Property Graphs: Neo4j and Cypher





Neo4j - the most commonly used graph DBMS, by far

https://db-engines.com/en/ranking/graph+dbms

Model: property graphs (graphs with nodes and edges carrying multiple data values arranged as key-value pairs)

Query language: Cypher.

ASCII-art pattern matching + usual database features

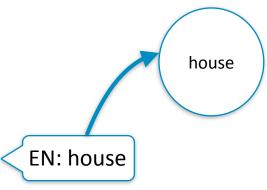




EN: house

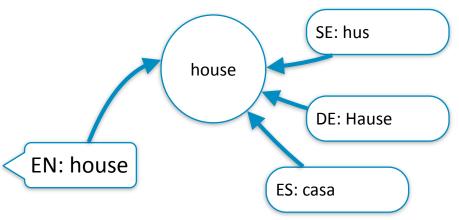






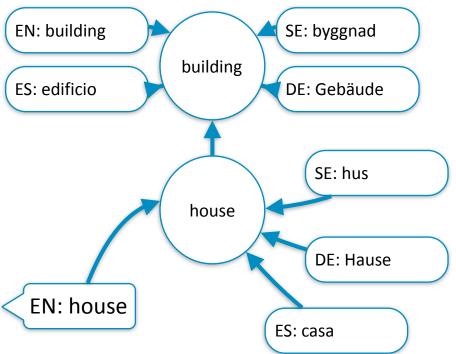


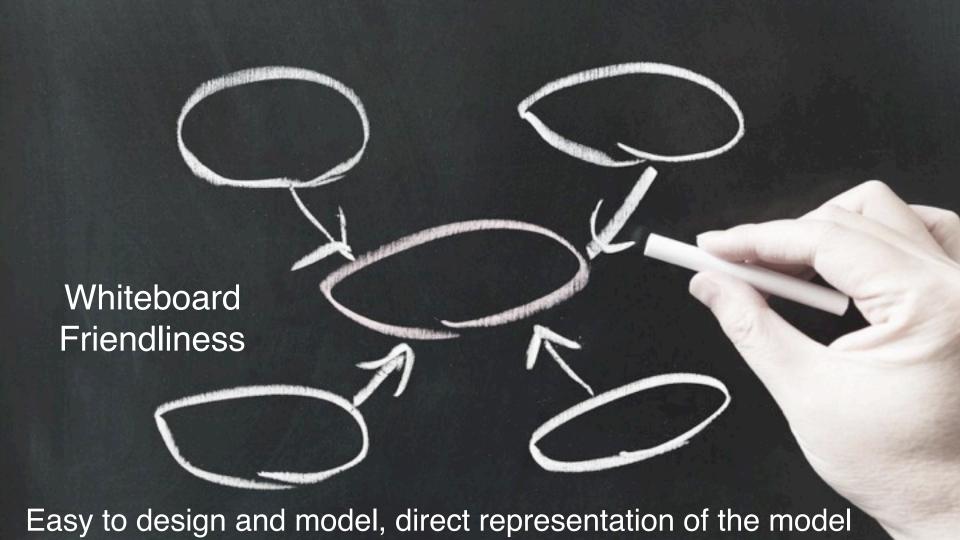




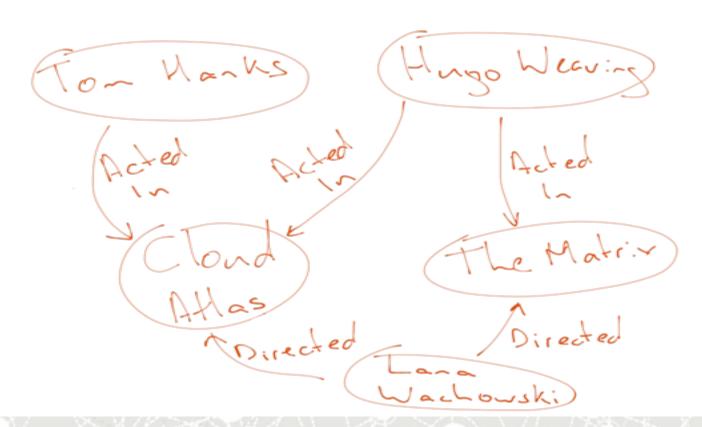




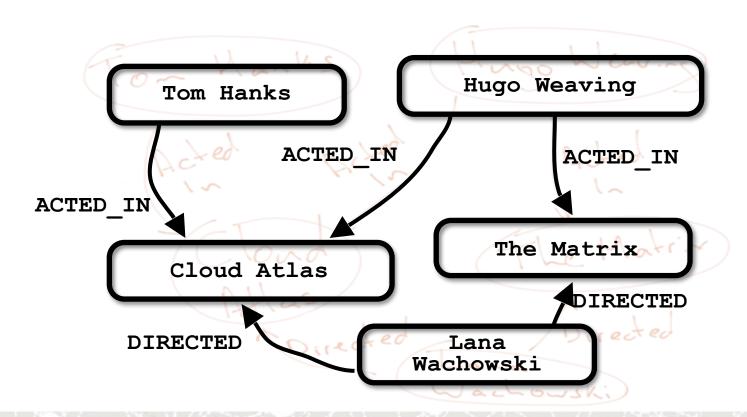




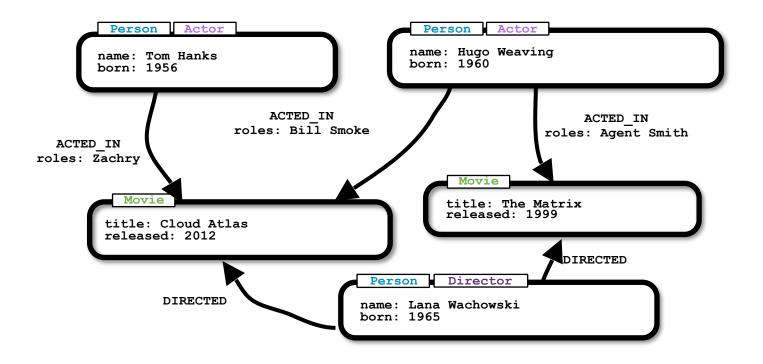




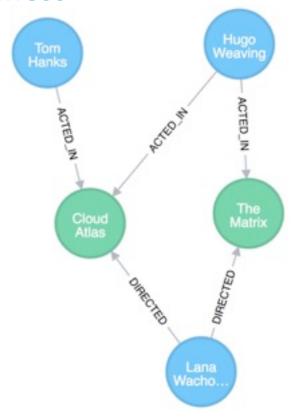












Intro to the property graph model



Neo4j Fundamentals



- Nodes
- Relationships
- Properties
- Labels

Property Graph Model Components



Nodes

- Represent the objects in the graph
- Can be labeled







Property Graph Model Components

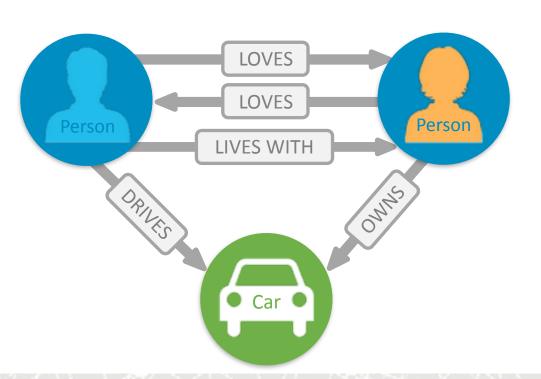


Nodes

- Represent the objects in the graph
