

Lecture 02

Relationships, Graph Databases, Neo4J

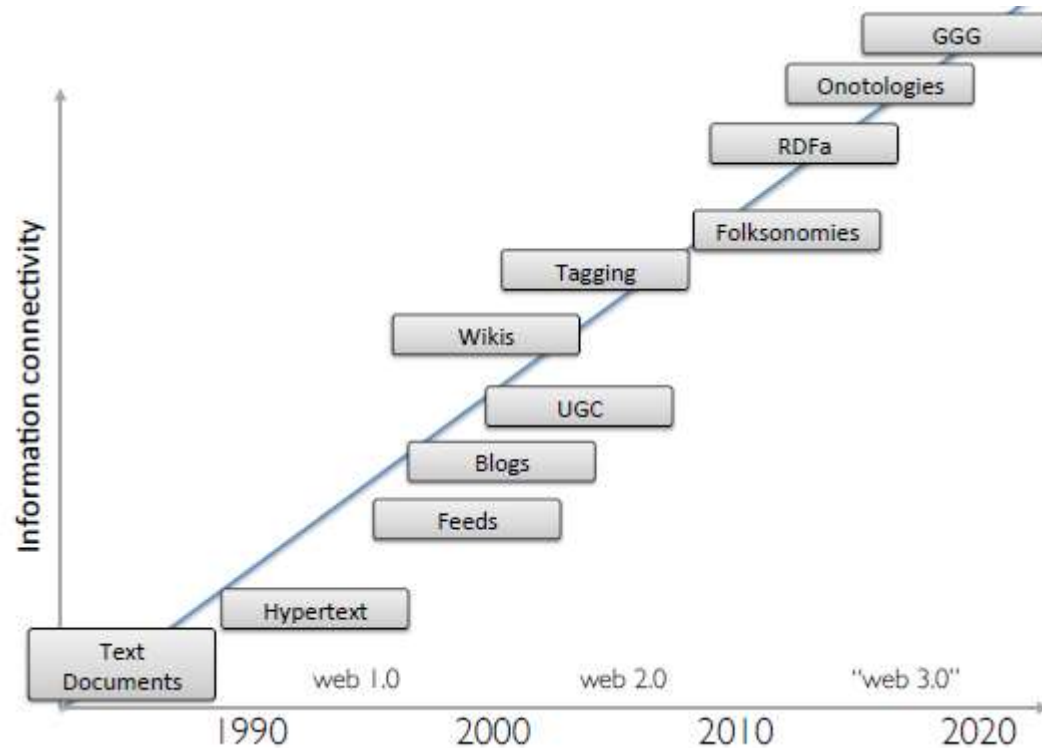
Zoran B. Djordjević

References

- These slides follow to a great measure Neo4J v3.1.1 User Manual:
<http://neo4j.com/docs/developer-manual/>

Why Graph Databases

- Data sizes are growing almost exponentially
- Large fraction of data exhibits a high degree of connectedness:



- Data are much more complex and relationships much more numerous.

RDBMS and other NoSQLs do not perform

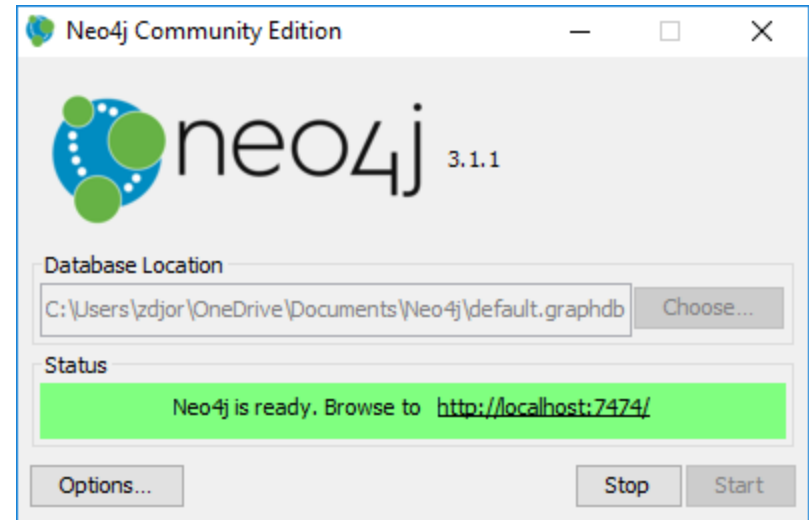
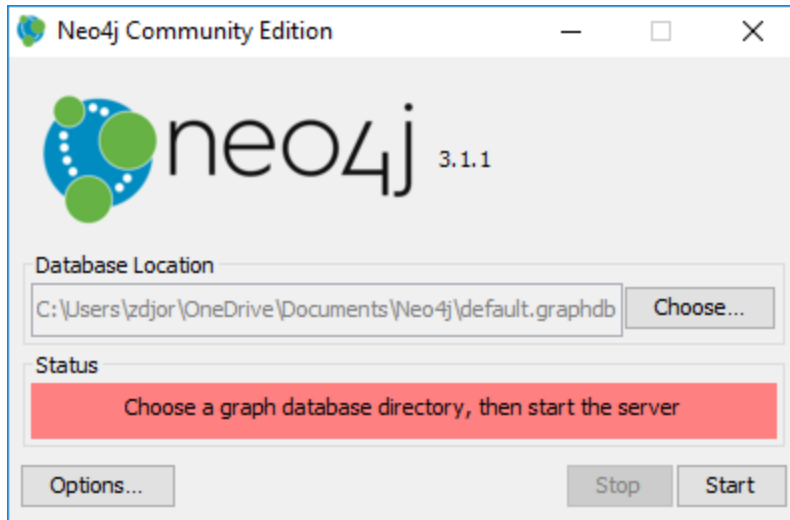
Social Network “path exists” Performance

- Experiment:
 - ~1k persons
 - Average 50 friends per person
 - `pathExists(a,b)` limited to depth 4
 - Caches warm to eliminate disk IO

	# persons	query time
Relational database	1000	2000ms
Neo4j	1000	2ms
Neo4j	1000000	2ms

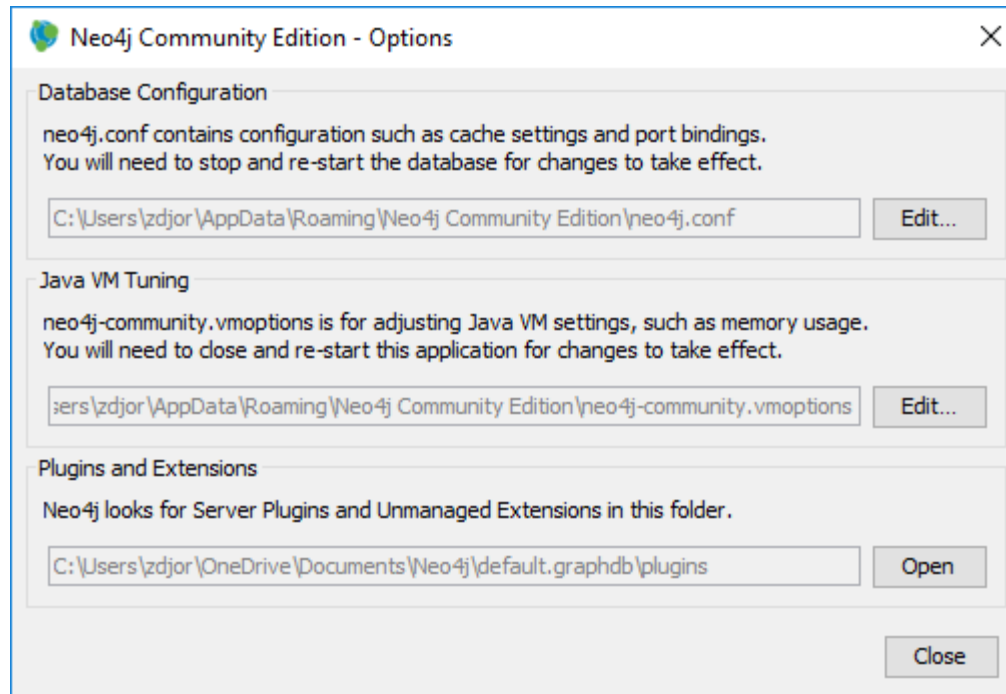
Installing Neo4J Community Edition

- Go to <https://neo4j.com/download/?ref=home>
- Select your operating system: Mac OS, Linux, Windows.
- Download `neo4j-community_windows-x64_3_2_3.exe` or similar
- Run the installation
- Select the directory
- Start



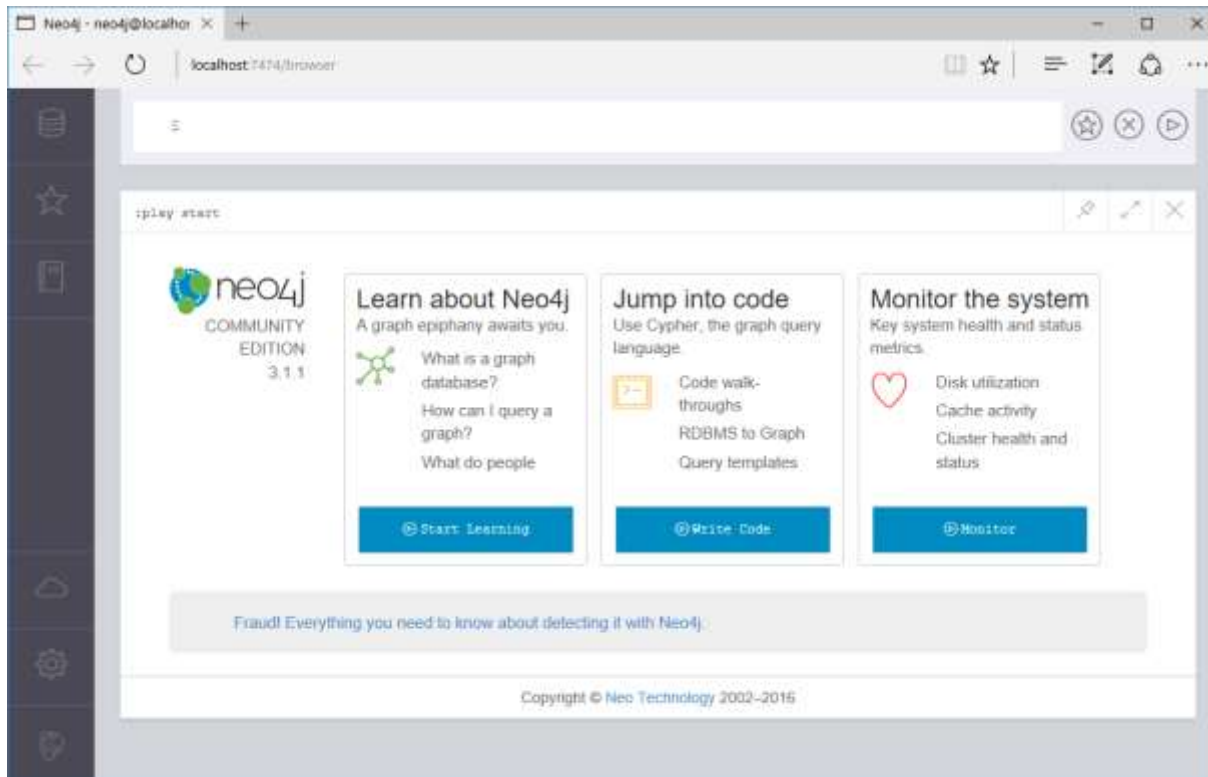
Properties files

- I did not pay attention when installing Neo4J so my files ended up in the default locations.



Port 7474, user/password: neo4j/neo4j

- Change your password on first login



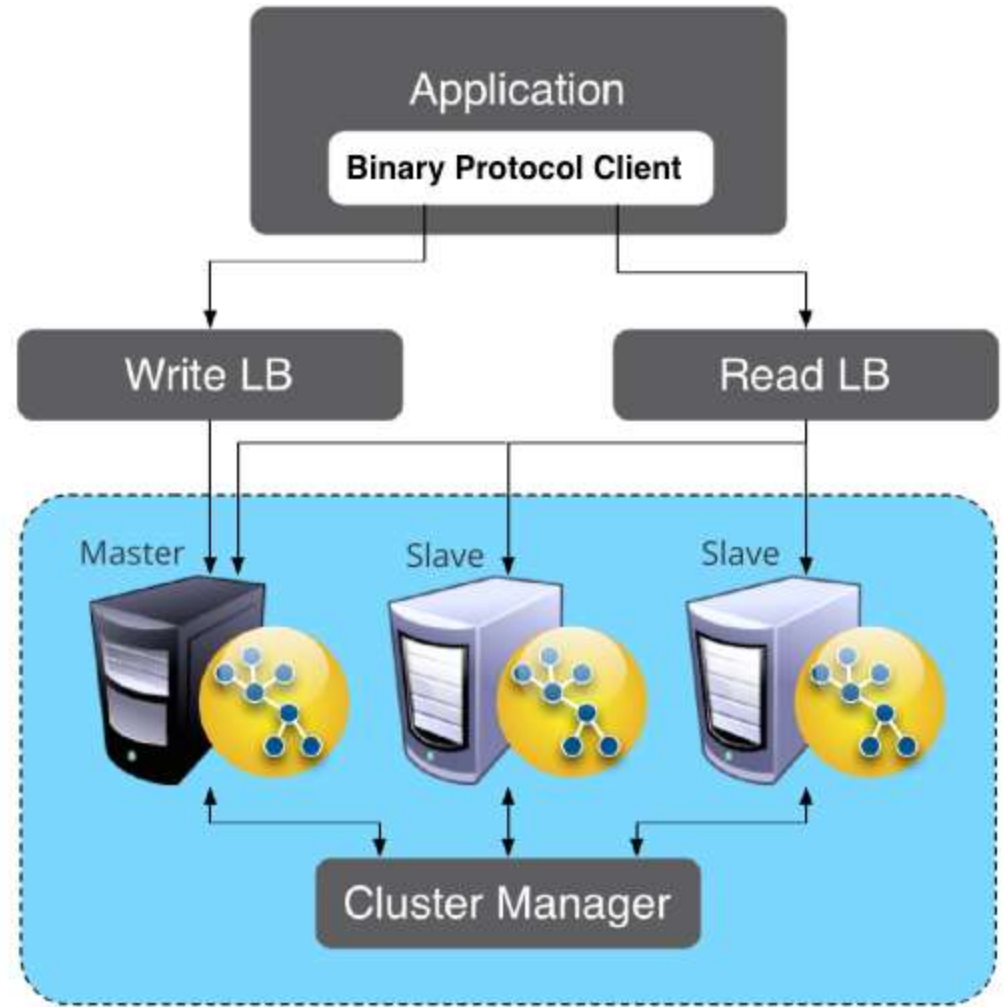
- Enter `:help` for a list of help topics. The command is preceded by ":" (colon)
- Use single line editing for brief queries or commands
- **Switch to multi-line editing with <shift-enter>**
- Run a query with <ctrl-enter>
- History is kept for easily retrieving previous commands (Ctrl –up/down arrow)

Neo4J

- As a robust, scalable and high-performance database, Neo4j is suitable for full enterprise deployment.
- Neo4J features:
 - true **ACID (Atomicity, Consistency, Isolation, Durability)** transaction properties
 - high availability,
 - scales to billions of nodes and relationships,
 - high speed querying through traversals,
 - declarative graph query language.
- Proper ACID behavior is the foundation of data reliability. Neo4j enforces that all operations that modify data occur within a transaction, guaranteeing consistent data. This robustness extends from single instance embedded graphs to multi-server high availability installations.
- Reliable graph storage can easily be added to any application. A graph can scale in size and complexity as the application evolves, with little impact on performance. Whether starting new development, or augmenting existing functionality, Neo4j is only limited by physical hardware.
- A single server instance can handle a graph of billions of nodes and relationships. When data throughput is insufficient, the graph database can be distributed among multiple servers in a high availability configuration
- The graph database storage shines when storing richly-connected data. Querying is performed through traversals, which can perform millions of traversal steps per second. A traversal step resembles a *join* in a RDBMS.

