMIT
3 LOAD CSV WITH HEADERS FROM 'http://neo4j.com/docs/2.2.5/csv/artists-with-headers.csv' AS
   ↪ line
4 CREATE (:Artist { name: line.Name, year: toInt(line.Year)})
5
6 MATCH (a:Movie { title: 'Wall Street' })
7 OPTIONAL MATCH (a)-->(x)
8 RETURN x
9
10 # return a movie with all its properties
11 MATCH (movie:Movie { title: 'The Matrix' })
12 RETURN movie;
13
14 # return certain attributes
15 MATCH (movie:Movie { title: 'The Matrix' })
16 RETURN movie.title, movie.year;
17
18 # show all actors sorted by name
19 MATCH (actor:Actor)
20 RETURN actor ORDER BY actor.name;
21
22 # all actors whose name end with s
23 MATCH (actor:Actor)
24 WHERE actor.name =~ ".*s$"
25 RETURN actor.name;

```

# Cypher Query Language: Examples [33]

```

1 # List all nodes together with their relationships
2 MATCH (n)-[r]->(m) RETURN n AS from, r AS '->', m AS to;
3
4 # Return number of movies for actors acting in "The Matrix"
5 MATCH (:Movie { title: "The Matrix" })<-[:ACTS_IN]-(actor)-[:ACTS_IN]->(movie)
6 RETURN movie.title, collect(actor.name), count(*) AS count
7 ORDER BY count DESC;
8
9 # Filtering
10 MATCH (p:Person)-[r:ACTED_IN]->(m:Movie)
11 WHERE p.name =~ "K.+" OR m.released > 2000 OR "Neo" IN r.roles
12 RETURN p,r,m;
13
14 # Filtering based on graph structure
15 # Here: Search for people that are actors in any movie but never directed any movie
16 MATCH (p:Person)-[:ACTED_IN]->(m)
17 WHERE NOT (p)-[:DIRECTED]->()
18 RETURN p,m;
19
20 # Identify how often actors and directors worked together
21 MATCH (actor:Person)-[:ACTED_IN]->(movie:Movie)<-[:DIRECTED]-(director:Person)
22 RETURN actor,director,count(*) AS collaborations;

```

# Cypher Query Language: Examples [33]

```

1 # Use UNION to combine results
2 MATCH (p:Person)-[r:ACTED_IN]->(m:Movie)
3 RETURN p,type(r) AS rel,m
4 UNION
5 MATCH (p:Person)-[r:DIRECTED]->(m:Movie)
6 RETURN p,type(r) AS rel,m
7
8 # Return five actors of each movie
9 MATCH (m:Movie)<-[:ACTED_IN]-(a:Person)
10 RETURN m.title AS movie, collect(a.name)[0..5] AS five_of_cast
11
12 # Use list predicates to restrict set further, here all relations must be acted in
13 MATCH path =(:Person)-->(:Movie)<--(:Person)
14 WHERE ALL (r IN rels(path) WHERE type(r) = 'ACTED_IN') AND ANY (n IN nodes(path) WHERE
    ↪ n.name = 'Clint Eastwood')
15 RETURN path
16 # Update values based on a query
17 MATCH (n {name: 'John'})-[:FRIEND]-(friend)
18 WITH n, count(friend) as friendsCount
19 WHERE friendsCount > 3
20 SET n.friendCount = friendsCount
21 RETURN n, friendsCount
22
23 # Update all nodes on possible paths between two nodes
24 MATCH p =(begin)-[*]->(end)
25 WHERE begin.name='A' AND end.name='D'
26 FOREACH (n IN nodes(p) | SET n.marked = TRUE )

```

# Schemas [33]

- Neo4j offers a few schema options to influence graph setup
- Simple constraints can be created using CREATE

```
1 CREATE CONSTRAINT ON (p:Person) ASSERT p.name IS UNIQUE
2 # And removed with DROP
3 DROP CONSTRAINT ON (p:Person) ASSERT p.name IS UNIQUE
```

- Indexes for lookup