- Can be labeled

Relationships

Relate nodes by type and direction



Property Graph Model Components



Nodes

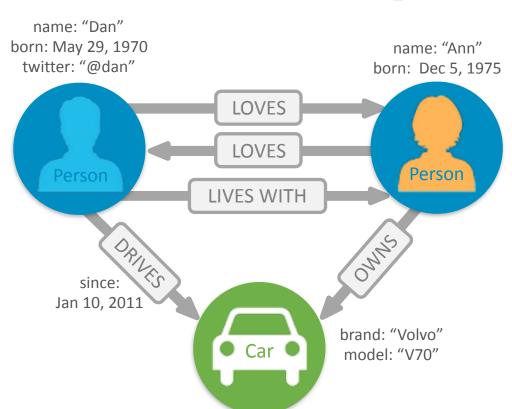
- Represent the objects in the graph
- Can be labeled

Relationships

Relate nodes by type and direction

Properties

 Name-value pairs that can go on nodes and relationships.



Summary of the graph building blocks



- Nodes Entities and complex value types
- Relationships Connect entities and structure domain
- Properties Entity attributes, relationship qualities, metadata
- Labels Group nodes by role

Graph Querying



Why not SQL?



SQL is inefficient in expressing graph pattern queries

In particular for queries that

- are recursive
- can accept paths of multiple different lengths
- Graph patterns are more intuitive and declarative than joins
- SQL cannot handle path values

Cypher

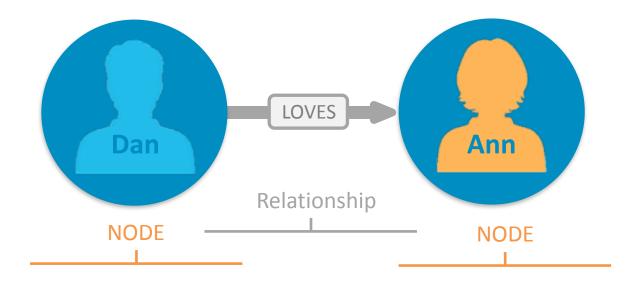


A pattern matching query language made for graphs

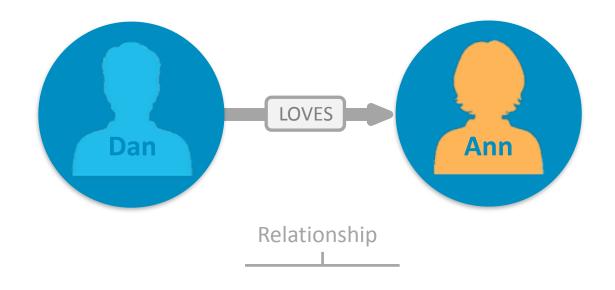
- Declarative
- Expressive
- Pattern Matching