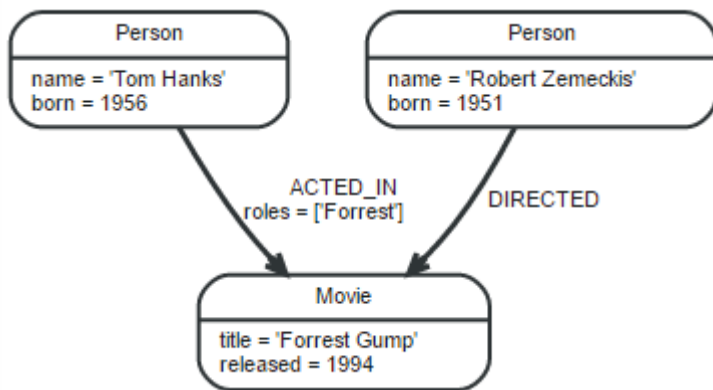
Scalability

- Until very recently Neo4J could only scale vertically. If you wanted more power, you could buy a more powerful machine, increase the number of CPU-s, increase memory and such.
- With release 3.x.x. Neo4J could scale over large clusters of machines meaning that it truly became a Big Data tool.



Key Concepts

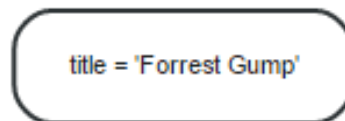
- Key concepts in a graph database are:
 - Nodes
 - Relationships
 - Properties
 - Labels
 - Traversal
 - Paths
 - Schema
- A graph records data in nodes and relationships. Both can have properties.
- This is sometimes referred to as the *Property Graph Model*.
- A graph database stores data in a graph, the most generic of data structures, capable of elegantly representing any kind of data in a highly accessible way. The following is an example of a simple graph:



- This graph contains three nodes, two with label Person and one with Label Movie.
- One person is related to the move since he DIRECTED it and the other ACTED_IN it.
- Nodes and relationships have properties.

Nodes

- The fundamental units that form a graph are nodes and relationships. In Neo4j, both nodes and relationships can contain properties.
- Nodes are often used to represent entities, but depending on the domain relationships may be used for that purpose as well.
- Apart from properties and relationships, nodes can also be labeled with zero or more labels.
- The simplest possible graph is a single Node. A Node can have zero or more named values referred to as properties.
- One graph could have a single node with a single property named title:

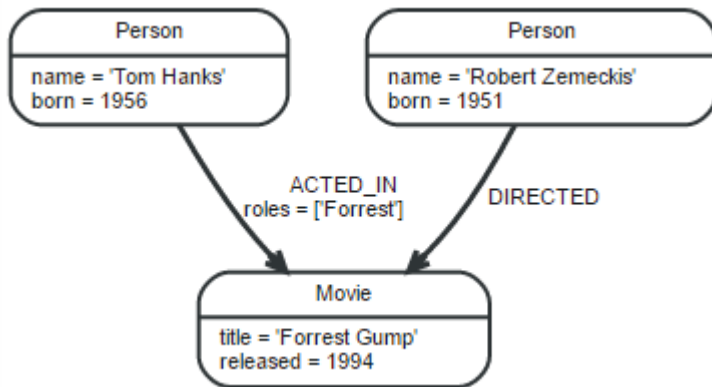


- More complex graphs could have two or more nodes. One could extend the previous graph with two more nodes and one more property on the first node:



Relationships

- Relationships organize the nodes by connecting them. A relationship connects two nodes — a start node and an end node. Just like nodes, relationships can have properties.
- Relationships between nodes are a key part of a graph database. They allow for finding related data. Just like nodes, relationships can have properties.
- A relationship connects two nodes, and is guaranteed to have valid start and end nodes.
- Relationships organize nodes into arbitrary structures, allowing a graph to resemble a list, a tree, a map, or a compound entity — any of which can be combined into yet more complex, richly inter-connected structures.



- Graph to the left uses ACTED_IN and DIRECTED as relationship types.
 - The roles property on the ACTED_IN relationship has an array value with a single item in it.
 - ACTED_IN relationship has Tom Hanks node as start node and Forrest Gump as end node.
 - Tom Hanks node has an outgoing relationship, while the Forrest Gump node has an incoming relationship
- Relationships are equally well traversed in either direction.
 - There is little need to add duplicate relationship in the opposite direction.

Relationships

- While relationships always have a direction, you can ignore the direction where it is not useful in your application. A node could have relationships to itself as well.
- We use relationship direction and type to extract information:

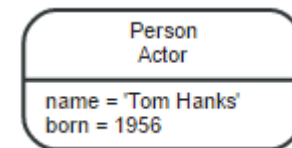
What we want to know	Start from	Relationship type	Direction
get actors in movie	movie node	ACTED_IN	incoming
get movies with actor	person node	ACTED_IN	outgoing
get directors of movie	movie node	DIRECTED	incoming
get movies directed by	person node	DIRECTED	outgoing

Properties

- Both nodes and relationships can have properties.
- Properties are named values where the name is a string. The supported property values are:
 - Numeric values,
 - String values,
 - Boolean values,
 - Collections of any other type of value.
- NULL is not a valid property value.
- NULLs can instead be modeled by the absence of a key.

Labels

- Labels assign roles or types to nodes.
- A label is a named graph construct that is used to group nodes into sets; all nodes labeled with the same label belongs to the same set. Many database queries can work with these sets instead of the whole graph, making queries easier to write and more efficient to execute. A node may be labeled with any number of labels, including none, making labels an optional addition to the graph.
- Labels are used when defining constraints and adding indexes for properties
- An example would be a label named User that you label all your nodes representing users with. With that in place, you can ask Neo4j to perform operations only on your user nodes, such as finding all users with a given name.
- However, you can use labels for much more. Labels can be added and removed during runtime, and can be used to mark temporary states for your nodes. You might create an Offline label for phones that are offline, a Happy label for happy pets, and so on.
- Person and Movie are two labels in our graph.
- A node could have two or more labels.
- Tom Hanks node could be labeled with Person and Actor

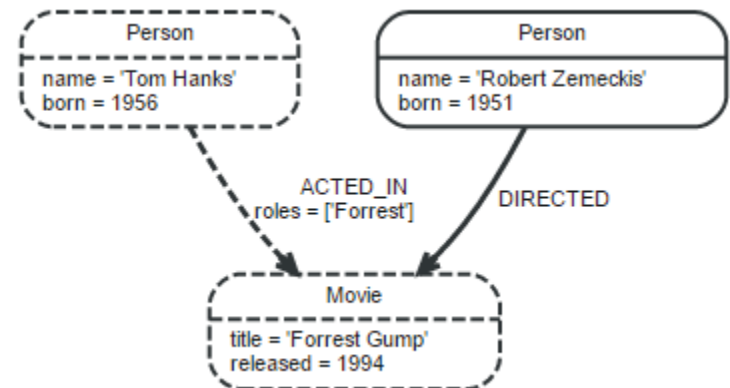


Label Names

- Any non-empty Unicode string can be used as a label name.
- In Cypher, you may need to use the backtick (``) syntax to avoid clashes with Cypher identifier rules or to allow non-alphanumeric characters in a label.
- By convention, labels are written with CamelCase notation, with the first letter in upper case. For instance, User or CarOwner.
- Labels have an id space of an `int`, meaning the maximum number of labels the database can contain is roughly 2 billion.

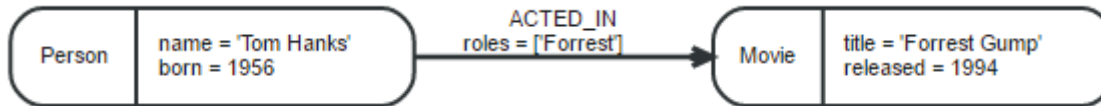
Traversal

- A traversal navigates through a graph to find paths.
- A traversal is how you query a graph, navigating from starting nodes to related nodes, finding answers to questions like "what music do my friends like that I don't yet own," or "if this power supply goes down, what web services are affected?"
- Traversing a graph means visiting its nodes, following relationships according to some rules. In most cases only a subgraph is visited, as you already know where in the graph the interesting nodes and relationships are found.
- Cypher provides a declarative way to query the graph powered by traversals and other techniques. Use Cypher Query Language.
- When writing server plugins or using embedded Neo4j, Neo4j provides a callback based traversal API which lets you specify the traversal rules. At a basic level there's a choice between traversing breadth- or depth-first.
- To find out which movies Tom Hanks acted, start the traversal from the Tom Hanks node, follow any ACTED_IN relationships connected to the node, and end up with Forrest Gump as the result (see the dashed lines)



Paths

- A path is one or more nodes with connecting relationships, typically retrieved as a query or traversal result.
- In the previous example, the traversal result could be returned as a path



- The path above has length one.
- The shortest possible path has length zero — that is it contains only a single node and no relationships.
- Self referenced node has path of length one.



Schema

- Neo4j is a schema-optional graph database.
- You can use Neo4j without any schema. Optionally you can introduce schema in order to gain performance or modeling benefits. This allows a way of working where the schema does not get in your way until you are at a stage where you want to reap the benefits of having one.

Indexes

- Performance is gained by creating indexes, which improve the speed of looking up nodes in the database.
- Once you've specified which properties to index, Neo4j will make sure your indexes are kept up to date as your graph evolves. Any operation that looks up nodes by the newly indexed properties will see a significant performance boost.
- Indexes in Neo4j are "eventually available". That means that when you first create an index the operation returns immediately. The index is populating in the background and so is not immediately available for querying. When the index has been fully populated it will eventually come online. That means that it is now ready to be used in queries.
- If something should go wrong with the index, it can end up in a failed state. When it is failed, it will not be used to speed up queries. To rebuild it, you can drop and recreate the index. Look at logs for clues about the failure.
- You can track the status of your index by asking for the index state through the API you are using. Note, however, that this is not yet possible through Cypher.

Constraints

- Neo4j can help you keep your data clean. It does so using constraints, that allow you to specify the rules for what your data should look like. Any changes that break these rules will be denied.

- As of version 3.2.x there are 4 types of constraints:

- **Unique constraint** makes sure that your database will never contain more than one node with a specific label and one property value. Use the IS UNIQUE syntax:

```
CREATE CONSTRAINT ON (book:Book) ASSERT book.isbn IS UNIQUE
```

- **Node property existence constraint** makes sure that all nodes with a certain label have a certain property, use the ASSERT exists(variable.propertyName) syntax.

```
CREATE CONSTRAINT ON (book:Book) ASSERT exists(book.isbn)
```

- **Relationship property existence constraint** makes sure that all relationships with a certain type have a certain property, use the ASSERT exists(variable.propertyName) syntax.

```
CREATE CONSTRAINT ON ()-[like:LIKED]-() ASSERT exists(like.day).
```

- **Node Key** ensuring that all nodes with a particular label have a set of defined properties whose combined value is unique, and where all properties in the set are present, use the ASSERT (variable.propertyName_1, ..., variable.propertyName_n) IS NODE KEY syntax.

```
CREATE CONSTRAINT ON (n:Person) ASSERT (n.firstname, n.surname) IS NODE KEY
```

- All constraints can be removed by using DROP CONSTRAINT command, for example:

```
DROP CONSTRAINT ON (n:Person) ASSERT (n.firstname, n.surname) IS NODE KEY
```

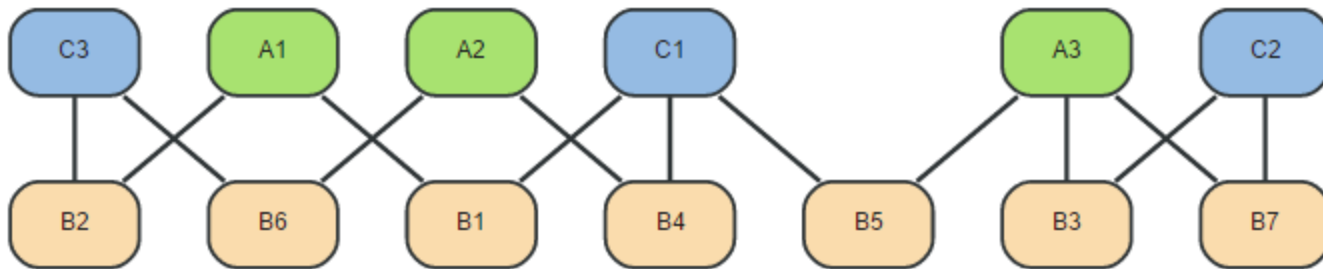
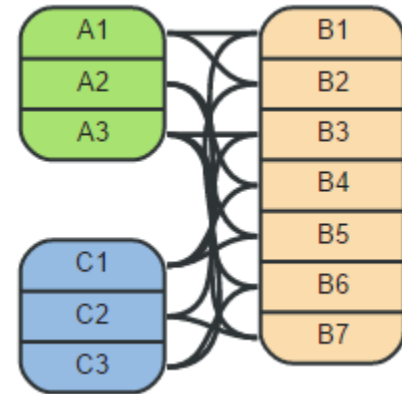
- One can impose other constraints through different APIs.

Release 3.2.x

- Neo4j 3.0, came with a "completely redesigned architecture."
- 3.0 introduced a new binary protocol dubbed Bolt, intended for speedier graph access, though the architectural redesign is centered on a new data store, where dynamic pointer compression is intended to expand Neo4j's available address space as needed.
- Bolt uses binary encoding over TCP or web sockets to snap up higher throughput and lower latency, and comes with built-in TLS which is enabled by default. Language drivers for JavaScript, Python, Java and .NET exist.
- Neo4j offers a companion cloud service, Neo4j Browser Sync, which – at no cost – allow developers to save and synchronize their favorite scripts and settings.
- Neo4j has support for Java Stores Procedures, allowing schema introspection to be added to the database, loading data from external sources, running graph algorithms and more, according to a canned statement. This combines with Bolt and Cypher – the query language for Neo4j – whose keywords ENDS, WITH, and CONTAINS are now also index-backed.
- A new cost-based query optimizer now supports write and read queries, and adds a parallel indexes capability to populate indexes.
- Neo4j claims that its custom are using graphs with hundreds-of-billions of nodes.