```
1 CREATE INDEX ON :Person(name)
2 DROP INDEX ON :Person(name)
```

1 Overview

2 Cypher Query Language (CQL)

**3 Interfaces**

4 Architecture

5 Summary

# Overview of the Interfaces

- Neo4j shell [38]
  - Create, import, export, execute Cypher
  - Present results as ASCII tables
- Web interface
  - Provides a shell for Cypher
  - Visualizes query results
  - Allows (performance) monitoring of Neo4j
  - Ships with Examples/Tutorials!
  - HTTPS support
- Java API
  - Core Java API offers graph algorithms & is faster than CQL
  - JCypher: DSL for higher abstraction level
  - Automatic object-graph mapping via annotations
- Relational mapping with JDBC driver
- REST, Python, ...



# Web Interface: Example Queries

The screenshot displays a web interface for a graph database. At the top, a query editor contains the Cypher query: `$ MATCH (n) RETURN n LIMIT 5`. To the right of the query are three icons: a star, a plus sign, and a play button. Below the query editor, a toolbar includes icons for download, pin, share, and refresh. The main area shows a graph visualization with five purple circular nodes. The nodes are labeled: `*(5)`, `Article(5)`, `WP.CO...`, `2101`, `WP.HO...`, `Nudge`, and `title stri...`. On the left side, there is a sidebar with icons for Graph, Rows, and a search icon. At the bottom, a status bar indicates "Displaying 5 nodes, 0 relationships." and an "AUTO-COMPLETE" toggle switch is set to "ON".

# Web Interface: Example Queries

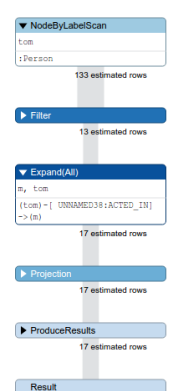
The screenshot displays a web interface for a graph database query. At the top, a Cypher query is entered:

```
$ MATCH (n:Article {title: "Bee"})-[r:links]->
(m:Article) RETURN n,r, m LIMIT 5
```

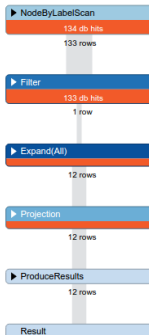
Below the query, a graph visualization is shown. The graph features several nodes representing articles: Xyloid..., Xyloidas, Xyphyr..., Bee, Wyvern, and Yale\_... The nodes are interconnected by directed edges labeled 'links'. The 'Bee' node is centrally located and has multiple outgoing links to other nodes. The 'Yale\_...' node is highlighted with a yellow border. The interface includes a sidebar with navigation icons (home, star, info) and a top bar with query execution controls (star, plus, play).

# Clauses for Debugging of Queries

- EXPLAIN: shows the execution plan
- PROFILE: runs the statement and shows where time is spend



EXPLAIN ...



Cypher version: CYPHER 2.3,  
planner: COST. 292 total db hits in 78 ms.

PROFILE ...

MATCH (tom:Person name:"Tom Hanks")-[:ACTED\_IN]->(m) RETURN m.name

# Java API: Example for our Student Table. See [37]