Pattern in our Graph Model

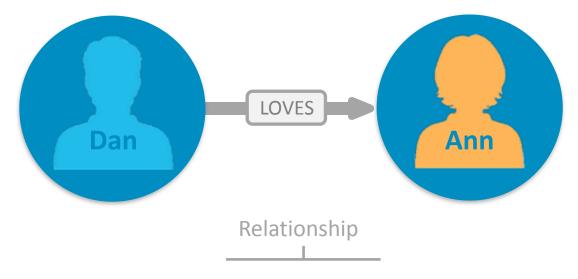






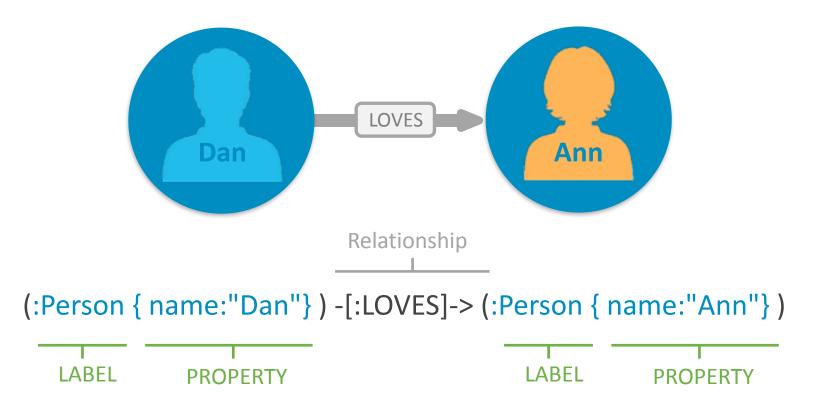




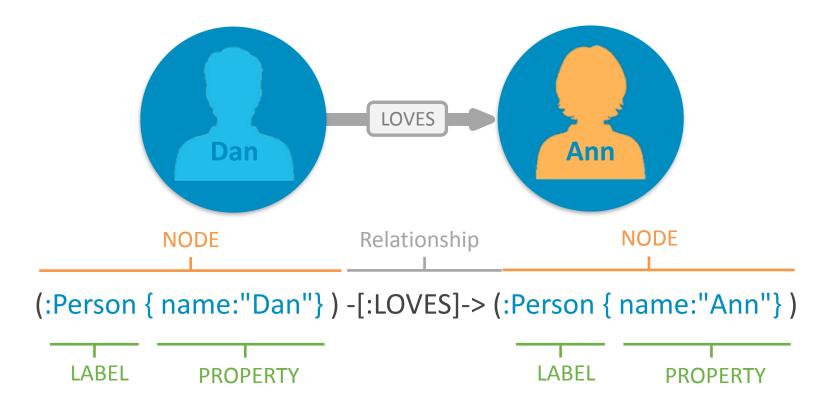


(:Person { name:"Dan"}) -[:LOVES]-> (:Person { name:"Ann"})

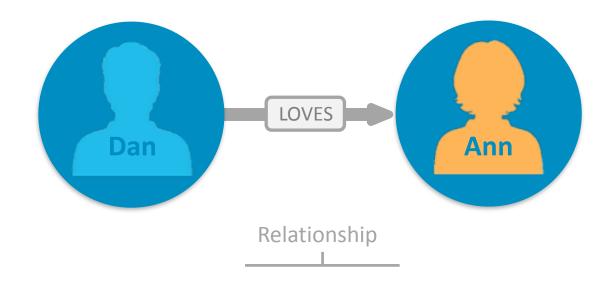




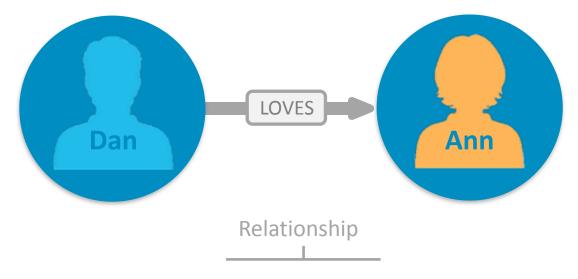






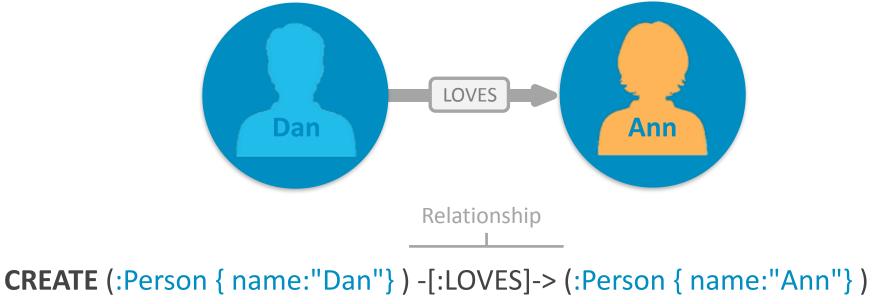






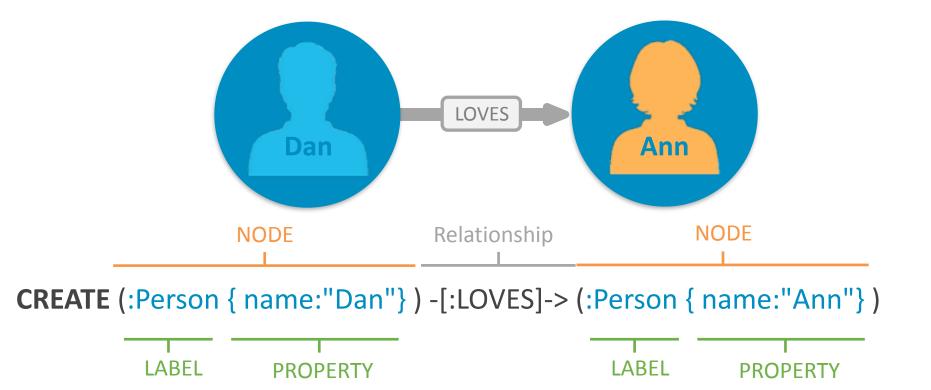
CREATE (:Person { name:"Dan"}) -[:LOVES]-> (:Person { name:"Ann"})