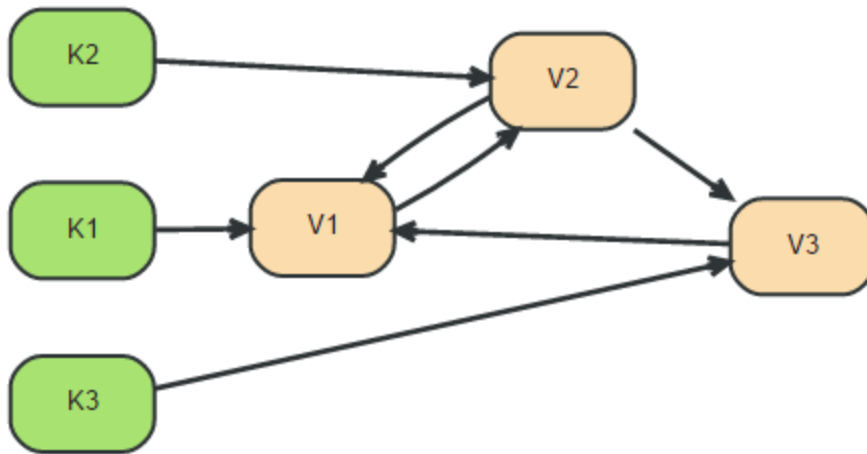
Graph DBs vs. RDBMS

- Relational databases are optimized for aggregate data.
- Graph databases are optimized for connected data.
- Extraction of information in relational databases depends on so called primary-foreign key relationships. Those relationships could also have properties.
- Deep navigation through RDBMS models require multiple joins which are expensive performance wise.
- Graph databases are optimized for localizing and returning deep and branched graph objects.



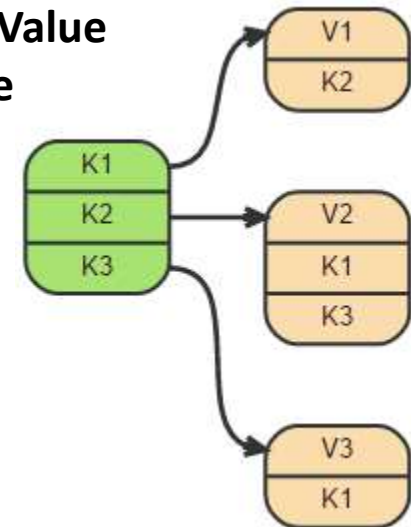
Graph DBs vs. Key-Value Stores

- A Key-Value model is great for lookups of simple values or lists. When the values are themselves interconnected, we get a graph. Neo4j lets you elaborate the simple data structures into more complex, interconnected data.
- K^* represents a key, V^* a value. Note that some keys point to other keys as well as plain values.



Graph Database as Key-Value Store

Key-Value Store

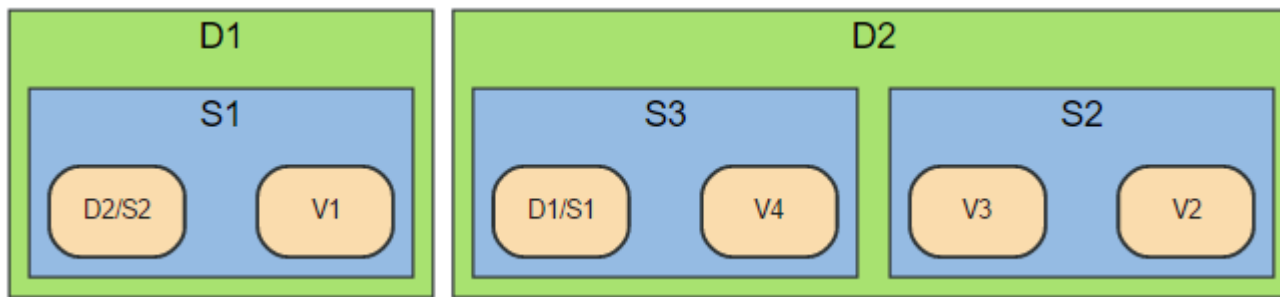


Graph Database vs. Column-Family

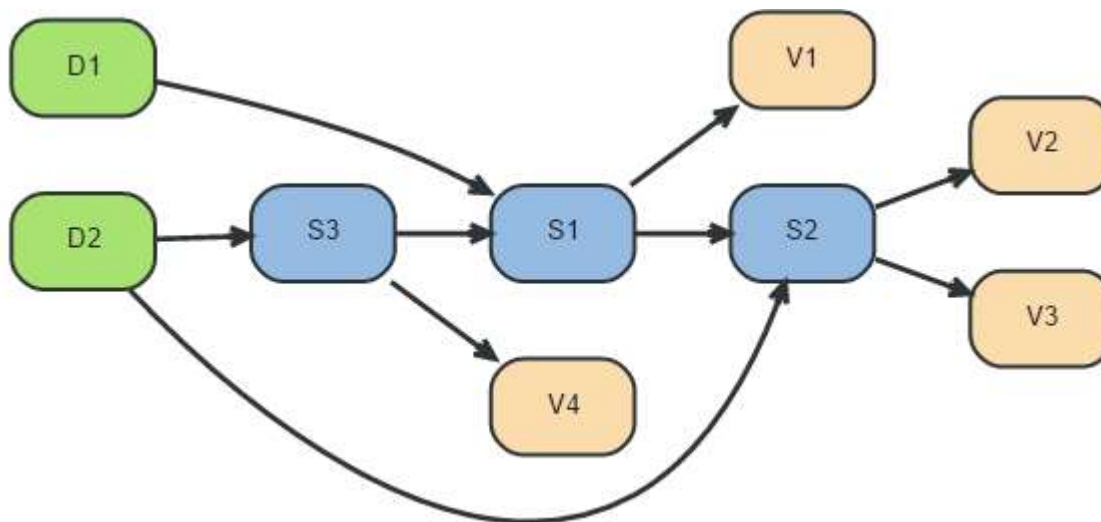
Column Family (BigTable-style) databases are an evolution of key-value, using "families" to allow grouping of rows. Stored in a graph, the families could become hierarchical, and the relationships among data becomes explicit.

Graph Database vs. Document Store

- The container hierarchy of a document database accommodates schema-free data, documents. A document could reference other documents (or document elements). Mutually referencing documents can be represented as a tree. That tree is a more expressive representation of the same data. When in Neo4j, those trees are graphs and relationships are easily navigable.



D=Document,
S=Subdocument,
V=Value,
D2/S2 = reference to
subdocument in
(other) document



**Graph Database as
Document Store**

Cypher

- Besides programming API-s for high level languages, users of Neo4j need a convenience tool for inspection and manipulation of objects (graphs) stored in the database. Cypher is a language and a tool for ad hoc analysis.
- Cypher provides a convenient way to express queries and other Neo4j actions. Although Cypher is particularly useful for exploratory work, it is fast enough to be used in production.
- Cypher performs the following:
 - Parse and validate the query.
 - Generate the execution plan.
 - Locate the initial node(s).
 - Select and traverse relationships.
 - Change and/or return values.
- Query Preparation is accomplished by Parsing and validating the query. Subsequent steps include generating an optimal search strategy.
- The execution plan must tell the database how to locate initial node(s), select relationships for traversal, etc. This involves complex optimization problems (e.g., which actions should happen first),

Cypher Concepts

- Like SQL (used in relational databases), Cypher is a textual, declarative query language. It uses a form of ASCII art to represent graph-related patterns.
- SQL-like clauses and keywords (eg, MATCH, WHERE, DELETE) are used to combine these patterns and specify desired actions.
- This combination tells Neo4j which patterns to match and what to do with the matching items (e.g., nodes, relationships, paths, collections).
- As a declarative language, Cypher does *not* tell Neo4j how to find nodes, traverse relationships, etc. (This level of control is available from Neo4j's Java APIs.

Locate Initial node, Traverse

Location of Initial Node

- Neo4j is highly optimized for traversing property graphs. Under ideal circumstances, it can traverse millions of nodes and relationships per second, following chains of pointers in the computer's memory.
- However, before traversal can begin, Neo4j must know one or more starting nodes. Unless the user (or, more likely, a client program) can provide this information, Neo4j will have to search for these nodes.
- A "brute force" search of the database can be *very* time consuming. Every node must be examined to see if it has the property, then to see if the value meets the desired criteria.
- To avoid this effort, Neo4j creates and uses indexes. Neo4j uses a separate index for each label/property combination.

Traversal and actions

- Once the initial nodes are determined, Neo4j can traverse portions of the graph and perform any requested actions.
- The execution plan helps Neo4j to determine which nodes are relevant, which relationships to traverse, etc.

Nodes, Relationships and Patterns

- Nodes and relationships are simply low-level building blocks. The real strength of the Property Graph lies in its ability to encode *patterns* of connected nodes and relationships.
- A single node or relationship typically encodes very little information, but a pattern of nodes and relationships can encode arbitrarily complex ideas.
- Cypher, Neo4j's query language, is strongly based on patterns. Patterns are used to match desired graph structures.
- A simple pattern, which has only a single relationship, connects a pair of nodes (or, occasionally, a node to itself). For example, *a Person LIVES_IN a City* or *a City is PART_OF a Country*.
- Complex patterns, using multiple relationships, can express arbitrarily complex concepts and support a variety of interesting use cases.
- For example, we might want to match instances where *a Person LIVES_IN a Country*. The following Cypher code combines two simple patterns into a (mildly) complex pattern which performs this match:

```
(:Person) -[:LIVES_IN]-> (:City) -[:PART_OF]-> (:Country)
```
- Pattern recognition is fundamental to the way that the brain works. Humans are very good at working with patterns.
- When patterns are presented visually (e.g., in a diagram or map), humans can use them to recognize, specify, and understand concepts. As a pattern-based language, Cypher takes advantage of this capability.

Node Syntax

- Cypher uses a pair of parentheses (usually containing a text string) to represent a node, e.g.: `()`, `(foo)`.
`()`
`(matrix)`
`(:Movie)`
`(matrix:Movie)`
`(matrix:Movie {title: "The Matrix"})`
`(matrix:Movie {title: "The Matrix", released: 1997})`
- The simplest form, `()`, represents an anonymous, uncharacterized node. If we want to refer to the node elsewhere, we can add an identifier, e.g.: `(matrix)`.
- Identifiers are restricted (i.e., scoped) to a single statement: an identifier may have different (or no) meaning in another statement.
- The `Movie` label (prefixed in use with a colon, `:`) declares the node's type. This restricts the pattern, keeping it from matching (say) a structure with an `Actor` node in this position.
- Neo4j's node indexes also use labels: each index is specified to the combination of a label and a property.
- The node's properties (e.g., `title`) are represented as a list of key/value pairs, enclosed within a pair of braces, eg: `{...}`.
- Properties can be used to store information and/or restrict patterns. For example, we could match nodes whose `title` is `"The Matrix"`.

Relationship Syntax

- Cypher uses a pair of dashes (--) to represent an undirected relationship. Directed relationships have an arrowhead at one end (eg, <--, -->). Bracketed expressions (eg: [...]) can be used to add details. This may include identifiers, properties, and/or type information, e.g.:

```
-->
```

```
-[role]->
```

```
-[:ACTED_IN]->
```

```
-[role:ACTED_IN]->
```

```
-[role:ACTED_IN {roles: ["Neo"]}]->
```

- The syntax and semantics found within a relationship's bracket pair are very similar to those used between a node's parentheses.
- An identifier (e.g., `role`) can be defined, to be used elsewhere in the statement.
- The relationship's type (e.g., `ACTED_IN`) is analogous to the node's label.
- Like with node labels, the relationship type `ACTED_IN` is added as a symbol, prefixed with a colon: `:ACTED_IN`
- Relationship properties (e.g., `roles`) are entirely equivalent to node properties. (Note that the value of a property may be an array.) Identifiers (e.g., `role`) can be used elsewhere in the statement to refer to the relationship.

Pattern Syntax, Clauses

- Combining the syntax for nodes and relationships, we can express patterns.

```
(keanu:Person{name: "Keanu Reeves"} )  
-[role:ACTED_IN {roles: ["Neo"] } ]->  
(matrix:Movie {title: "The Matrix"} )
```

- Cypher allows patterns to be assigned to identifiers. This allow the matching paths to be inspected, used in other expressions, etc.

```
acted_in = (:Person)-[:ACTED_IN]->(:Movie)
```

- The `acted_in` variable would contain two nodes and the connecting relationship for each path that was found or created. There are a number of functions to access details of a path, including `nodes(path)`, `rels(path)` (same as `relationships(path)`), and `length(path)`.

Clauses

- Cypher statements typically have multiple *clauses*, each of which performs a task,
 - create and match patterns in the graph
 - filter, project, sort, or paginate results
 - connect/compose partial statements
- By combining Cypher clauses, we can compose more complex statements that express what we want to know or create.
- Neo4j then figures out how to achieve the desired goal in an efficient manner.

Creating Data

- You can create individual nodes and relationships:

```
CREATE (:Movie { title:"The Matrix",released:1997 })
```

Added 1 label, created 1 node, set 2 properties, statement executed in 353 ms

- We could ask for the node to be returned to us, as well:

```
CREATE (p:Person { name:"Keanu Reeves", born:1964 }) RETURN p;
```

- Cypher browser tool would return:

The screenshot shows the Cypher browser interface. At the top, the Cypher query is entered: `$ CREATE (p:Person { name:"Keanu Reeves", born:1964 }) RETURN p`. Below the query, the results are displayed. On the left, there is a sidebar with icons for 'Graph', 'Rows', and 'Code'. The main area shows a graph view with a single green circular node labeled 'Keanu Reeves'. Below the graph, there is a table view showing the result of the query. The table has one row with the following columns: 'Person', '<id>: 1', 'born: 1964', and 'name: Keanu Reeves'.

Person	<id>: 1	born: 1964	name: Keanu Reeves
Person	<id>: 1	born: 1964	name: Keanu Reeves

Creating Multiple Elements

- If we want to create more than one element, we can separate the elements with commas or use multiple CREATE statements without comma separation
- We can also create more complex structures, including an ACTED_IN relationship with information about the character, or DIRECTED relationship for the director.

```
CREATE (a:Person { name:"Tom Hanks",  
born:1956 })-[r:ACTED_IN { roles: ["Forrest"]}]>(m:Movie { title:"Forrest  
Gump",released:1994 })  
CREATE (d:Person { name:"Robert Zemeckis", born:1951 })-[:DIRECTED]>(m)  
RETURN a,d,r,m
```



</> Code Introspection

- If you click on </> symbol you will get an insight into traffic to the server and back

\$ CREATE (a:Person { name:"Tom Hanks", born:1956 })-[r:ACTED_IN { roles: ["Forrest"]}]>(m...

Graph

Rows

</>
Code

▼ Request

Header	Value
Accept	application/json, text/plain, */*
X-stream	true
Content-Type	application/json; charset=utf-8
Authorization	Basic bmVvNGo6YWRTaW4=

Payload

```
{  "statements": []}
```

▼ Response

Header	Value
Location	http://localhost:7474/db/data/transaction/7
Date	Fri, 08 Apr 2016 13:50:03 GMT
Server	Jetty(9.2.z-SNAPSHOT)
Access-	*

Request finished in 70 ms.

