



BDA Report: Neo4j

Tobias Kneidinger, Rafael Hochedlinger

December 10, 2024

Contents

Introduction to Graph Databases	2
Core Concepts in Neo4j	3
Graph Algorithms with GDS	4
Performance and Scalability	5
Use Cases and Applications	5
References	6

Introduction to Graph Databases

Graph Database Management Systems (GDBMS) are a type of NoSQL database which - instead of traditional tables with rows and columns - uses graph structures with nodes, edges, and properties to represent and store data, as seen in Figure 1. This is especially useful for data that is highly interconnected, such as:

- Social networks
- Recommendation engines
- Fraud detection
- Network and IT operations
- Master data management
- Real-time pricing and promotions
- Identity and access management
- etc. . .

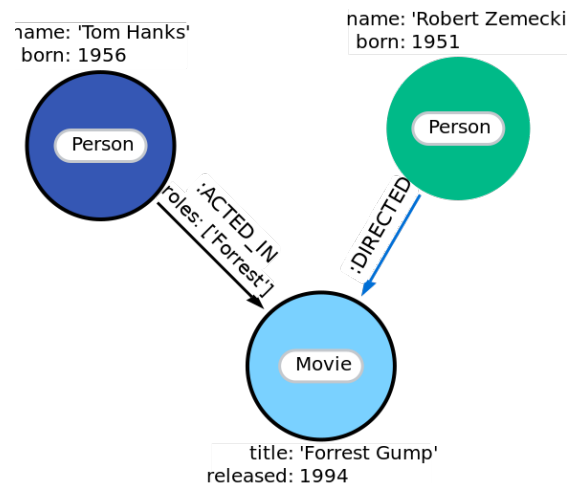


Figure 1: Graph Database Structure by Neo4j, 2024a

Graph databases are optimized for traversing and querying graph structures, which makes them very efficient for certain types of queries. As they are NoSQL databases, they are also schema-less and can store data in a flexible way. This makes them also very suitable for machine learning and data science applications, as the data can be stored in a way that is directly usable for the algorithms.

Neo4j

Neo4j is a popular GDBMS that is used by many companies and organizations. It is developed by Neo Technology and is written in Java, the first version was released in 2007. Neo4j is an open-source database, but there is also a commercial version available with additional features and support. (Neo4j, 2024a)

It uses the property graph model, which consists of the following elements, as described by Lal, 2015:

- **Nodes:** Entities that represent objects or entities in the graph and can have multiple properties and labels
- **Relationships:** Connect nodes and represent the connections between them, can have a direction and a type and can also have properties
- **Properties:** Key-value pairs that can be attached to nodes and relationships
- **Labels:** Represents the type of node, can be used to group nodes together

The main query language for Neo4j is Cypher, which is a declarative query language that is similar to SQL. These queries can be used to create, read, update, and delete data in the database. The language is more thoroughly described in the next section.

For data science and machine learning applications, Neo4j also provides the Graph Data Science (GDS) library, which contains a collection of graph algorithms that can be used to analyze and process graph data. This library is also described in more detail in a later section.

Core Concepts in Neo4j

Neo4j offers a flexible and powerful data model in which the semantic structure of the information is mapped directly in the graph (Robinson et al., 2015). The query language of Neo4j is Cypher, with which data can be intuitively described as a pattern of nodes and relationships. Cypher queries are more like the visual representation of a graph than the traditional table logic of relational systems, which allows complex connections and paths to be found efficiently.

An essential core aspect of Neo4j is *pattern matching*. Instead of linking several tables via joins, queries in Cypher are described with nodes in round brackets () and relationships in square brackets [], often with arrows to show the direction. This enables a natural form of navigation in the graph.

The query Figure 2 finds all persons who are linked to a specific company via a `WORKS_FOR` relationship. `MATCH` defines the pattern, `WHERE` restricts the results, and `RETURN` returns the desired data. This makes it easy to query even complex graph structures.

```
MATCH (p:Person)-[r:WORKS_FOR]->(c:Company)
WHERE c.name = "Example Company Ltd."
RETURN p.name, c.name
```

Figure 2: Query to find people working for Example Company Ltd.

Neo4j also provides indices and constraints to speed up queries and ensure data integrity. For example, this guarantees unique identifiers and prevents empty values. The flexible modeling makes it possible to introduce new relationships, labels or properties at any time without having to adapt an existing schema. These concepts - from cypher and indices to constraints and flexible schema - are the basis for data analyzes in Neo4j.