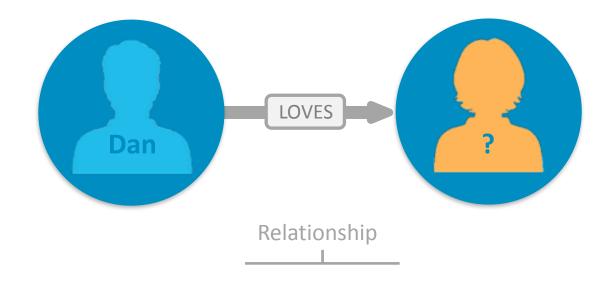


LABEL PROPERTY LABEL PROPERTY

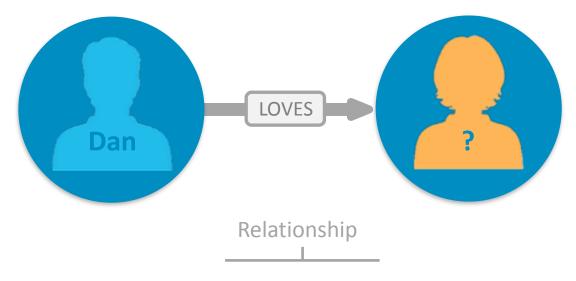






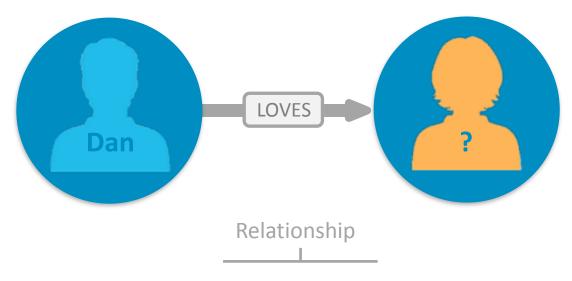






MATCH (:Person { name:"Dan"}) -[:LOVES]-> (whom) RETURN whom

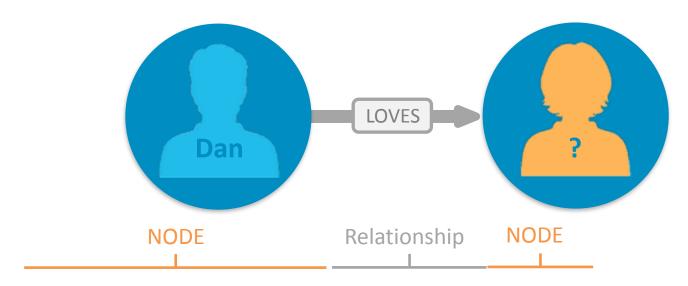




MATCH (:Person { name:"Dan"}) -[:LOVES]-> (whom) RETURN whom







MATCH (:Person { name:"Dan"}) -[:LOVES]-> (whom) RETURN whom



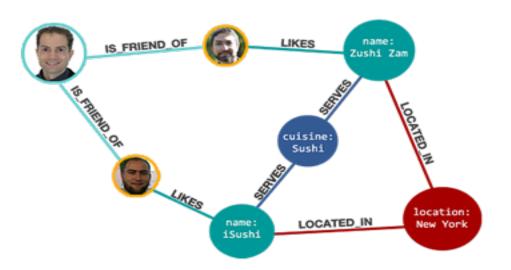
A graph query example



A social recommendation







A social recommendation







The Syntax

Nodes



Nodes are drawn with parentheses.

()

Relationships



Relationships are drawn as arrows, with additional detail in brackets.

Patterns



Patterns are drawn by connecting nodes and relationships with hyphens, optionally specifying a direction with > and < signs.

The components of a Cypher query



MATCH (m:Movie)
RETURN m

MATCH and RETURN are Cypher keywords

m is a variable

:Movie is a node label

The components of a Cypher query



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

MATCH and RETURN are Cypher keywords

p, r, and m are variables

:Movie is a node label

:ACTED_IN is a relationship type

The components of a Cypher query



```
MATCH path = (:Person)-[:ACTED_IN]->(:Movie)
RETURN path
```

MATCH and RETURN are Cypher keywords

path is a variable

:Movie is a node label

:ACTED_IN is a relationship type

Graph versus Tabular results



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

Graph versus Tabular results



```
MATCH (m:Movie)
RETURN m.title, m.released
```

Properties are accessed with {variable}.{property_key}

Case sensitivity



Case sensitive

Case insensitive

Node labels

Cypher keywords

Relationship types

Property keys

Case sensitivity



Case sensitive

Case insensitive

:Person

MaTcH

:ACTED_IN

return

name



Write queries

The CREATE Clause



```
CREATE (m:Movie {title:'Mystic River', released:2003})
RETURN m
```

The SET Clause



```
MATCH (m:Movie {title: 'Mystic River'})
SET m.tagline = 'We bury our sins here, Dave. We wash them clean.'
RETURN m
```

The CREATE Clause



```
MATCH (m:Movie {title: 'Mystic River'})

MATCH (p:Person {name: 'Kevin Bacon'})

CREATE (p)-[r:ACTED_IN {roles: ['Sean']}]->(m)

RETURN p, r, m
```



```
MERGE (p:Person {name: 'Tom Hanks'})
RETURN p
```



```
MERGE (p:Person {name: 'Tom Hanks', oscar: true})
RETURN p
```



```
MERGE (p:Person {name: 'Tom Hanks', oscar: true})
RETURN p
```

There is not a :Person node with name: 'Tom Hanks' and oscar:true in the graph, but there is a :Person node with name: 'Tom Hanks'.

What do you think will happen here?



```
MERGE (p:Person {name: 'Tom Hanks'})
SET p.oscar = true
RETURN p
```





There is not a :Movie node with title: "The Terminal" in the graph, but there is a :Person node with name: "Tom Hanks".