▼ Request


Header	Value
Accept	application/json, text/plain, */*
X-stream	true
Content-Type	application/json; charset=utf-8
Authorization	Basic bmVvNGo6YWRTaW4=

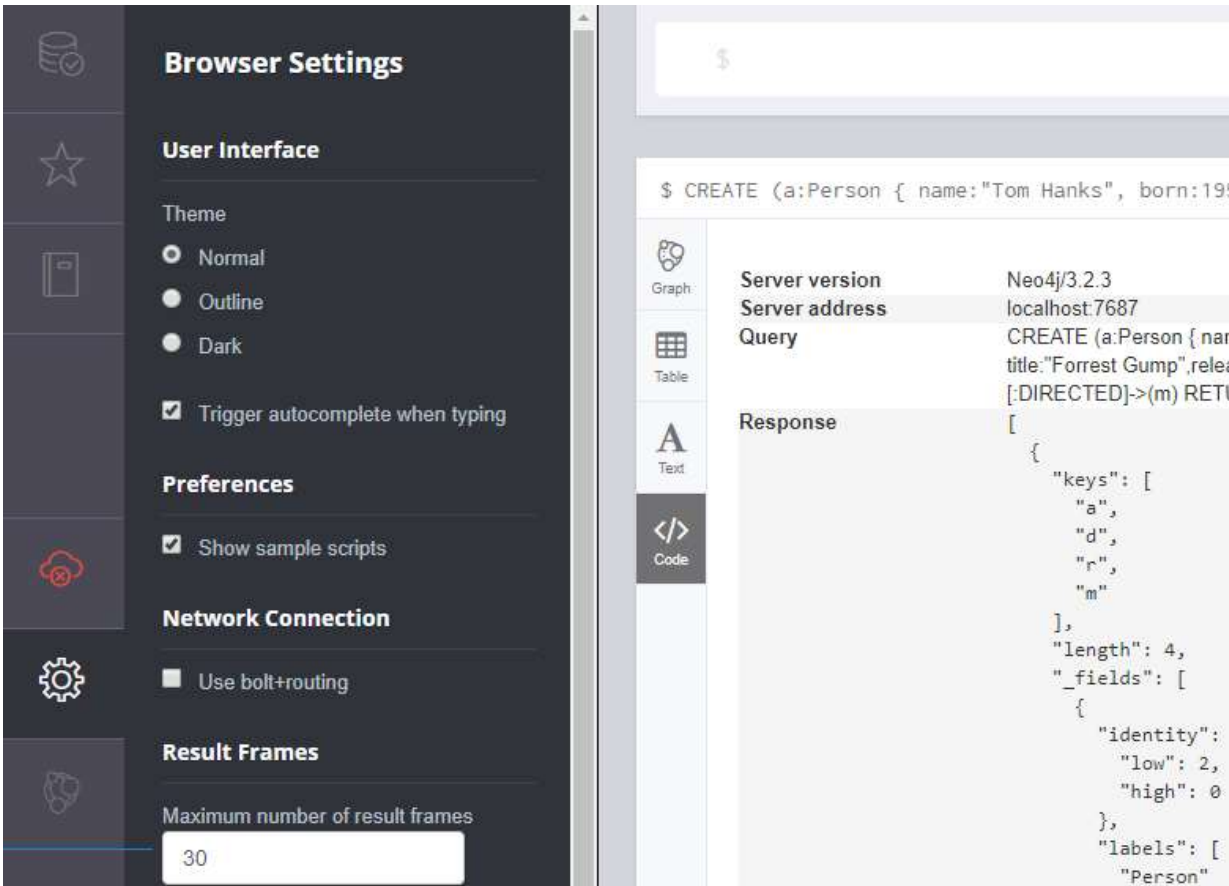
Payload

```
{  "statements": [    {      "statement": "CREATE (a:Person { name:\"Tom\", \"resultDataContents\": [        \"row\",        \"graph\"      ]},      \"includeStats\": true    )"}  ]}
```

▼ Response

Setting, Bolt vs. HTTP

- On previous diagram we see details of HTTP Request/Response dialogue.
- To see HTTP dialogue you select Code `</>` icon. However, you must also go to the bottom of the left navigation bar, select  (Settings) and make sure that box next to "Use bolt+routing" is not checked.
- Bolt traffic is binary and does not have HTTP Request format.



The screenshot shows the Neo4j Browser interface. On the left is a dark sidebar with navigation icons and settings sections. The main area on the right displays a query and its result.

Browser Settings


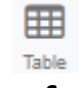
- User Interface**
 - Theme: ☐ Normal, ☐ Outline, ☐ Dark
 - ☒ Trigger autocomplete when typing
- Preferences**
 - ☒ Show sample scripts
- Network Connection**
 - ☐ Use bolt+routing
- Result Frames**
 - Maximum number of result frames: 30

Query and Response

Query: `$ CREATE (a:Person { name:"Tom Hanks", born:195`

Response:

```
[
  {
    "keys": [
      "a",
      "d",
      "r",
      "m"
    ],
    "length": 4,
    "_fields": [
      {
        "identity":
          "low": 2,
          "high": 0
      },
      "labels": [
        "Person"
      ]
    }
  }
]
```

- Select  to return a graph
- Select  to return a table of values

Fetching Values, MATCH-ing Patterns

- Matching patterns is a task for the MATCH statement. We pass the patterns we've used so far to the MATCH statement to describe what we're looking for.
- A MATCH statement will search for the patterns we specify and return one row per successful pattern match.
- We can ask for all nodes labeled with the Movie label or a Person named Keanu Reeves

```
MATCH (m:Movie) RETURN m;
```

```
MATCH (p:Person { name:"Keanu Reeves" }) RETURN p;
```

- We can also find more interesting connections, like for instance the movies titles that *Tom Hanks* acted in and the roles he played.

```
MATCH (p:Person { name:"Tom Hanks" })-[r:ACTED_IN]->(m:Movie)
RETURN m.title, r.roles
```

m.title	r.roles
Forrest Gump	[Forrest]

- Cypher returned the properties of the nodes and relationships that we were interested in. You can access properties of nodes or relationships using the dot notation: `identifier.property`

Extending the Graph, Attaching Structures

- To extend the graph with new information, we first match the existing connection points and then attach the newly created nodes to them with relationships.
- Adding *Cloud Atlas* as a new movie for *Tom Hanks* could be achieved like this:

```
MATCH (p:Person { name:"Tom Hanks" })  
CREATE (m:Movie { title:"Cloud Atlas",released:2012 })  
CREATE (p)-[r:ACTED_IN { roles: ['Zachry']}]>(m) RETURN p,r,m
```



- We can assign identifiers to both nodes and relationships and use them later on, no matter if they were created or matched.
- A tricky aspect of the combination of MATCH and CREATE is that we get *one row per matched pattern*. This causes subsequent CREATE statements to be executed once for each row. In many cases this is what you want. If that's not intended, move the CREATE statement before the MATCH, or change the cardinality of the query.

Completing Patterns, MERGE (update)

- Whenever we get data from external systems or are not sure if certain information already exists in the graph, we want to be able to express a repeatable (idempotent) update operation.
- In Cypher MERGE command has this idempotent property. Acts like a combination of MATCH or CREATE, which checks for the existence of data before creating it.
- With MERGE you define a pattern to be found or created. Usually, as with MATCH you only want to include the key property to look for in your core pattern. MERGE allows you to provide additional properties you want to set ON CREATE.
- For example, even if graph already had *Cloud Atlas* we could merge it in again:

```
MERGE (m:Movie { title:"Cloud Atlas" })  
ON CREATE SET m.released = 2012 RETURN m
```

- Clause SET m.released = ... acted as UPDATE statement of standard SQL
- MERGE makes sure that you can't create duplicate information or structures, but it comes with the cost of needing to check for existing matches. Especially on large graphs it can be costly to scan a large set of labeled nodes for a certain property. You can alleviate some of that by creating supporting indexes or constraints.

```
MERGE (m:Movie { title:"Cloud Atlas" })  
ON MATCH SET m.released = 2013 RETURN m
```

- If Movie exists, above statement will change released to 2013.

MERGE

- MERGE can also assert that a relationship is only created once. For that to work *you have to pass in* both nodes from a previous pattern match.

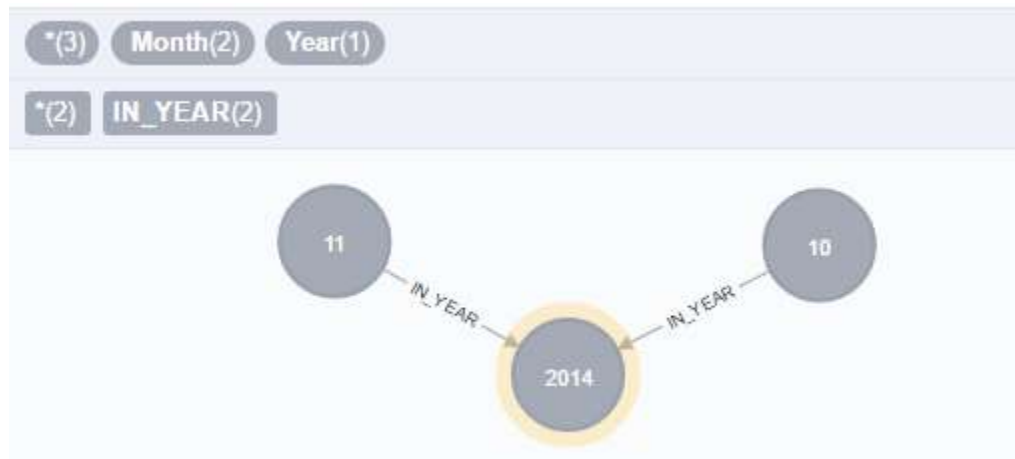
```
MATCH (m:Movie { title:"Cloud Atlas" })
MATCH (p:Person { name:"Tom Hanks" })
MERGE (p)-[r:ACTED_IN]->(m)
ON CREATE SET r.roles =['Zachry']
ON MATCH set r.rewarded = ['No'] RETURN p,r,m
```

- If relationship `ACTED_IN` does not exist it will be created.
- If relationship exists it will get a new property: `rewarded` with value of `'No'` .
- In case the direction of a relationship is arbitrary, you can leave off the arrowhead. MERGE will then check for the relationship in either direction, and create a new directed relationship if no matching relationship was found.

Creating Branching Points

- If you choose to pass in only one node from a preceding clause, MERGE offers an interesting functionality. It will only match within the direct neighborhood of the provided node for the given pattern, and, if not found create it. This can be used for creating tree structures.

```
CREATE (y:Year { year:2014 })
MERGE (y)<-[:IN_YEAR]-(m10:Month { month:10, name="October" })
MERGE (y)<-[:IN_YEAR]-(m11:Month { month:11, name="November" })
RETURN y,m10,m11
```



Deleting Nodes and Relationships

- To delete unattached nodes with label Label, use DELETE clause:

```
MATCH (n:Label) DELETE n;
```

- To delete all nodes and relationships you can do this:

```
MATCH (n) DETACH DELETE n;
```

Deleted 6 nodes, deleted 6 relationships, statement executed in 49 ms.

Selecting Results we Need

- In normal, relational, SQL we perform many operations in order to restrict and select results we need. Cypher performs almost all of those operations:
 - Filtering
 - Returning
 - Aggregating
 - Ordering and Paginating
 - Collecting Aggregation

- Let us create (extend) a more elaborate graph, first:

```
CREATE (matrix:Movie { title:"The Matrix",released:1997 })
CREATE (cloudAtlas:Movie { title:"Cloud Atlas",released:2012 })
CREATE (forrestGump:Movie { title:"Forrest Gump",released:1994 })
CREATE (keanu:Person { name:"Keanu Reeves", born:1964 })
CREATE (robert:Person { name:"Robert Zemeckis", born:1951 })
CREATE (tom:Person { name:"Tom Hanks", born:1956 })
CREATE (tom)-[:ACTED_IN { roles: ["Forrest"]}]>(forrestGump)
CREATE (tom)-[:ACTED_IN { roles: ['Zachry']}]>(cloudAtlas)
CREATE (robert)-[:DIRECTED]>(forrestGump)
```

- Note, there are no semicolons between statements. This is a single statement.

Filtering, WHERE Clause

- Quite often there are conditions on what we want to see. Similarly to *SQL*, Cypher filter conditions are expressed in a *WHERE* clause.
- *WHERE* clause allows use of any number of Boolean expressions (predicates) combined with *AND*, *OR*, *XOR* and *NOT*. The simplest predicates are comparisons.

```
MATCH (m:Movie)
```

```
WHERE m.title = "The Matrix" RETURN m
```

- Or in a slightly more compact notation

```
MATCH (m:Movie { title: "The Matrix" }) RETURN m
```

- *WHERE* clause could include numeric comparisons, matching regular expressions and checking the existence of values within a collection.

```
MATCH (p:Person) WHERE p.name =~ "K.+" RETURN p; //or
```

```
MATCH (p:Person)-[r:ACTED_IN]->(m:Movie)
```

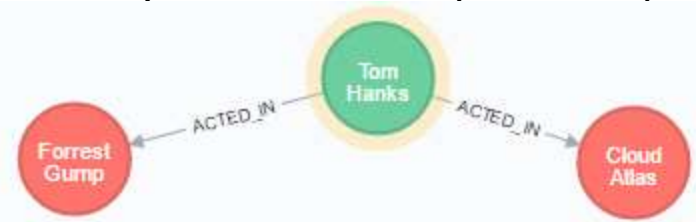
```
WHERE p.name =~ "K.+" OR m.released > 2000 OR "Neo" IN r.roles RETURN p,r,m
```

- We can use patterns as predicates. A pattern predicate restricts the current result set. It only allows the paths to pass that satisfy the additional patterns (or NOT)

```
MATCH (p:Person)-[:ACTED_IN]->(m)
```

```
WHERE NOT (p)-[:DIRECTED]->()
```

```
RETURN p,m
```




RETURN-ing Results

- RETURN clause can actually return any number of expressions and not only nodes, relationships, or paths .
- What are actually expressions in Cypher? The simplest expressions are literal values like numbers, strings and arrays as `[1,2,3]`, and maps like `{name:"Tom Hanks", born:1964, movies:["Forrest Gump", ...], count:13}`.
- We can access individual properties of any node, relationship, or map with a dot-syntax like `n.name`. Individual elements or slices of arrays can be retrieved with subscripts like: `names[0]` or `movies[1..-1]`
- Each function evaluation like: `length(array)`, `toInteger("12")`, `substring("2014-07-01",0,4)`, or `coalesce(p.nickname, "n/a")` is also an expression.
- By default, the expression itself will be used as label for the column, in many cases you want to alias with a more understandable name using `expression AS alias`. You can later on refer to that expression (column) using its alias.

```
MATCH (p:Person)
```

```
RETURN DISTINCT p, p.name AS name, upper(p.name),  
coalesce(p.nickname, "n/a") AS nickname, { name: p.name,  
label:head(labels(p))} AS person
```

- Here we used key word `DISTINCT` to eliminate duplicate rows.
- Function `coalesce()` returns the first non-null value in the list of its arguments
- To see all properties, use  (text output)

Grouping, Aggregating

- In many cases you want to aggregate or group the data that you encounter while traversing patterns in your graph. Aggregation happens in the RETURN clause while computing your final results.
- Many common aggregation functions are supported, e.g. `count`, `sum`, `avg`, `min`, and `max`, and several more.
- Counting the number of people in your database could be achieved by this query:

```
MATCH (:Person) RETURN count(*) AS people
```

3

```
MATCH (p:Person) RETURN count(DISTINCT p.name) AS people
```

3

- NULL values are skipped during aggregation. For aggregating only unique values use `DISTINCT`, like in `count(DISTINCT role)`.
- *Cypher uses all non-aggregated columns as grouping keys.* Aggregation affects which data is still visible in ordering or later query parts.
- For example, to find out how often an actor and director worked together, you'd run this statement:

```
MATCH (actor:Person)-[:ACTED_IN]->(movie:Movie)<-[:DIRECTED]-(director:Person)
RETURN actor,director,count(*) AS collaborations
```

"actor"	"director"	"collaborations"
{ "born":1956,"name":"Tom Hanks" }	{ "born":1951,"name":"Robert Zemeckis" }	1

Ordering and Pagination

- Ordering is imposed with `ORDER BY expression [ASC|DESC]` clause.
- If you return `person.name` you can still `ORDER BY person.age` as both are accessible from the person reference. You cannot order by things that you can't infer from the information you return. This is especially important with aggregation and `DISTINCT` return values as both remove the visibility of data that is aggregated.
- Pagination is achieved by use of `SKIP {offset} LIMIT {count}`.
- A common pattern is to aggregate for a `count` (score or frequency), order by that result (`count`) and only return the top-n entries.
- For example, the most prolific actors could be found as:

```
MATCH (a:Person)-[:ACTED_IN]->(m:Movie)
RETURN a,count(*) AS appearances
ORDER BY appearances DESC SKIP 2 LIMIT 10;
```

- We have one row in the database so we use `SKIP 0` (skip no rows). Had we had many rows, we could have skipped a finite number of them and then display the first 10.
- `appearance DESC(ending)` tells the query engine to display rows in descending order of number `appearance`.