Graph Algorithms with GDS

To efficiently analyze and process graph data for data science and machine learning applications, the Graph Data Science (GDS) library can be used. The algorithms in this library are optimized for graph data and can be used to perform various tasks, as described by Neo4j, 2024b:

- Community detection
- Centrality analysis
- Link prediction
- Similarity analysis
- Pathfinding
- etc...

To be able to use conventional Neo4j graphs with the GDS library, the graph must be projected into a named graph. This projected graph is a subgraph of the original graph and can be created with the query shown in Figure 3, where 'socialgraph' is the name of the newly projected graph, 'Person' is the node label to project, and 'FRIENDS' is the relationship type to project.

```
CALL gds.graph.create('socialgraph', 'Person', 'FRIENDS')
```

Figure 3: Projected graph creation

After the graph has been projected, the GDS algorithms can be used to analyze the graph. In the example shown in Figure 4, the Louvain algorithm is used to detect communities in the graph. The algorithm is called with the 'socialgraph' as the graph name, the 'Person' node label, and the 'FRIENDS' relationship type. The algorithm yields `nodeId` and `communityId`, which represent the node name and community ID. The results are then ordered by `communityId`.

```
CALL gds.louvain.stream('socialgraph')
YIELD nodeId, communityId
RETURN gds.util.asNode(nodeId).name AS nodeName, communityId
ORDER BY communityId;
```

Figure 4: GDS algorithm execution

Instead of `stream` the algorithm can also be executed with `write` to write the results back to the graph. The results can then be used for further analysis or visualization and don't have to be recalculated every time.

Performance and Scalability

Neo4j uses various techniques to ensure both performance and scalability (“Neo4j Operations Manual - Performance”, 2023). By using indexes, constraints and intelligent cache strategies, Neo4j can process queries efficiently. In addition, the use of graph algorithms directly in the database makes it possible to perform complex calculations faster without having to export data. The asynchronous processing of queries as well as optimized storage and traversal strategies ensure high performance with growing data volumes.

For the operation of distributed systems, Neo4j supports cluster configurations in which multiple instances of the database work together. This increases both reliability and throughput performance. The causal consistency used here ensures that write and read operations can be traced in a logical sequence, which guarantees consistent data processing, especially in distributed environments.

Use Cases and Applications

As specified in the introduction, graph databases are especially useful for data that is highly interconnected. This makes them very suitable for a wide range of applications, as described by Neo4j, 2024a:

- **Social networks:** To store and analyze user data, relationships, and interactions
- **Recommendation engines:** To provide personalized recommendations based on user behavior
- **Fraud detection:** To detect fraudulent activities and patterns in financial transactions
- **Network and IT operations:** To monitor and analyze network traffic and IT infrastructure
- **Master data management:** To store and manage master data and relationships between entities
- **Real-time pricing and promotions:** To calculate and provide real-time pricing and promotions
- **Identity and access management:** To manage and analyze user identities and access rights

The flexibility and scalability of graph databases make them very versatile and suitable for many different use cases, as they can also be combined with other technologies and databases such as relational databases, document databases, or key-value stores.

However, not all use cases are suitable for graph databases. If the data is changing frequently and the relationships don't have to be queried often, another technology - like a relational database - might be the better choice. Also, if the data is not highly interconnected and the queries are simple, a key-value store or document database might be more efficient.

References

- Lal, M. (2015). *Neo4j graph data modeling*. Packt Publishing. <https://books.google.it/books?id=WTlECgAAQBAJ>. (Cit. on p. 3)
- Neo4j operations manual - performance [Accessed: 2024-12-09]. (2023). (Cit. on p. 5).
- Robinson, I., Webber, J., & Eifrem, E. (2015). *Graph databases: New opportunities for connected data*. “ O’Reilly Media, Inc.” (Cit. on p. 3).
- Neo4j. (2024a). *Neo4j* [Accessed: 2024-12-02]. <https://neo4j.com/>. (Cit. on pp. 2, 5)
- Neo4j. (2024b). *Neo4j graph data science documentation* [Accessed: 2024-12-02]. <https://neo4j.com/product/graph-data-science/>. (Cit. on p. 4)