What do you think will happen here?



```
MERGE (p:Person {name: 'Tom Hanks'})
MERGE (m:Movie {title: 'The Terminal'})
MERGE (p)-[r:ACTED_IN]->(m)
RETURN p, r, m
```

ON CREATE and ON MATCH



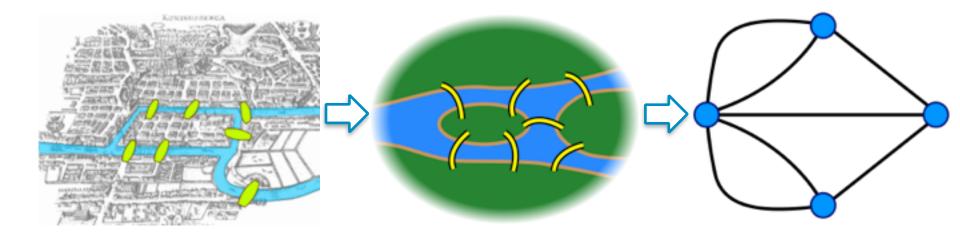
```
MERGE (p:Person {name: 'Your Name'})
ON CREATE SET p.created = timestamp(), p.updated = 0
ON MATCH SET p.updated = p.updated + 1
RETURN p.created, p.updated;
```

Graph Modeling



Models

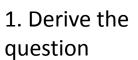




The modeling workflow









2. Obtain the data



3. Develop a model





4. Ingest the data



5. Query/ Prove our model

Developing the model and the query



- 1. Identify application/end-user goals
- 2. Figure out what questions to ask of the domain
- 3. Identify entities in each question
- 4. Identify relationships between entities in each question
- 5. Convert entities and relationships to paths
 - These become the basis of the data model
- 6. Express questions as graph patterns
 - These become the basis for queries

1. Application/End-User Goals



As an employee

I want to know who in the company has similar skills to me

So that we can exchange knowledge

2. Questions to ask of the Domain



As an employee

I want to know who in the company has similar skills to me

So that we can exchange knowledge

Which people, who work for the same company as me, have similar skills to me?

3. Identify Entities



Which people, who work for the same company as me, have similar skills to me?

- Person
- Company
- Skill

4. Identify Relationships Between Entities



Which people, who work for the same company as me, have similar skills to me?

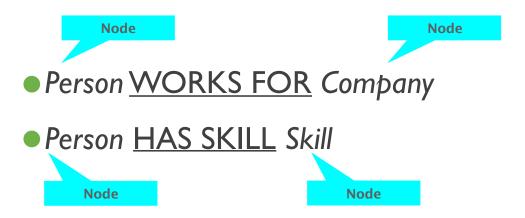
- Person WORKS FOR Company
- Person HAS SKILL Skill



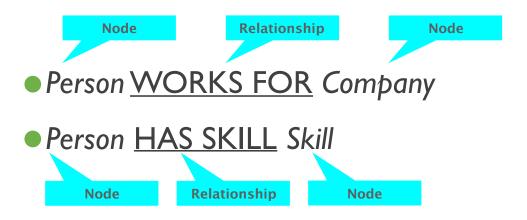


- Person WORKS FOR Company
- Person HAS SKILL Skill

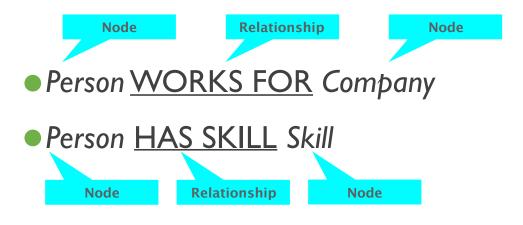












- (:Person)-[:WORKS_FOR]->(:Company),
- (:Person)-[:HAS_SKILL]->(:Skill)

5. Convert to Cypher Paths



Relationship Node Node Person WORKS FOR Company Person HAS SKILL Skill Node Relationship **Node** Label Label (:Person)-[:WORKS_FOR]->(:Company), (:Person)-[:HAS_SKILL]->(:Skill) Label Label

5. Convert to Cypher Paths



Relationship Node Node Person WORKS FOR Company Person HAS SKILL Skill Relationship Node Node **Relationship Type** Label Label (:Person)-[:WORKS FOR]->(:Company), (:Person)-[:HAS_SKILL]->(:Skill) Label **Relationship Type** Label

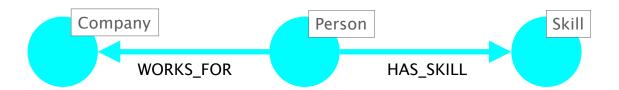
Consolidate Pattern



```
(:Person)-[:WORKS_FOR]->(:Company),
(:Person)-[:HAS_SKILL]->(:Skill)
```



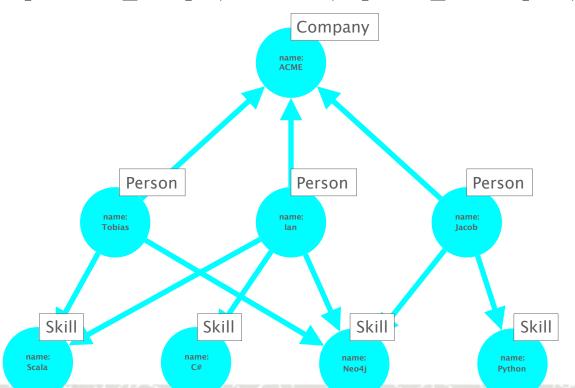
(:Company)<-[:WORKS_FOR]-(:Person)-[:HAS_SKILL]->(:Skill)



Candidate Data Model



(:Company)<-[:WORKS_FOR]-(:Person)-[:HAS_SKILL]->(:Skill)

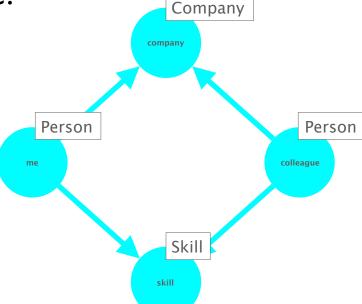


6. Express Question as Graph Pattern



Which people, who work for the same company as me, have

similar skills to me?