Collecting Aggregation

- The most helpful aggregation function is `collect()`, which, as the name says, collects all aggregated values into a *real* array or list. With COLLECT we don't lose the detail information while aggregating.
- `collect()` is well suited for retrieving the typical parent-child structures, where one core entity (parent, root or head) is returned per row with all its dependent information in associated collections created with collect. This means there's no need to repeat the parent information per each child-row or even running 1+n statements to retrieve the parent and its children individually.
- To retrieve the cast of each movie in our database you could use this statement:

```
MATCH (m:Movie)<-[:ACTED_IN]-(a:Person)
RETURN m.title AS movie, collect(DISTINCT a.name) AS cast, count(*) AS
actors
```

movie	cast	actors
Forrest Gump	[Tom Hanks]	1
Cloud Atlas	[Tom Hanks]	1

Composing Large Statements

- If you want to combine the results of two statements that have the same result structure, you can use UNION [ALL].
- For instance if you want to list both actors and directors without using the alternative relationship-type syntax `()-[:ACTED_IN| :DIRECTED]->()` we can do:

```
MATCH (actor:Person)-[r:ACTED_IN]->(movie:Movie)
```

```
RETURN actor.name AS name, type(r) AS acted_in, movie.title AS title
```

```
UNION
```

```
MATCH (director:Person)-[r:DIRECTED]->(movie:Movie)
```

```
RETURN director.name AS name, type(r) AS acted_in, movie.title AS title
```

name	acted_in	title
Tom Hanks	ACTED_IN	Cloud Atlas
Tom Hanks	ACTED_IN	Forrest Gump
Robert Zemeckis	DIRECTED	Forrest Gump

Chain Statements with `WITH` clause

- It is possible to chain fragments of statements together, much like you would do within a data flow pipeline. Each fragment works on the output from the previous one and its results can feed into the next one.
- You use the `WITH` clause to combine the individual parts and declare which data flows from one to the other. `WITH` is very much like `RETURN` with the difference that it doesn't finish a query but prepares the input for the next part.
- You can use the same expressions, aggregations, ordering and pagination as in the `RETURN` clause. However, you *must* alias all columns as they would otherwise not be accessible. Only columns that you declare in your `WITH` clause are available in subsequent query parts.
- Below, we collect the movies someone appeared in, and then filter out (remove) those actors that appear in only one movie.

```
MATCH (person:Person)-[:ACTED_IN]->(m:Movie)
WITH person, count(*) AS appearances, collect(m.title) AS movies
WHERE appearances > 1
RETURN person.name, appearances, movies
```

person.name	appearances	movies
Tom Hanks	2	[Cloud Atlas, Forrest Gump]

- In SQL, if you want to filter by an aggregated value, you would have to use `HAVING`. That's a single purpose clause for filtering aggregated information. In Cypher, `WHERE` can be used in both case

Labels, Constraints

- **Labels** are a convenient way to group nodes together. They are used to restrict queries, define constraints and create indexes.
- Unique constraints are used to guarantee uniqueness of a certain property on nodes with a specific label.
- These constraints are also used by the MERGE clause to make certain that a node only exists once.
- Let's add a constraint. For example, we decided that all Movie node `titles` should be unique.

```
CREATE CONSTRAINT ON (movie:Movie) ASSERT movie.title IS UNIQUE
```

Added 1 constraint, statement executed in 314 ms.

- Note that adding the unique constraint will add an index on that property, so we won't do that separately.
- If we drop a constraint, and still want an index on the same property, we have to create such an index.
- Constraints can be added after a label is already in use, but that requires that the existing data complies with the constraints.

Indexes

- For a graph query to scan a graph quickly we don't need indexes. However, we need indexes to find the starting points of our graphs.
- So, the reason for using indexes in a graph database is to find the starting points in the graph as fast as possible. After the initial index seek you rely on in-graph structures and the first class citizenship of relationships in the graph database to achieve high performance.
- In this case we want an index to speed up finding actors by name in the database:

```
CREATE INDEX ON :Actor(name)
```

Added 1 index, statement executed in 44 ms.

- If we add new data, indexes are applied automatically.

```
CREATE (actor:Actor { name:"Tom Hanks" }), (movie:Movie { title:'Sleepless  
IN Seattle' }), (actor)-[:ACTED_IN]->(movie);
```

- Next time we ask for Tom Hanks node, the index will kick in behind the scenes to boost performance.

```
MATCH (actor:Actor { name: "Tom Hanks" })  
RETURN actor;
```

Labels

- We can add more than one label to a node or relationship. For example, so far Tom Hanks was only an Actor. We could add another label, American:

```
MATCH (actor:Actor { name: "Tom Hanks" })  
SET actor:American;
```

- We could remove Labels from nodes and relationships:

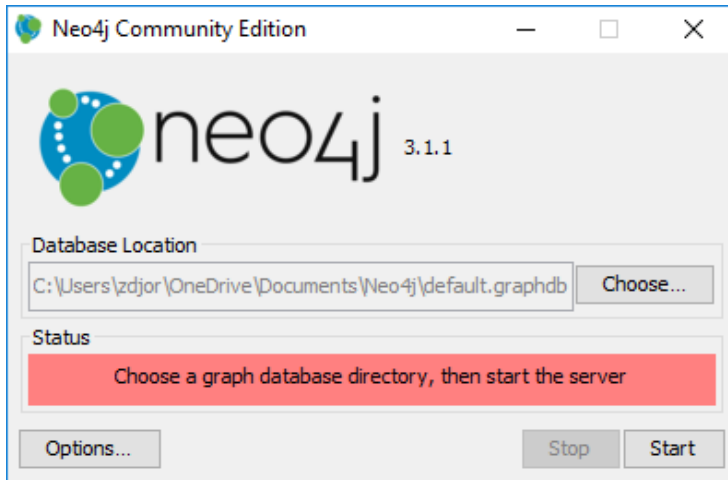
```
MATCH (actor:Actor { name: "Tom Hanks" })  
REMOVE actor:American;
```

Loading Data, LOAD CSV

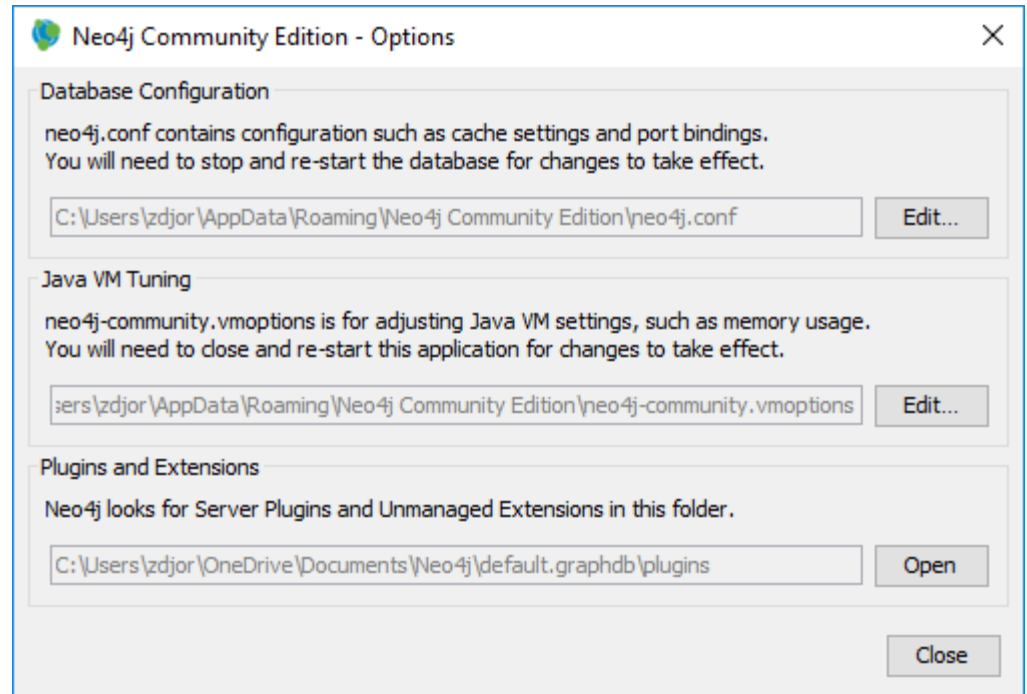
- Create statements we have seen so far are not practical if we have a very large number of objects. Often we need to use an existing data source to drive graph generation or modification process.
- Cypher provides an elegant built-in way to import tabular CSV data into graph structures.
- The LOAD CSV clause parses a local or remote file into a stream of rows which represent maps (with headers) or lists. Then you can use whatever Cypher operations you want to apply to either create nodes or relationships or to merge with existing graph structures.
- As CSV files usually represent either node- or relationship-lists, you run multiple passes to create nodes and relationships separately.
- We will use LOAD CSV command to import data from CSV files
- The URL of the CSV file is specified by using FROM followed by an arbitrary expression evaluating to the URL in question.
- It is required to specify an identifier for the CSV data using AS.
- LOAD CSV supports resources compressed with gzip, Deflate, as well as ZIP archives.
- CSV files can be stored on the OS of Neo4J server and accessed using a file:/// URL.
- Alternatively, LOAD CSV also supports accessing CSV files via HTTPS, HTTP, and FTP.

Configure LOAD CSV

- Stop your server



- Select Options and then Edit next to Database Configuration
- Add values listed on next slide
- Save edited file.
- Start the server.



neo4j.conf File

- Server properties file needs new lines in red. Change

`dbms.security.auth_enabled` to `false`

```
*****
# Server configuration
*****
# This setting constrains all `LOAD CSV` import files to be
under the `import` directory. Remove or uncomment it to
# allow files to be loaded from anywhere in filesystem; this
introduces possible security problems. See the `LOAD CSV`
# section of the manual for details.
# Allow CSV file loading
allow_file_urls=true
dbms.directories.import=C:\\Zoran\\code\\csv_data
dbms.security.allow_csv_import_from_file_urls=true

# Require(or disable the requirement of) auth to access Neo4j
dbms.security.auth_enabled=false
```


Note on position of csv-files directory

- My Neo4J database resides in C:\Program Files\Neo4j CE 3.1.1
- My neo4j.conf configuration files is in :
c:\Users\073621\AppData\Roaming\Neo4j Community Edition\
- The directory, in which I placed 4 csv files, is presented in neo4j.conf file as
C:\\Zoran\\code\\csv_data # double slashes are because of Java

persons.csv

```
id,name
1,Charlie Sheen
2,Oliver Stone
3,Michael Douglas
4,Martin Sheen
5,Morgan Freeman
```

movies.csv

```
id,title,country,year
1,Wall Street,USA,1987
2,The American President,USA,1995
3,The Shawshank Redemption,USA,1994
```

roles.csv

```
personId,movieId,role
1,1,Bud Fox
4,1,Carl Fox
3,1,Gordon Gekko
4,2,A.J. MacInerney
3,2,President Andrew Shepherd
5,3,Ellis Boyd 'Red' Redding
```

movie_actor_role.csv

```
title;released;actor;born;characters
Back to the Future;1985;Michael J. Fox;1961;Marty McFly
Back to the Future;1985;Christopher Lloyd;1938;Dr. Emmet Brown
```

Loading data from CSV files

- To load data in those 4 files I run the following 4 commands:

```
USING PERIODIC COMMIT LOAD CSV WITH HEADERS FROM "file:///persons.csv" AS line
MERGE (a:Person { id:line.id })
ON CREATE SET a.name=line.name;
```

Added 5 labels, created 5 nodes, set 10 properties, statement executed in 568 ms.

```
LOAD CSV WITH HEADERS FROM "file:///movies.csv" AS line
CREATE (m:Movie { id:line.id,title:line.title, released:toInteger(line.year)});
Added 3 labels, created 3 nodes, set 9 properties, statement executed in 467 ms.
```

```
LOAD CSV WITH HEADERS FROM "file:///roles.csv" AS line
MATCH (m:Movie { id:line.movieId })
MATCH (a:Person { id:line.personId })
CREATE (a)-[:ACTED_IN { roles: [line.role]}]->(m);
Set 6 properties, created 6 relationships, completed after 123 ms.
```

```
LOAD CSV WITH HEADERS FROM "file:///movie_actor_roles.csv" AS line
FIELDTERMINATOR ";";
MERGE (m:Movie { title:line.title })
ON CREATE SET m.released = toInteger(line.released)
MERGE (a:Person { name:line.actor })
ON CREATE SET a.born = toInteger(line.born)
MERGE (a)-[:ACTED_IN { roles:split(line.characters,",")}]->(m);
Added 3 labels, created 3 nodes, set 8 properties, created 2 relationships, statement
executed in 235 ms
```

Notes

- If you import a large amount of data (more than 10000 rows), it is recommended to prefix your LOAD CSV clause with a `USING PERIODIC COMMIT` hint. This allows Neo4j to regularly commit the import transactions and avoid memory issues.
- You can load files through HTTP or HTTPS protocol as well. One of previous commands could read like:

LOAD CSV WITH HEADERS FROM

["http://neo4j.com/docs/2.3.3/csv/intro/movies.csv"](http://neo4j.com/docs/2.3.3/csv/intro/movies.csv) AS line

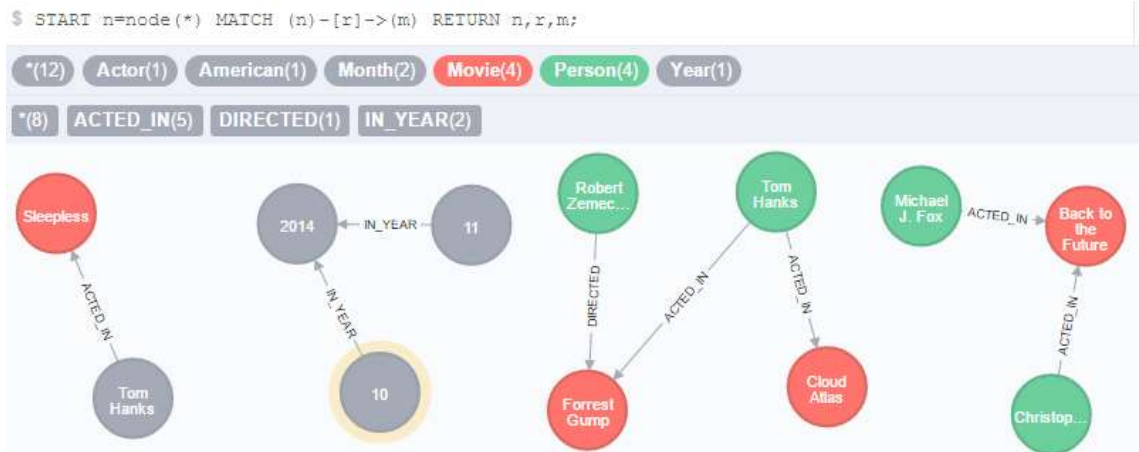
```
CREATE (m:Movie { id:line.id,title:line.title,  
released:toInteger(line.year) });
```