```
1 private static enum MyRelationTypes implements RelationshipType
2 { ATTENDS } // we can use enums for relation types
3
4 public static void main(String [ ] args){
5     GraphDatabaseService graphDb; // start database server
6     graphDb = new GraphDatabaseFactory().newEmbeddedDatabaseBuilder(File("x"));
7     registerShutdownHook( graphDb );
8
9     Node student;   Node lecture; Relationship attends;
10    // encapsulate operations into a transaction
11    try ( Transaction tx = graphDb.beginTx() ){
12        student = graphDb.createNode();
13        student.setProperty( "Name", "Hans" );
14        lecture = graphDb.createNode();
15        lecture.setProperty( "Lecture", "Big Data Analytics" );
16        attends = student.createRelationshipTo( lecture, RelTypes.ATTENDS );
17        attends.setProperty( "Semester", "1718" );
18        tx.success();
19    }
20    graphDb.shutdown(); // shutdown application server
21 }
```

1 Overview

2 Cypher Query Language (CQL)

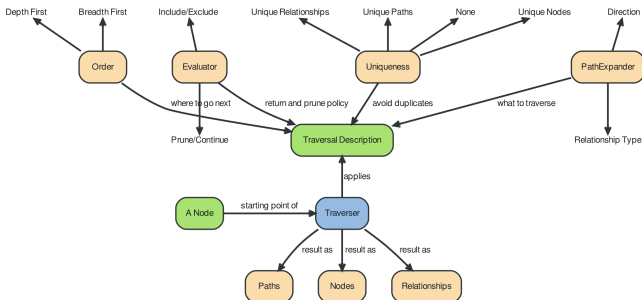
3 Interfaces

4 Architecture

5 Summary

# Evaluation of Cypher expressions [33]

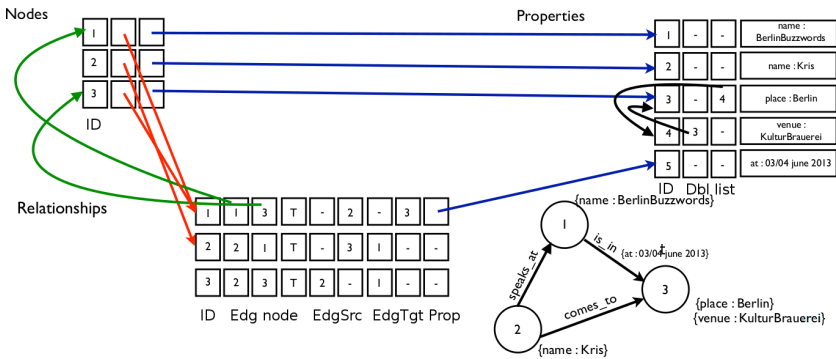
- An execution planner transforms a query into a plan
  - Rule-based planner uses indexes
  - Cost-based planner uses statistical information
- Use indices if available
- Order (DFS or BFS)
- Uniqueness: avoid duplicates
- Evaluator: decide what to return and when to stop
- Recursive matching with backtracking



Source: The Neo4j Manual 2.2.5 (36.1. Main Concepts) [33]

# Neo4j Architecture: On-Disk Format [32]

- Physically, multiple “store files” are used
- Data is stored as double linked lists of records
- Storage for nodes, relationships and properties
  - Long values are persisted in separate array and string stores



Source: K. Geusebroek. I MapReduced a Neo store [34] (modified)

# Neo4j Consistency [32]

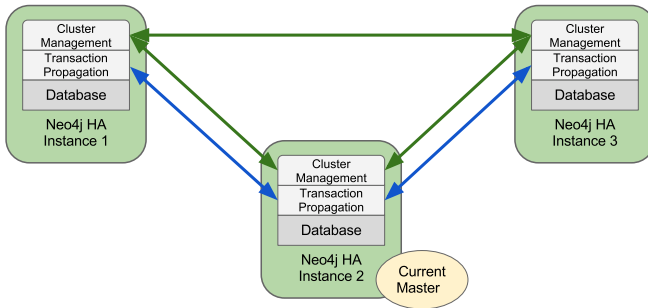
- ACID transaction support
  - Isolation of concurrent operations until transaction is completed
  - All write operations are sorted (before stored/communicated) to ensure predictable update order
  - Write changes are applied in sorted order to the transaction log
  - Apply the changes to the store files
  - Implemented via locking of Nodes/Relationships during transaction
- Upon completion of transaction, changes are persisted
- Recovery: re-applies the transaction log



# Neo4j High-Availability [32, 33, 35]

- Neo4j clustering replicates the database across servers
- One master multiple slaves provides
  - Data redundancy
  - Service fault tolerance
- A master election protocol is used
- A quorum (majority) of servers must be up to serve writes
- Transactions are first committed to master
  - Creating an incrementing transaction id (txid)
  - Eventually applied to slaves sending streams
  - Update interval defines delay
- Applying transactions to a slave
  - The master coordinates locking
  - After applying transaction on master
  - The slave uses the same txid

# Neo4j High-Availability Architecture [33]



Source: The Neo4j Manual 2.2.5 (25.1. Architecture) [33]

# Neo4j Performance Aspects [32]

- Remember: Data is completely replicated across servers
- Clustered Neo4j allows horizontal scaling of reads
- Writes are always coordinated by the master
  - Transactions can be speed up with batch inserts and periodic commits
  - The file format is optimized for graph-local operations
  - Indexing and caching speed up access
- Fine lock granularity (on node/relationship level)
- Consistency: Nodes/Relationships have an unique ID
  - Blocks for IDs are pre-allocated from the master
  - Creation of nodes/relationships does not require a lock

# Performance Aspects [32]

## Indexing

- Index: Node labels, relation type and property values
- Eventually available, populated in background
- Handled via Apache Lucene search library
- Automatic indexing possible

## Caches

- Filesystem cache: for blocks of store files
  - LFU eviction policy
  - Uses *mmap()* to map data blocks into memory
- Node/Relationship cache

# Summary

- Neo4j is a powerful graph database
- ACID transaction semantics
- Other data models can be converted to graphs
- Many interfaces for accessing graph
- CypherQL is the SQL for the Neo4j graph DB
- Interactive web interface processes CQL
- Simple file format with linked lists
- Clustering increases read scalability

# Bibliography

10 Wikipedia

31 Interactive Online Course

<http://neo4j.com/graphacademy/online-training/>

32 <http://de.slideshare.net/thobe/an-overview-of-neo4j-internals>

33 The Neo4j Manual v2.2.5. <http://neo4j.com/docs/stable/>

34 I MapReduced a Neo store. Kris Geusebroek.

<http://2013.berlinbuzzwords.de/sites/2013.berlinbuzzwords.de/files/slides/CreatingLargeNeo4jDatabasesWithHadoop.pdf>

35 D. Montag. Understanding Neo4j Scalability.

36 <http://neo4j.com/use-cases/>

37 <http://neo4j.com/docs/stable/tutorials-java-embedded-hello-world.html>

38 <http://neo4j.com/docs/stable/tools.html>