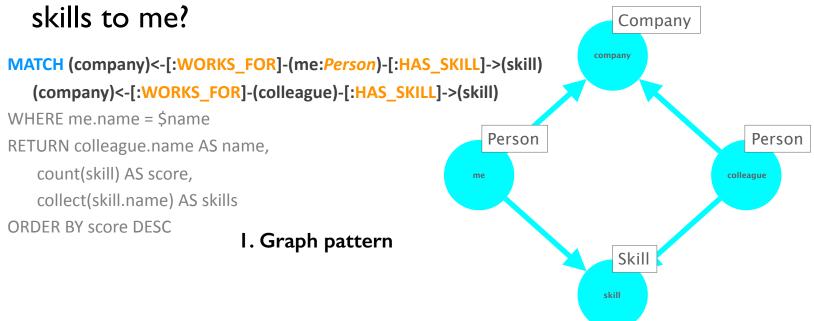


Which people, who work for the same company as me, have similar

skills to me? Company MATCH (company)<-[:WORKS FOR]-(me:Person)-[:HAS SKILL]->(skill) (company)<-[:WORKS FOR]-(colleague)-[:HAS SKILL]->(skill) WHERE me.name = \$name Person Person RETURN colleague.name AS name, count(skill) AS score, me colleague collect(skill.name) AS skills **ORDER BY score DESC** Skil skill

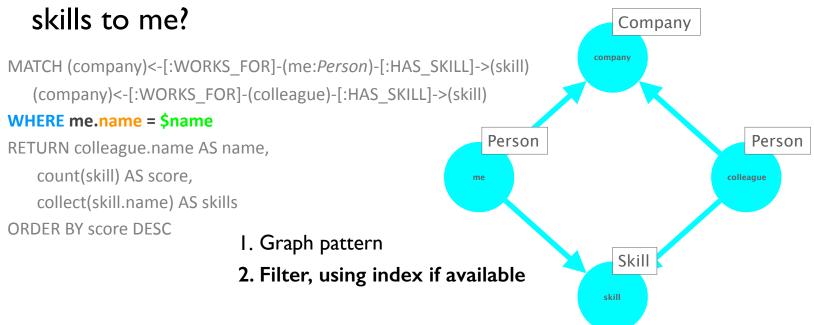


Which people, who work for the same company as me, have similar



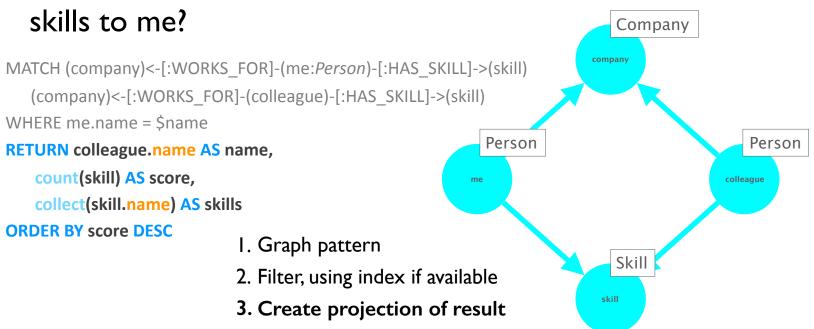


Which people, who work for the same company as me, have similar



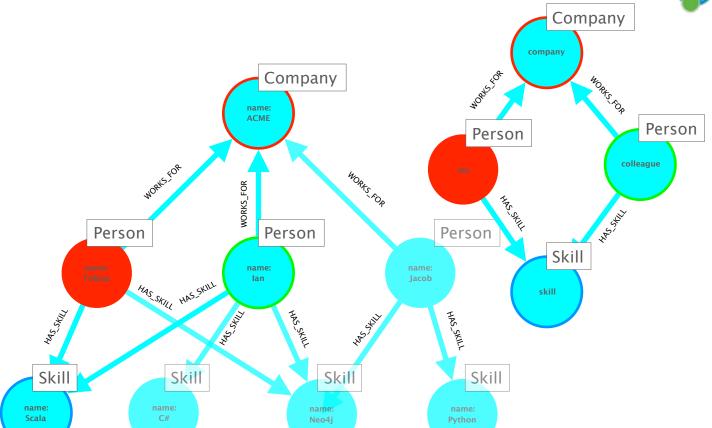


Which people, who work for the same company as me, have similar



First Match





Second Match Company company Company name: **ACME** Person Person colleague WORKS_FOR Person Person Person Skill name: lan lacob skill HAS SKILL HAS_SKILL Skill Skill Skill Skill name: name:

Neo4j

Python

Scala

C#

Third Match Company company Company name: **ACME** Person Person colleague WORKS_FOR Person Person Person Skill name: Jacob lan skill HAS SKILL HAS_SKILL Skill Skill Skill Skill

name:

Neo4j

Python

Scala

C#

Result of the Query





Modeling exercise: Movie genres

Adding movie genres



The question: should we model them as properties or as nodes?

genre: ["Action", "Sci-Fi"] Movie name: "Action" Genre Kname: "Sci-Fi"

Genres as properties



```
MATCH (m:Movie {title: 'The Matrix'})
SET m.genre = ['Action', 'Sci-Fi']
RETURN m
```

Genres as properties



```
MATCH (m:Movie {title: 'Mystic River'})
SET m.genre = ['Action', 'Mystery']
RETURN m
```

The good side of properties



Accessing a movie's genres is quick and easy.

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

The bad side of properties



Finding movies that share genres is painful and we have a disconnected pattern in the MATCH clause - a sure sign you have a modeling issue.

```
MATCH (m1:Movie), (m2:Movie)
WHERE any(x IN m1.genre WHERE x IN m2.genre)
AND m1 <> m2
RETURN m1, m2;
```

Genres as nodes



```
MATCH (m:Movie {title:"The Matrix"})
MERGE (action:Genre {name:"Action"})
MERGE (scifi:Genre {name:"Sci-Fi"})
MERGE (m)-[:IN_GENRE]->(action)
MERGE (m)-[:IN_GENRE]->(scifi)
```

Genres as nodes



```
MATCH (m:Movie {title:"Mystic River"})
MERGE (action:Genre {name:"Action"})
MERGE (mystery:Genre {name:"Mystery"})
MERGE (m)-[:IN_GENRE]->(action)
MERGE (m)-[:IN_GENRE]->(mystery)
```

The good side of nodes



Finding movies that share genres is a natural graph pattern.