- After all this work, we would like to see what we got.
- To display all nodes associated with some relationships we could ask:

```
START n=node(*) MATCH (n)-[r]->(m) RETURN n,r,m;
```

- If we want all the nodes whether they are in relationship or not, we could ask:

```
START n=node(*) RETURN n;
```



Data Structures

- Cypher can create and consume more complex data structures out of the box.
- You can create literal lists (`[1,2,3]`) and maps (`{name: value}`) within a statement.
- There are a number of functions that work with lists. For example, function `size(list)` returns the size of a list. Function `reduce()` runs an expression, a function, against all elements of the list and accumulates the results.
- Let us collect the names of actors in every movie, and return up to 3 of them:

```
MATCH (movie:Movie)<-[:ACTED_IN]-(actor:Person)
RETURN movie.title AS movie, collect(actor.name)[0..2] AS the_cast;
```

movie	the_cast
Forrest Gump	[Tom Hanks]
Cloud Atlas	[Tom Hanks]
Back to the Future	[Christopher Lloyd, Michael J. Fox]

- In the last line we are accessing elements 0 to 2 of list `collect(actor.name)[0..2]`.
- You can also access individual elements or slices of a list with `list[1]` or `list[5..-5]`. Other functions one could use to access parts of a list are `head(list)`, `tail(list)` and `last(list)`

List Processing

- Often you want to process lists to filter, aggregate (reduce) or transform (extract) their values. Those transformations can be done within Cypher or in the calling code. This kind of list-processing can reduce the amount of data handled and returned, so it might make sense to do it within the Cypher statement. For example:

```
WITH range(1,10) AS numbers
WITH extract(n IN numbers | n*n) AS squares
WITH filter(n IN squares WHERE n > 16) AS large_squares
RETURN reduce(a = 0, n IN large_squares | a + n) AS sum_large_squares;
```

sum_large_squares

355

- In a graph-query we can filter or aggregate collected values or work on array properties.

```
MATCH (m:Movie)<-[r:ACTED_IN]-(a:Person)
WITH m.title AS movie, collect({ name: a.name, roles: r.roles }) AS cast
RETURN movie, filter(actor IN cast WHERE actor.name STARTS WITH "M") Starts_M
```

Unwind Lists

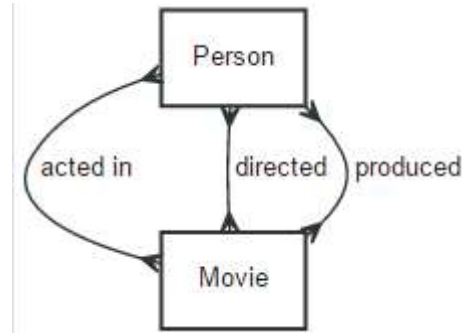
- Sometimes you have collected information into a list, but want to use each element individually as a row. For instance, you might want to further match patterns in the graph. Or you passed in a collection of values but now want to create or match a node or relationship for each element.
- You can use the UNWIND clause to unroll a list into a sequence of row values.
- For instance, a query to find the top 3 `co_actors` and then follow their movies and again list the `cast` for each of those movies:

```
MATCH (actor:Person)-[:ACTED_IN]->(movie:Movie)<-[:ACTED_IN]-(
co_actor:Person)
WHERE actor.name < co_actor.name
WITH actor, co_actor, count(*) AS frequency, collect(movie) AS movies
ORDER BY frequency DESC LIMIT 3 UNWIND movies AS m
MATCH (m)<-[:ACTED_IN]-(a)
RETURN m.title AS movie, collect(a.name) AS cast
```

movie	cast
Back to the Future	[Christopher Lloyd, Michael J. Fox]

Comparison of RDBMS and Graph Database

- Entity Relationship with use in relational models for a person who participates in movies looks like the one the right.
- We have Person and Movie entities, which are related in three different ways, each of which have many-to-many cardinality.
- In a RDBMS we would use tables for the entities as well as for the associative entities (join tables) needed. In this case we decided to go with the following tables: `movie`, `person`, `acted_in`, `directed`, `produced`. You'll find the SQL for this arrangement on the next slide.
- In Neo4j, the basic data units are nodes and relationships. Both can have properties, which correspond to attributes in a RDBMS.
- Nodes can be grouped by putting labels on them. In the example, we will use the labels `Movie` and `Person`.
- When using Neo4j, related entities can be represented directly by using relationships. There's no need to deal with foreign keys to handle the relationships, the database will take care of such mechanics.
- Also, the relationships always have full referential integrity. There are no constraints to enable for this, as it's not optional; it's really part of the underlying data model.
- Relationships always have a type, and we will differentiate the different kinds of relationships by using the types `ACTED_IN`, `DIRECTED`, `PRODUCED`.



Relational Tables

```
CREATE TABLE movie (
  id INTEGER,
  title VARCHAR(100),
  released INTEGER,
  tagline VARCHAR(100)
);

CREATE TABLE person (
  id INTEGER,
  name VARCHAR(100),
  born INTEGER
);

CREATE TABLE acted_in (
  role varchar(100),
  person_id INTEGER,
  movie_id INTEGER
);

CREATE TABLE directed (
  person_id INTEGER,
  movie_id INTEGER
);

CREATE TABLE produced (
  person_id INTEGER,
  movie_id INTEGER
);

INSERT INTO movie (id, title, released, tagline)
VALUES (
  (1, 'The Matrix', 1999, 'Welcome to the Real World'),
  (2, 'The Devil''s Advocate', 1997, 'Evil has its winning ways'),
  (3, 'Monster', 2003, 'The first female serial killer of America')
);

INSERT INTO person (id, name, born)
VALUES (
  (1, 'Keanu Reeves', 1964), (2, 'Carrie-Anne Moss', 1967),
  (3, 'Laurence Fishburne', 1961), (4, 'Hugo Weaving', 1960),
  (5, 'Andy Wachowski', 1967), (6, 'Lana Wachowski', 1965),
  (7, 'Joel Silver', 1952), (8, 'Charlize Theron', 1975),
  (9, 'Al Pacino', 1940), (10, 'Taylor Hackford', 1944) );

INSERT INTO acted_in (role, person_id, movie_id)
VALUES (
  ('Neo', 1, 1), ('Trinity', 2, 1), ('Morpheus', 3, 1),
  ('Agent Smith', 4, 1), ('Kevin Lomax', 1, 2),
  ('Mary Ann Lomax', 8, 2), ('John Milton', 9, 2),
  ('Aileen', 8, 3) );

INSERT INTO directed (person_id, movie_id)
VALUES (
  (5, 1), (6, 1), (10, 2));

INSERT INTO produced (person_id, movie_id)
VALUES (
  (7, 1), (8, 3) );
```


Neo4J

- In Neo4j we won't create any schema up front. Labels can be used right away without declaring them. In other words, there is no predefined schema.
- In the CREATE statements below, we tell Neo4j what data we want to have in the graph. Simply put, the parentheses denote nodes, while the arrows (`-->`, or in our case with a relationship type included `-[:DIRECTED]->`) denote relationships.
- For the nodes we set identifiers like `TheMatrix` so we can easily refer to them later on in the statement. The identifiers are scoped to the statement, and not visible to other Cypher statements.

Create Statement in Neo4J

```
CREATE (TheMatrix:Movie { title:'The Matrix', released:1999, tagline:'Welcome to the
Real World' })
CREATE (Keanu:Person { name:'Keanu Reeves', born:1964 })
CREATE (Carrie:Person { name:'Carrie-Anne Moss', born:1967 })
CREATE (Laurence:Person { name:'Laurence Fishburne', born:1961 })
CREATE (Hugo:Person { name:'Hugo Weaving', born:1960 })
CREATE (AndyW:Person { name:'Andy Wachowski', born:1967 })
CREATE (LanaW:Person { name:'Lana Wachowski', born:1965 })
CREATE (JoelS:Person { name:'Joel Silver', born:1952 })
CREATE (Keanu)-[:ACTED_IN { roles: ['Neo']}]>(TheMatrix),
      (Carrie)-[:ACTED_IN { roles: ['Trinity']}]>(TheMatrix),
      (Laurence)-[:ACTED_IN { roles: ['Morpheus']}]>(TheMatrix),
      (Hugo)-[:ACTED_IN { roles: ['Agent Smith']}]>(TheMatrix), (AndyW)-[:DIRECTED]>
>(TheMatrix),
      (LanaW)-[:DIRECTED]>(TheMatrix), (JoelS)-[:PRODUCED]>(TheMatrix)
CREATE (TheDevilsAdvocate:Movie { title:"The Devil's Advocate", released:1997,
      tagline: 'Evil has its winning ways' })
CREATE (Monster:Movie { title: 'Monster', released: 2003,
      tagline: 'The first female serial killer of America' })
CREATE (Charlize:Person { name:'Charlize Theron', born:1975 })
CREATE (Al:Person { name:'Al Pacino', born:1940 })
CREATE (Taylor:Person { name:'Taylor Hackford', born:1944 })
CREATE (Keanu)-[:ACTED_IN { roles: ['Kevin Lomax']}]>(TheDevilsAdvocate),
      (Charlize)-[:ACTED_IN { roles: ['Mary Ann Lomax']}]>(TheDevilsAdvocate),
      (Al)-[:ACTED_IN { roles: ['John Milton']}]>(TheDevilsAdvocate),
      (Taylor)-[:DIRECTED]>(TheDevilsAdvocate), (Charlize)-[:ACTED_IN { roles:
['Aileen']}]>(Monster),
      (Charlize)-[:PRODUCED { roles: ['Aileen']}]>(Monster)
```

Read the Data

- In RDBMS:

```
SELECT movie.title
FROM movie;
```

```
SELECT movie.title
FROM movie
WHERE movie.released > 1998;
```

```
// names of actors and their movies
SELECT person.name, movie.title
FROM person
JOIN acted_in AS acted_in ON
acted_in.person_id = person.id
JOIN movie ON acted_in.movie_id =
movie.id;
```

```
// actors in a movie with K. Reeves
SELECT DISTINCT co_actor.name
FROM person AS keanu
    JOIN acted_in AS acted_in1 ON
acted_in1.person_id = keanu.id
    JOIN acted_in AS acted_in2 ON
acted_in2.movie_id = acted_in1.movie_id
    JOIN person AS co_actor
    ON acted_in2.person_id =
co_actor.id AND co_actor.id <> keanu.id
WHERE keanu.name = 'Keanu Reeves';
```

- In Neo4J

```
MATCH (movie:Movie)
RETURN movie.title;
```

```
MATCH (movie:Movie)
WHERE movie.released > 1998
RETURN movie.title;
```

```
// names of actors and their movies
MATCH (person:Person)-[:ACTED_IN]->
movie:Movie)
RETURN person.name, movie.title;
```

```
// actors in a movie with K. Reeves
MATCH (keanu:Person)-[:ACTED_IN]->
(movie:Movie), (coActor:Person)-[:ACTED_IN]->
(movie)
WHERE keanu.name = 'Keanu Reeves'
```

Read the Data

- In RDBMS:

```
//who both acted and produced
movies.
SELECT person.name
FROM person
WHERE person.id IN (SELECT person_id
FROM acted_in)
    AND person.id IN (SELECT person_id
FROM produced);

//directors of K. Reeves movies
SELECT director.name, count(*)
FROM person keanu
    JOIN acted_in ON keanu.id =
acted_in.person_id
    JOIN directed ON acted_in.movie_id
= directed.movie_id
    JOIN person AS director ON
directed.person_id = director.id
WHERE keanu.name = 'Keanu Reeves'
GROUP BY director.name
ORDER BY count(*) DESC
```

- In Neo4J

```
//who both acted and produced movies.
MATCH (person:Person)
WHERE (person)-[:ACTED_IN]->() AND (person)-
[:PRODUCED]->()
RETURN person.name;

//directors of K. Reeves movies
MATCH (keanu:Person { name: 'Keanu Reeves'
})-[:ACTED_IN]->(movie:Movie),
    (director:Person)-[:DIRECTED]->(movie)
RETURN director.name, count(*)
ORDER BY count(*) DESC
```

Programmatic Access to Neo4J

- Starting with release 3.0 Neo4j supports a binary protocol called Bolt.
- Bolt is based on the PackStream serialization and supports protocol versioning, authentication and TLS via certificates.
- The binary protocol is enabled in Neo4j by default, so you can use any language driver that supports it.
- Neo4j officially provides drivers for:
 - .Net,
 - Java,
 - JavaScript and
 - Python
- It appears that there are also community drivers for
 - Ruby,
 - PHP and
 - R

Community Drivers

- There are several Python packages developed outside of Neo4j that provide full access to Neo4J database. Some of them are:
- Py2neo <http://py2neo.org/2.0/essentials.html>
- Neomodel (Object Graph Mapper)
<http://neomodel.readthedocs.io/en/latest/#>
- Neo4jRestClient <http://neo4j-rest-client.readthedocs.io/en/latest/index.html>
- Bulbflow <http://bulbflow.com/>
- Package RNeo4J (<https://neo4j.com/developer/r/>) allows you to access Neo4J database from R. The package is written by Nicole White who works at Neo4J.
- You can access Neo4J from Scala code. See these links:
<https://neo4j.com/developer/java/# scala anormcypher>

<https://www.packtpub.com/mapt/book/big-data-and-business-intelligence/9781783287253/2/ch02lv11sec34/accessing-neo4j-from-scala>

Neo4J Python Driver

- Neo4j can be accessed via its binary (Bolt) and HTTP APIs.
- You can use the official binary driver for Python (neo4j-python-driver) or connect via HTTP with any of our community drivers.
- The Neo4j Python driver is **officially supported** by Neo4j and connects to the database using the binary protocol. It aims to be minimal, while being idiomatic to Python.

```
C:\Zoran\code\neo4j> pip install neo4j-driver
```

```
Collecting neo4j-driver
```

```
  Downloading neo4j_driver-1.1.0-py2.py3-none-any.whl (43kB)
```

```
    100% |#####| 51kB 304kB/s
```

```
Installing collected packages: neo4j-driver
```

```
Successfully installed neo4j-driver-1.1.0
```

- Another popular Python driver is Py2Neo. Installs similarly:

```
C:\Zoran\code\neo4j> pip install py2neo
```

```
Collecting py2neo
```

```
  Downloading py2neo-3.1.2.tar.gz (100kB)
```

```
    100% |#####| 102kB 664kB/s
```

```
Building wheels for collected packages: py2neo
```

```
  Running setup.py bdist_wheel for py2neo ... done
```

```
Successfully built py2neo
```

```
Installing collected packages: py2neo
```

```
Successfully installed py2neo-3.1.2
```

Neo4j Python Client

- Python Clients could be used to pass Cypher queries and retrieve result.
- For example

```
from neo4j.v1 import GraphDatabase, basic_auth
```

```
driver =  
GraphDatabase.driver("bolt://localhost:7687",auth=basic_auth("neo4j",  
"neo4j"))
```

```
session = driver.session()
```

```
session.run("CREATE (a:Person {name: {name}, title: {title}})",  
{ "name": "Arthur", "title": "King" })
```

```
result = session.run("MATCH (a:Person) WHERE a.name = {name} "  
                    "RETURN a.name AS name, a.title AS title",  
                    { "name": "Arthur" })
```

```
for record in result:
```

```
    print "%s %s" % (record["title"], record["name"])
```

```
session.close()
```

```
C:\Zoran\code\neo4j>python test-driver.py
```

```
King Arthur
```


Py2Neo Client

```
from py2neo import Graph, Path, Node
graph = Graph("http://neo4j:neo4j@localhost:7474")
for name in ["Alice", "Bob", "Carol"]:
    person = Node("person", name=name)
    graph.create(person)
    print name
```

```
friends = Path(Node("Person", name="Alice"), "KNOWS",
Node("Person", name="Bob"), "KNOWS", Node("Person", name="Carol"))
graph.create(friends)
```

- If we run this Python script and then make a query in Neo4J Browser:

```
match(n:Person) where n.name="Alice" or n.name="Bob" or n.name = "Carol"
return n;
```

- We would see "friends"

```
$ match(n:Person) where n.name="Alice" or
n.name="Bob" or n.name = "Carol" return n;
```

```
match(n:Person) where n.name="Alice" or n.name="Bob" or n.name = "Carol" return n;
```

*(3) Person(3)

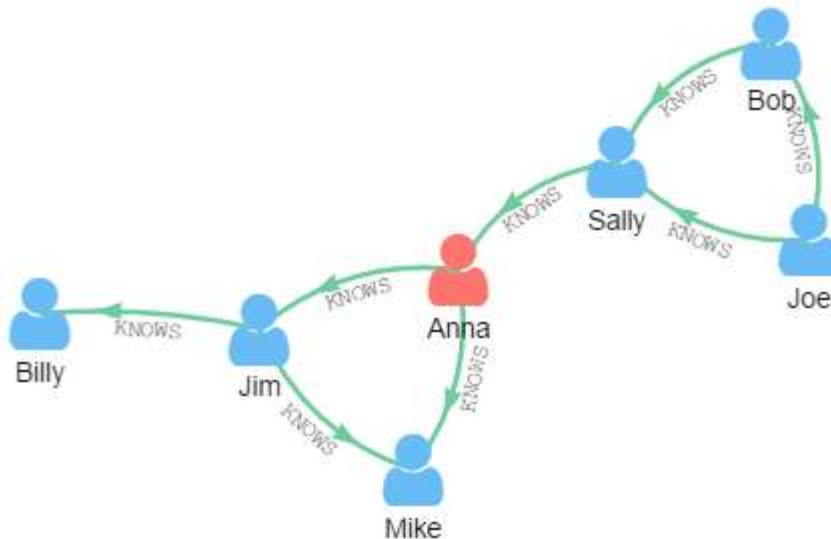
*(2) KNOWS(2)



Friends of Friends

- Find all of Joe's second-degree friends

```
MATCH (person:Person)-[:KNOWS]-(friend:Person)-[:KNOWS]-(foaf:Person)
WHERE
person.name = "Joe"
AND NOT (person)-[:KNOWS]-(foaf) RETURN foaf;
```



Java APIs

- There are several ways of using Neo4j from Java and other languages.
- You can see a list of options on this page:

<http://neo4j.com/developer/java/# neo4j for java developers>

- The standalone Neo4j Server can be installed on any machine and then accessed via its HTTP API from any language. There appear to be REST libraries for many languages: Java, JavaScript, PHP, Ruby, Scala, .Net, ...
- The dedicated Neo4j drivers go beyond that by offering comprehensive APIs for integrating with graph-based applications.
- One can also run Neo4j embedded in your JVM process, much like HSQL or Derby. This is great for unit testing, but also for high performance / no-network setups.
- If you use [Neo4j Embedded](#), you can use the Neo4j Core-Java-API directly.
- Besides an object-oriented approach to the graph database, working with Nodes, Relationships, and Paths, it also offers highly customizable high-speed traversal- and graph-algorithm implementations.
- You can also choose from any useful drivers wrapping that API, which exist either for specific programming languages or that add interesting functionality.

Use REST API to communicate with the Server

- The Neo4j REST API is designed to be used with any client that can send HTTP requests and receive HTTP responses.
- Existing Neo4J REST API is designed with discoverability in mind, so that you can start with a GET and from there discover URIs to perform other requests.
- The existing APIs are subject to change in the future, so for future-proofness discover URIs where possible, instead of relying on the current layout.
- The default representation is json, both for responses and for data sent with POST/PUT requests.
- To interact with the JSON interface you must explicitly set the request header [Accept:application/json](#) for those requests that responds with data.
- You should also set the header [Content-Type:application/json](#) if your request sends data, for example when you're creating a relationship.
- The server supports streaming results, with better performance and lower memory overhead.

Neo4J REST API

- The default way to interact with Neo4j is by using REST endpoint.
- If you go to Cypher Browser and elect `</>` Code screen you will see that Cypher Browser is submitting its queries as REST requests.
- The Neo4j transactional HTTP endpoint allows you to execute a series of Cypher statements within the scope of a transaction. The transaction may be kept open across multiple HTTP requests, until the client chooses to commit or roll back. Each HTTP request can include a list of statements, and for convenience you can include statements along with a request to begin or commit a transaction.
- The server guards against orphaned transactions by using a timeout. If there are no requests for a given transaction within the timeout period, the server will roll it back.

Transaction Endpoint

- If there is no need to keep a transaction open across multiple HTTP requests, you can begin a transaction, execute statements, and commit with just a single HTTP request.
- Example request. Note that endpoint contains /transaction/commit, instructing the server to commit whatever it receives with this statement.

```
POST http://localhost:7474/db/data/transaction/commit
Accept: application/json; charset=UTF-8
Content-Type: application/json
{
  "statements" : [ {
    "statement" : "CREATE (n:Apple) RETURN id(n) "
  } ]
}
```

- Example response

```
201: OK
Content-Type: application/json
{
  "results" : [ {
    "columns" : [ "id(n)" ],
    "data" : [ {
      "row" : [ 18 ]
    } ]
  } ],
  "errors" : [ ]
}
```

Use curl to submit REST Statements

- We need to find a way to send POST request to neo4j server. Curl is one such tool. You can run curl on any OS. You can install it on Cygwin. It comes with many Python distributions, like Anaconda. Some of many curl options are

```
C:\> curl -help
```

```
Usage: curl [options...] <url>
```

```
Options: (H) means HTTP/HTTPS only, (F) means FTP only
```

```
-K, --config FILE    Read config from FILE
```

```
-d, --data DATA      HTTP POST data (H)
```

```
    --data-raw DATA  HTTP POST data, '@' allowed (H)
```

```
    --data-ascii DATA HTTP POST ASCII data (H)
```

```
    --data-binary DATA HTTP POST binary data (H)
```

```
    --data-urlencode DATA HTTP POST data url encoded (H)
```

```
-f, --fail           Fail silently (no output at all) on HTTP errors (H)
```

```
    --false-start    Enable TLS False Start.
```

```
-F, --form CONTENT   Specify HTTP multipart POST data (H)
```

```
    --form-string STRING Specify HTTP multipart POST data (H)
```

```
    --ftp-account DATA Account data string (F)
```

```
    --ftp-create-dirs Create the remote dirs if not present (F)
```

```
    --ftp-method [MULTICWD/NOCWD/SINGLECWD] Control CWD usage (F)
```

```
-P, --ftp-port ADR   Use PORT with given address instead of PASV (F)
```

```
-G, --get            Send the -d data with a HTTP GET (H)
```

```
-H, --header LINE    Pass custom header LINE to server (H)
```

```
-i, --include        Include protocol headers in the output (H/F)
```

```
-o, --output FILE    Write to FILE instead of stdout
```

```
-X, --request COMMAND Specify request command to use
```

Creating a Node with curl

- Let us try running the following on a single line:

```
curl -i -H accept:application/json -H content-type:application/json -XPOST
http://localhost:7474/db/data/transaction/commit -d
'{"statements":[{"statement":"CREATE (p:Strawberry) RETURN p"}]}'
```

- We are passing 2 headers accept and content-type (-H options), submitting a POST request (-XPOST), providing auto-commit URL endpoint and submitting data (-d option).
- With curl I installed on Windows I am getting an error.
- With curl that came with my Anaconda Python 3.4 visible to my Cygwin, the command runs with response.

```
$ which curl
```

```
/cygdrive/e/Programs/Anaconda3/Library/bin/curl
```

```
$ curl -i -H accept:application/json -H content-type:application/json -XPOST
http://localhost:7474/db/data/transaction/commit -d
'{"statements":[{"statement":"CREATE (p:Peach) RETURN p"}]}'
```

% Total	% Received	% Xferd	Average Speed	Time	Time	Time	Current
			Dload	Upload	Total	Spent	Left
100	123	100	65	100	58	245	218
--:--:--	--:--:--	--:--:--	--:--:--	--:--:--	--:--:--	--:--:--	--:--:--

```
245HTTP/1.1 200 OK
```

```
Date: Sat, 09 Apr 2016 12:52:37 GMT
```

```
Content-Type: application/json
```

```
Access-Control-Allow-Origin: *
```

```
Content-Length: 65
```

```
Server: Jetty(9.2.z-SNAPSHOT)
```

```
{"results":[{"columns":["p"],"data":[{"row":[]}]}], "errors":[]}
```


Verifying the result

- We can verify that our Peach node is created in Cypher Browser or run a curl command:

```
$ curl -i -H accept:application/json -H content-type:application/json -XPOST http://localhost:7474/db/data/transaction/commit -d '{"statements":[{"statement":"MATCH (p:Strawberry) RETURN p"}]}'
```

% Total	% Received	% Xferd	Average	Speed	Time	Time	Time
Current			Dload	Upload	Total	Spent	Left
Speed							
100	122	100	65	100	57	260	228
--:--:--	--:--:--	--:--:--	--:--:--	--:--:--	--:--:--	--:--:--	--:--:--

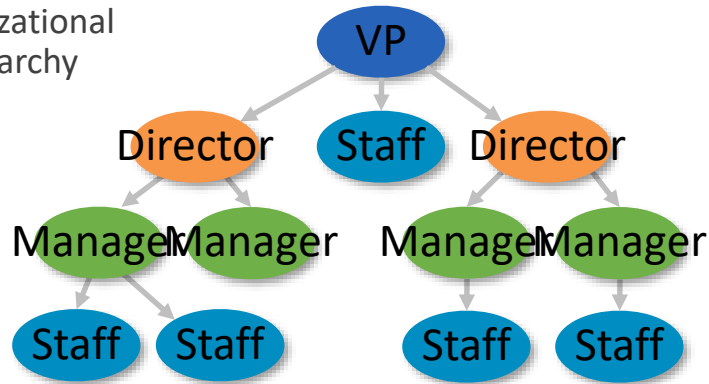
```
260HTTP/1.1 200 OK
Date: Sat, 09 Apr 2016 12:58:35 GMT
Content-Type: application/json
Access-Control-Allow-Origin: *
Content-Length: 65
Server: Jetty(9.2.z-SNAPSHOT)

{"results":[{"columns":["p"],"data":[{"row":[{"p":"Strawberry"}]}]}], "errors":[]}
```

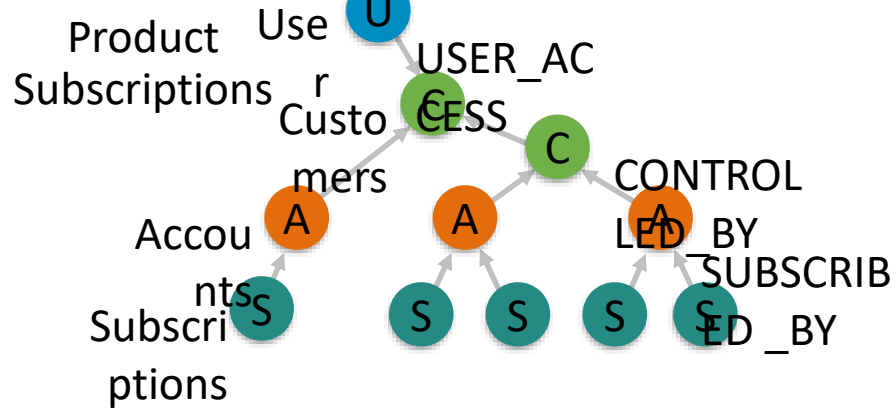
Use Case Graphs for Master Data Management

MDM Solutions with Graph Databases

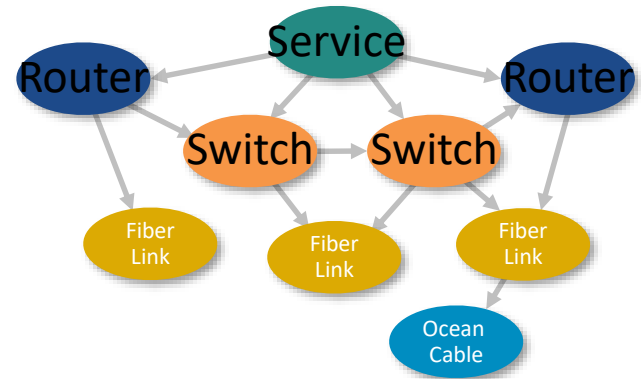
Organizational Hierarchy



Social Networks

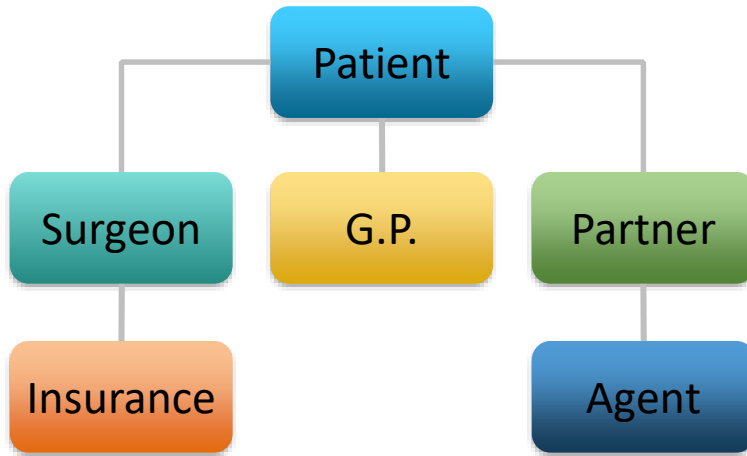


CMDB Network Inventory

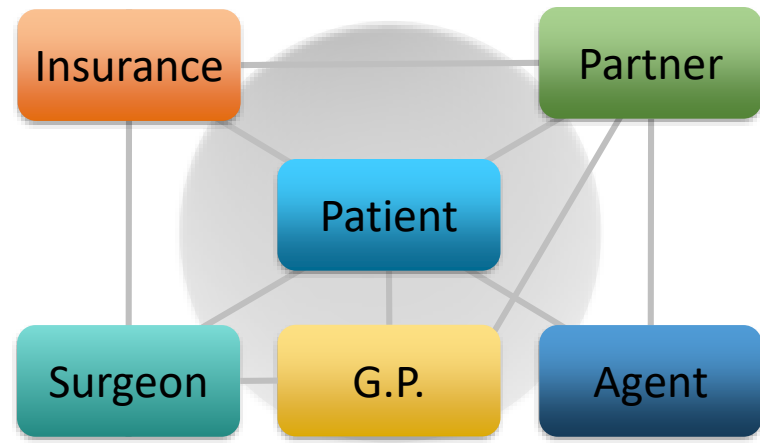


MDM Isn't Hierarchical

Typical MDM system structure



...but MDM is really a network



Gartner®

Challenges with Current MDM Systems

Lack of support for non-hierarchical or matrix data relationships

- Master data is never strictly hierarchical
- Systems are designed for fixed top-down hierarchy
- Non-hierarchical data is not supported