The (not too) bad side of nodes



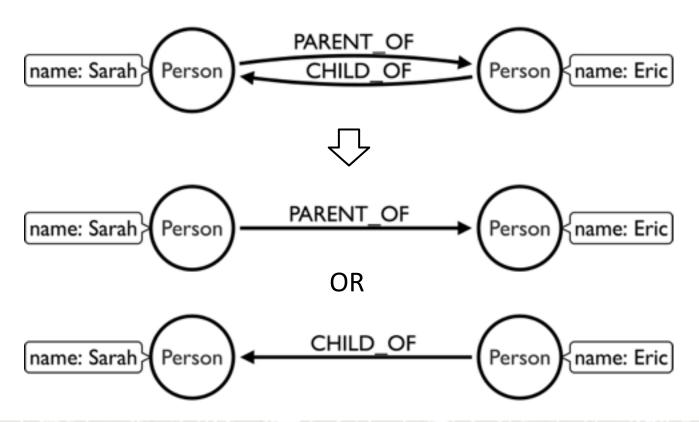
Accessing the genres of movies requires a bit more typing.



Symmetric Relationships

Symmetric relationships



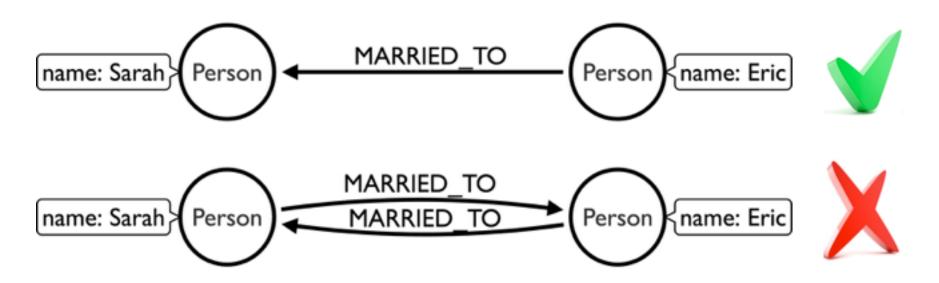




Bidirectional Relationships

Use single relationship and ignore direction in queries





MATCH (:Person {name:'Eric'})-[:MARRIED_TO]-(p2)
RETURN p2

Formal semantics

Nadime Francis,
 Paolo Guagliardo,
 Leonid Libkin

 Formally defines a (large) core of Cypher



4.3.2 Pattern Matching Clauses

• [MATCH
$$\pi$$
]_G(T) = $\bigcup_{\substack{u \in T \\ \pi' \in \operatorname{inpliff}) \\ \operatorname{trajes of paths } p}} \{(u, u') \mid u' \text{ is uniform with free}(u, \pi') \\ \operatorname{and} (\beta, G, (u, u')) \models \pi'$

Remark that even though both u' and $\bar{\pi}'$ range over infinite sets, only a finite number of values contribute to a non-empty set in the final union. Thus $[MATCS \#]_C(T)$ is finite, provided that T is finite.

• [OPTIONAL MATCH
$$\pi$$
]_G(T) = $\bigcup_{u \in T} T_u$
where $T_u = \begin{cases} [\text{MATCH } \pi]_G(\{u\}) & \text{if } [\text{MATCH } \pi]_G(\{u\}) \neq \emptyset \\ \{u, (\text{free}(u, \pi) : \text{mull})\} & \text{otherwise} \end{cases}$

Example 4. Let G be the property graph defined in Example 2 and consider the following clause:

EATCH
$$(x) - [:knows*] -> (y)$$

Let T be the table $T = \{(x : n_1); (x : n_4)\}$. The goal of this example is to compute [MITCH $(x) - [:knows] -> (y)]_{C}(T)$.

First, remark that rigid((x) - [:knows*] -> (y)) is the (infinite) set of all rigid paths π_d of the form (x) - [:knows*d] -> (y), for d>0. Notice however that these patterns can only be satisfied by paths which contains exactly d distinct relationships. Since G only contains 3 relationships, we can immediately deduce that only π_1 , π_2 and π_3 can contribute non-empty lags to the final result.

Let us first look at the case where $u=(x:n_1)$, $\pi'=\pi_1$ and $p=n_1r_1n_2$. Remark that first(u, π_1) = {y}, thus u', if it exists, is a record over the set of fields {y}. Additionally, we can check that $(n_1r_1n_2, G, (x:n_1, y:n_2)) \models \pi_1$. It is straightforward to check that n_2 is actually the only suitable value for y, and thus the contribution of this specific triple u, π' , p to the final result is precisely the bag { $(x:n_1), (y:n_2)$ }.

One can verify that no path p other than $n_1r_1n_2$ can contribute a record in the case where $u = (x : n_1)$ and $\pi' = \pi_1$. Indeed, π_1 requires both p to be of length 1 and to start at the node named x, which u evaluates to be n_1 .

Following a similar reasoning, we can now compute the contribution of the following triples:

(x:n₁, y:n₃), π₂, n₁r₁n₂r₂n₃ yields (x:n₁, y:n₃);



Future improvements

Regular Path Queries



- Conjunctive bi-directional Regular Path Queries with Data
 - Plus specification of a cost function for paths,
 which allows for more interesting notions of shortest path

MATCH (a)-/ \sim coauth* COST x/->(b) ORDER BY x LIMIT 10

Further improvements



- Support for querying from multiple graphs
- Support for returning graphs
- Support for defining views
- Integration with SQL

RDBMSs and Graphs



RDBMS can't handle relationships well



- Cannot model or store data and relationships without complexity
- Performance degrades with number and levels of relationships, and database size
- Query complexity grows with need for JOINs
- Adding new types of data and relationships requires schema redesign, increasing time to market

Express Complex Queries Easily with Cypher



Find all managers and how many people they manage, up to 3 levels down.