Inability to unlock value from data relationships

- Systems store only very simple data relationships
- Complex relationships and links not stored

Inflexible and expensive to maintain

- Changes to the model are expensive and time-consuming

Use Case

Graphs for Network and IT Operations Management

Network Graphs – Telco Example

PROBLEM

Need: Instantly diagnose problems in networks of 1B+ elements

But: Basing diagnosis solely on streaming machine data severely limits accuracy and effectiveness

SOLUTION

Real-time graph analytics provide actionable insight for the largest complex connected networks in the world

- The entire network lives in a graph
- Analyzes dependencies in real time
- Highly scalable with carrier-grade uptime requirements



User Case Graphs for Fraud Detection

Fraud Scenarios

Retail First Party Fraud

- Opening many lines of credit with no intention of paying back
- Accounts for \$10B+ in annual losses at US banks⁽¹⁾



Synthetic Identities and Fraud Rings

- Rings of synthetic identities committing fraud



Insurance – Whiplash for Cash

- Insurance scams using fake drivers, passengers and witnesses
- Increase network efficiency



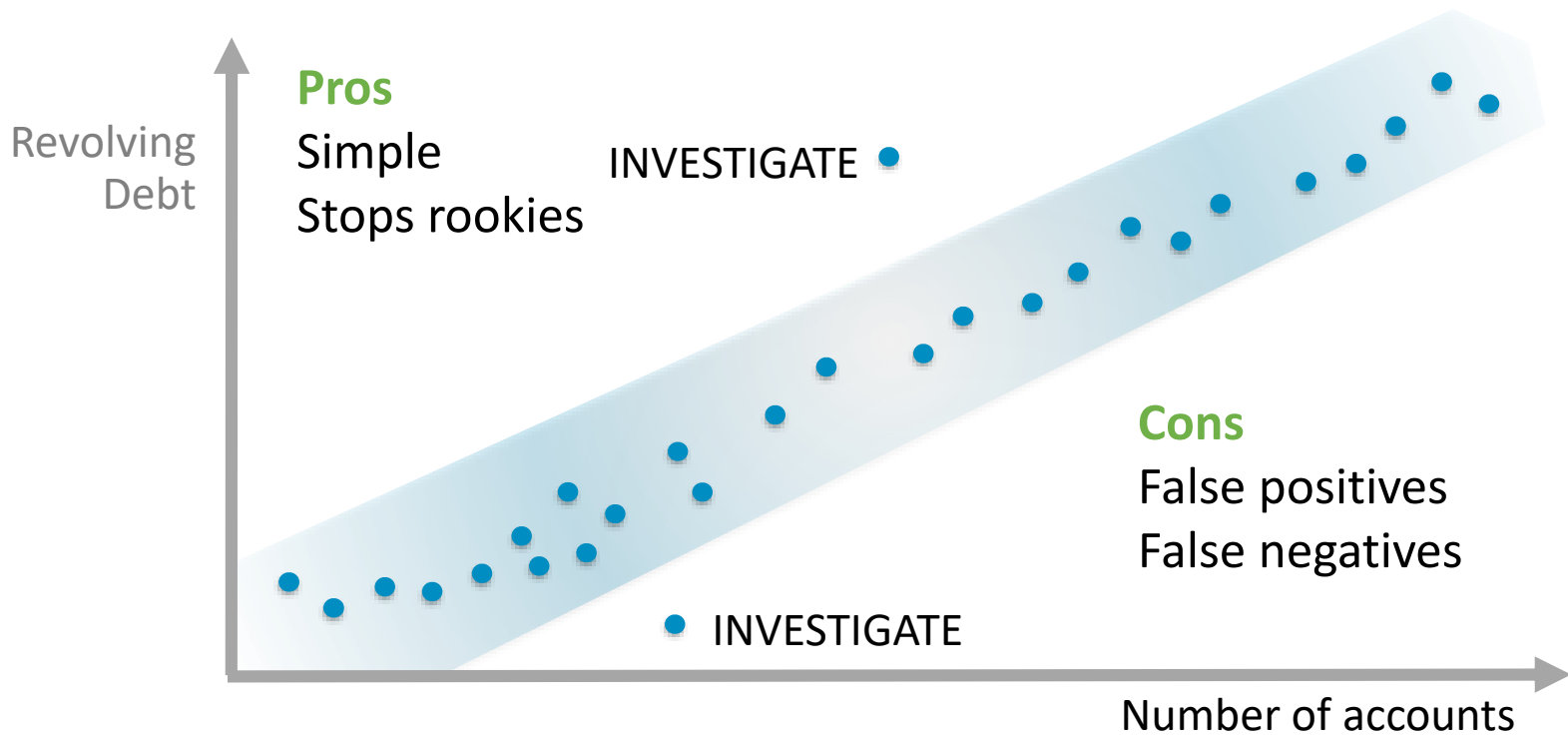
eCommerce Fraud

- Online payment fraud

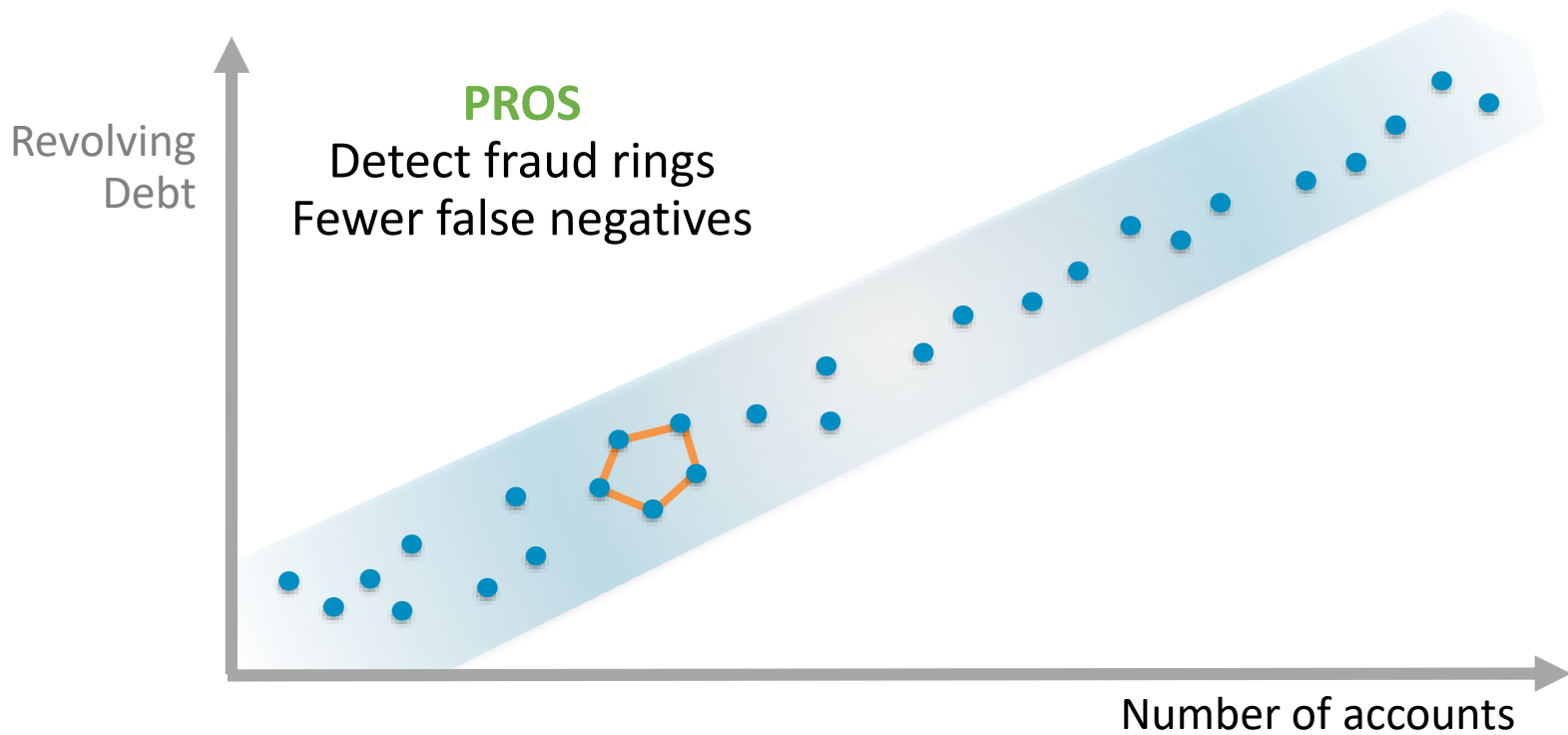


(1) Business Insider: <http://www.businessinsider.com/how-to-use-social-networks-in-the-fight-against-first-party-fraud-2011-3>

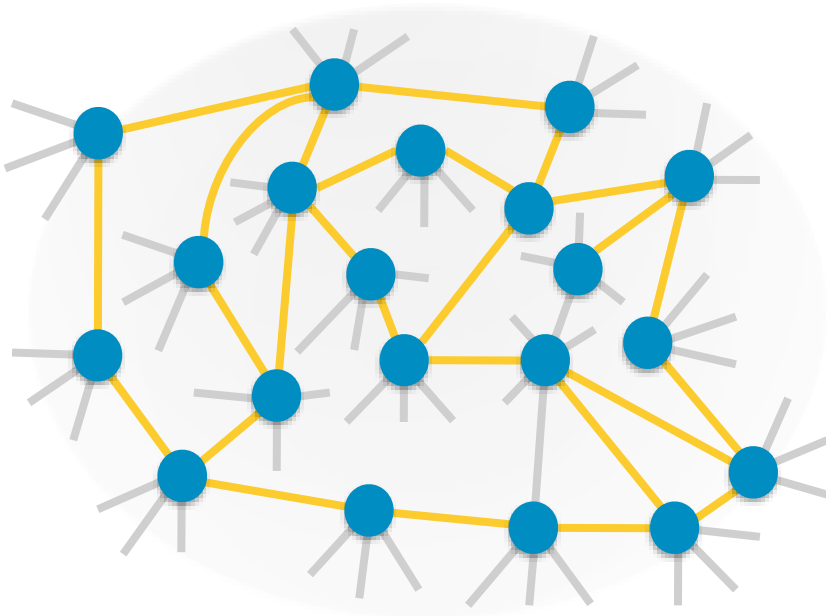
Discrete Data Analysis



Connected Analysis



Connected Analysis with Neo4j



Value

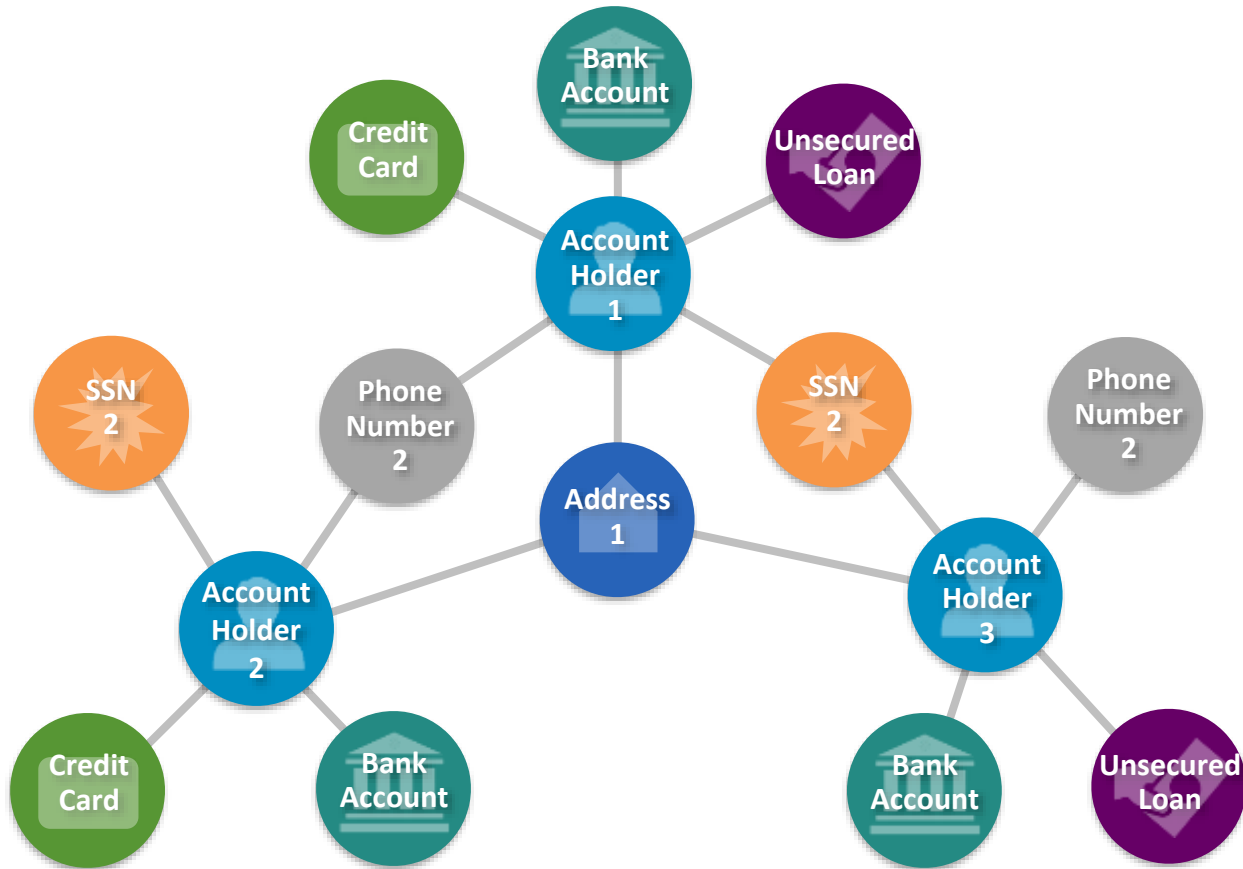
Effective in detecting some of the most impactful attacks, even from organized rings

Challenge

Extremely difficult with traditional technologies

*For example a **ten-person fraud bust-out** is **\$1.5M**, assuming 100 false identities and 3 financial instruments per identity, each with a \$5K credit limit*

Modeling a Fraud Ring as a Graph



Use Case Graphs for Real-time Recommendations

Real-Time Recommendations - Benefits



Online Retail

- Suggest related products and services
- Increase revenue and engagement



Media and Broadcasting

- Create an engaging experience
- Produce personalized content and offers



Logistics

- Recommend optimal routes
- Increase network efficiency

Real-Time Recommendations - Challenges

Make effective real-time recommendations

- Timing is everything in point-of-touch applications
- Base recommendations on current data, not last night's batch load

Process large amounts of data and relationships for context

- Relevance is king: Make the right connections
- Drive traffic: Get users to do more with your application

Accommodate new data and relationships continuously

- Systems get richer with new data and relationships
- Recommendations become more relevant



Walmart – Retail Recommendations

- Needed online customer recommendations to keep pace with competition
- Data connections provided predictive context, but were not in a usable format
- Solution had to serve many millions of customers and products while maintaining superior scalability and performance



World's largest company
by revenue

World's largest retailer and
private employer

SF-based global
e-commerce division
manages several websites

Found in 1969
Bentonville, Arkansas