Cypher

SQL

```
SELECT T. direct Reporties: A5 direct Reportiess, surri? Journ 2 A5 count
                                                                                      SQUECT depth SReportness pict AG direct/Reportness.
                                                                                      count/depth/Maportees directly, manageril 45 count.
SEIDCT manager pid IIS directReportees, 0 IIS count
                                                                                      FIGM person, reporter manager
 FROM partico, reported manager
                                                                                      JON person, reporter L1 Asporters.
 WHIRE manager pld = (SELECT of FROM partor WHIRE name = "Riame Name")
                                                                                      (Ni manager directly, manages + Lifesportses, pld.
                                                                                      (ON person, reporter L) Reporters
                                                                                     CN LI Naportero directly , manages + LDRoportero, pill
 SOUTH manager and At-direct Reportions, countil manager almostly, manager) At-count
FROM person, reported manager
                                                                                      WHERE manager and a SIGLEST of FROM pursue WHERE name a "Rising Warse").
WHERE manager and + (SOLICT to FROM person WHERE name + "Thome Rome")
                                                                                      ORCUP In direct Reportees
GROUP BY physiologicals
                                                                                      SRCUP BY SINKSReporteed
SEDCT manager pid K6 direct/aportees, count/reportees directly, manager) K6 count.
IROM person, reported manager
                                                                                      (SE)SCTT directReportoes AS directReportoes, sumC asset) AS count
JON person, reported reported
CRI manager directly, managers - reportee and
                                                                                       SSUECT reporter directly, manager All directReporters, CAS count.
WARK manager pid + (MUSC) of FROM person WHIRE name + "Name Name")
                                                                                      EROSE parties, reported manager
                                                                                      ION person reporter reporter
                                                                                      (% manager directly, manages a reporter pill
MUSCT manager pid KS directAsportees, count), 2Reportees directly_manages( AS count
                                                                                     WHERE manager and a SIGUECT of HIGH person WHERE name a "Risane Rame").
FROM person, reportee manager
                                                                                      GROUP BY directhesomess
JOT person, reporter L1Reporters
CRI manager directly, managers + LTReportees, and
                                                                                      MUSC? L2Reportees pid 85 directReportees, count(),2Reportees directly, managed
JON person, reporter Lifeporters.
                                                                                      All-sount
ON Lifeportees directly, manages + Lifeportees pld
                                                                                     FBOM parson_reportes manager
WHERE manager and + CRELECT in FROM person WHERE name + "Warne Warne").
                                                                                      JON person, reporter L1 Reporters
                                                                                      (Ni manager directly , manages + 1.5 Reportures, pid.
URCUP BY directReportees
                                                                                      20% person, reporter LEReporters.
GACUP IN direct/seporteed
                                                                                      DN L1Reportees directly manages = L2Reportees pid
                                                                                      WHERE manager pid + (56,607 of FROM person WHERE name + "Name Name")
(MESCT T directReporters AS directReporters, sum(T sount) AS sount
                                                                                      DROUP BY directReportees
                                                                                      1.865 0
SEIDCT manager directly, manages AS directReportses, GAS count
                                                                                      GROUP BY directReportees)
FROM server, reported manager
                                                                                      LAHON
MARKE manager and + CREACT of FROM person WHERE name + "Name Name").
                                                                                      (MILECT LIMportees directly, manager, NJ directReportees, II NJ count.
                                                                                      FBCM person_reporter manager
SECOT reported pill AG direct Reporteds, countil reported, directly, managed AG count
                                                                                     JOR person reporter Lifeportess
FROM person, reportee manager
                                                                                      DN manager directly, manages + LI Reportees pid
2019 person, reporter reporter
                                                                                      30% person, regioner Lifegoriees.
Oh manager directly, manages + reporter and
                                                                                      CN L1Augurtees.directly_manages + L2Augurtees.pid
WHERE manager pid + (SOUCT id FROM person WHERE name + "Name Name")
                                                                                      WHERE manager pid = SSLECT of FROM person WHERE name = "Name Name")
SACUP BY (invollegorises)
```

Unlocking Value from Your Data Relationships



- Model your data **naturally** as a graph of data and relationships
- Drive graph model from domain and use-cases
- Use relationship information in real-time to transform your business
- Add new relationships on the fly to adapt to your changing requirements

High Query Performance with a Native Graph DB



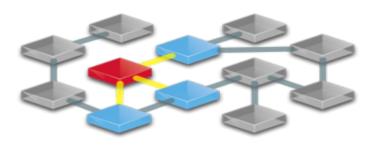
- Relationships are first class citizen
- No need for joins, just follow pre-materialized relationships of nodes
- Query & Data-locality navigate out from your starting points
- Only load what's needed
- Aggregate and project results as you go
- Optimized disk and memory model for graphs



The performance advantage of materialised relationships



- a sample social graph with ~1,000 persons
- average 50 friends per person
- pathExists(a,b) limited to depth 4
- caches warmed up to eliminate disk I/O



	# persons	query time
Relational database	1,000	2000ms
Neo4j	1,000	2ms
Neo4j	1,000,000	2ms

Getting started with Neo4j



- 1. Download Neo4j: http://neo4j.com/download/
- 2. Start the server.
- 3. It should be running on: http://localhost:7474
- 4. Log-in with default credentials

user: neo4j

password: neo4j

5. Choose a **new** password

We're good to go!

