

DATAVERSITY WEBINAR:
**GRAPH DATA
MODELING IN
FOUR DIMENSIONS**

JANUARY 7TH, 2020

INSTRUCTED BY THOMAS FRISENDAL

Photo by Ricardo Gomez Angel on [Unsplash](#)

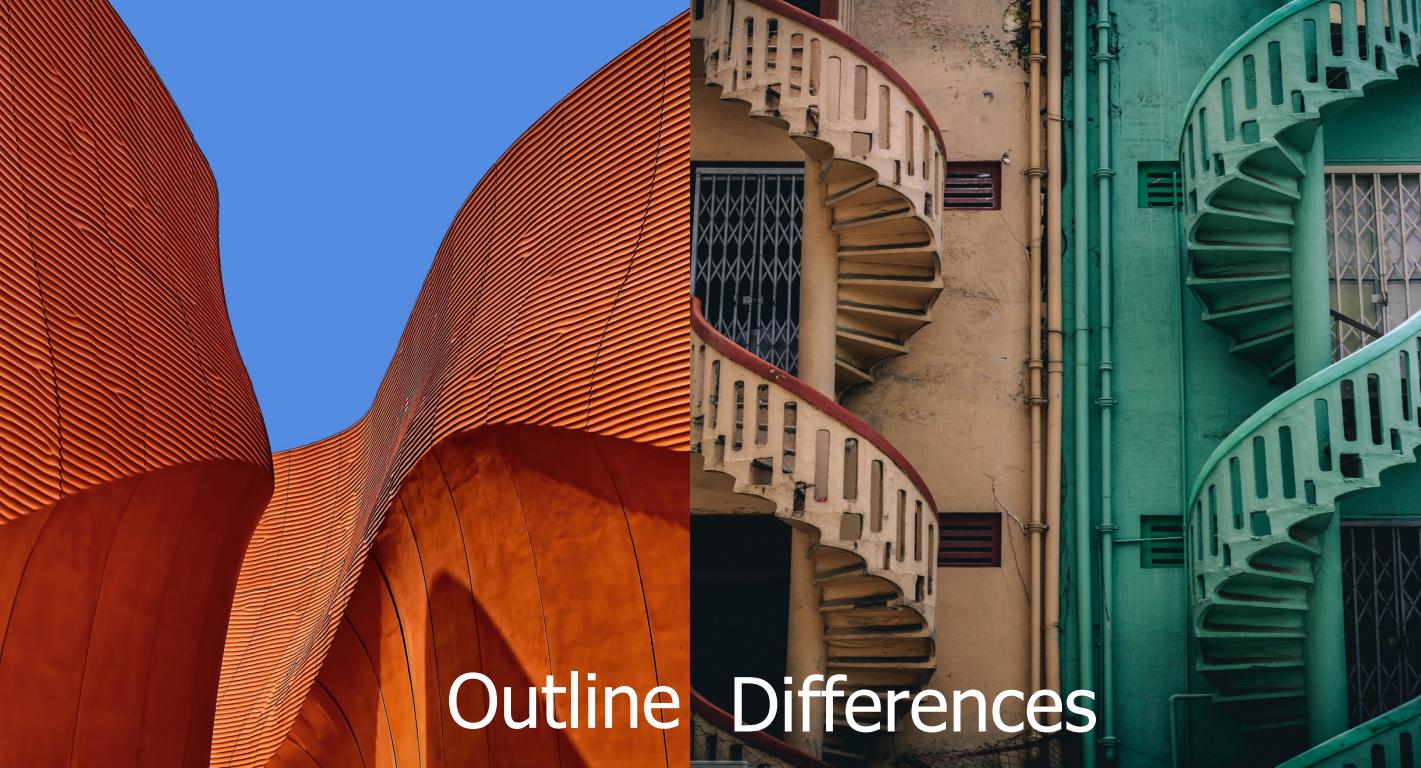


Photo by Greg Jeanneau on [Unsplash](#)



Photo by Dennis Cortés on [Unsplash](#)



Photo by david latorre romero on [Unsplash](#)

GRAPH DATA MODELING OUTLINE

- What?
- Why?

TWO MAJOR CATEGORIES: PROPERTY GRAPH VS. RDF

Property Graph – Heavily Connected Data

The screenshot shows the ICIJ website with a large banner for 'Three Years On The Panama Papers'. Below the banner, a sub-section highlights the impact of the investigation, stating 'Panama Papers Helps Recover More Than \$1.2 Billion Around The World'. The URL in the address bar is <https://www.icij.org/investigations/panama-papers/panama-papers-helps-recover-more-than-1-2-billion-around-the-world>.

RDF: Linked Semantics

The screenshot shows the EDM Council website with a banner for 'THE OPEN SEMANTIC STANDARD FOR THE FINANCIAL INDUSTRY'. It features a logo for 'FIBO' (Financial Industry Business Ontology) and a section titled 'What is the Financial Industry Business Ontology (FIBO™)?'. The text explains that FIBO is a conceptual model developed by members of the financial industry to harmonize data across repositories. The URL in the address bar is <https://edmcouncil.org/page/aboutfiboreview>.

Precise meaning translates into a common language between systems and sources, reduces the cost of doing business and promotes confidence in data among business users. FIBO is the standard for harmonization of data across repositories. It is a mechanism for validating data quality. It is the building block for business process automation and the pathway for flexible risk analysis.

DB DB-Engines Ranking - popularit X +

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

DB-ENGINES

Knowledge Base of Relational and NoSQL Database Management Systems

Home | DB-Engines Ranking | Systems | Encyclopedia | Blog | Search | Vendor Lo

Featured Products: [Couchbase](#) [Neo4j](#) [Redis](#) [AllegroGraph](#) [DataStax](#)

Select a ranking

- Complete ranking
- Relational DBMS
- Key-value stores
- Document stores
- Graph DBMS
- Time Series DBMS
- Object oriented DBMS
- RDF stores
- Search engines
- Wide column stores
- Multivalue DBMS
- Native XML DBMS
- Event Stores
- Content stores
- Navigational DBMS

Special reports

- Ranking by database model
- Open source vs. commercial

Ranking > Graph DBMS

DB-Engines Ranking of Graph DBMS

The DB-Engines Ranking ranks database management systems according to their popularity. The ranking is updated monthly.

This is a partial list of the [complete ranking](#) showing only graph DBMS.

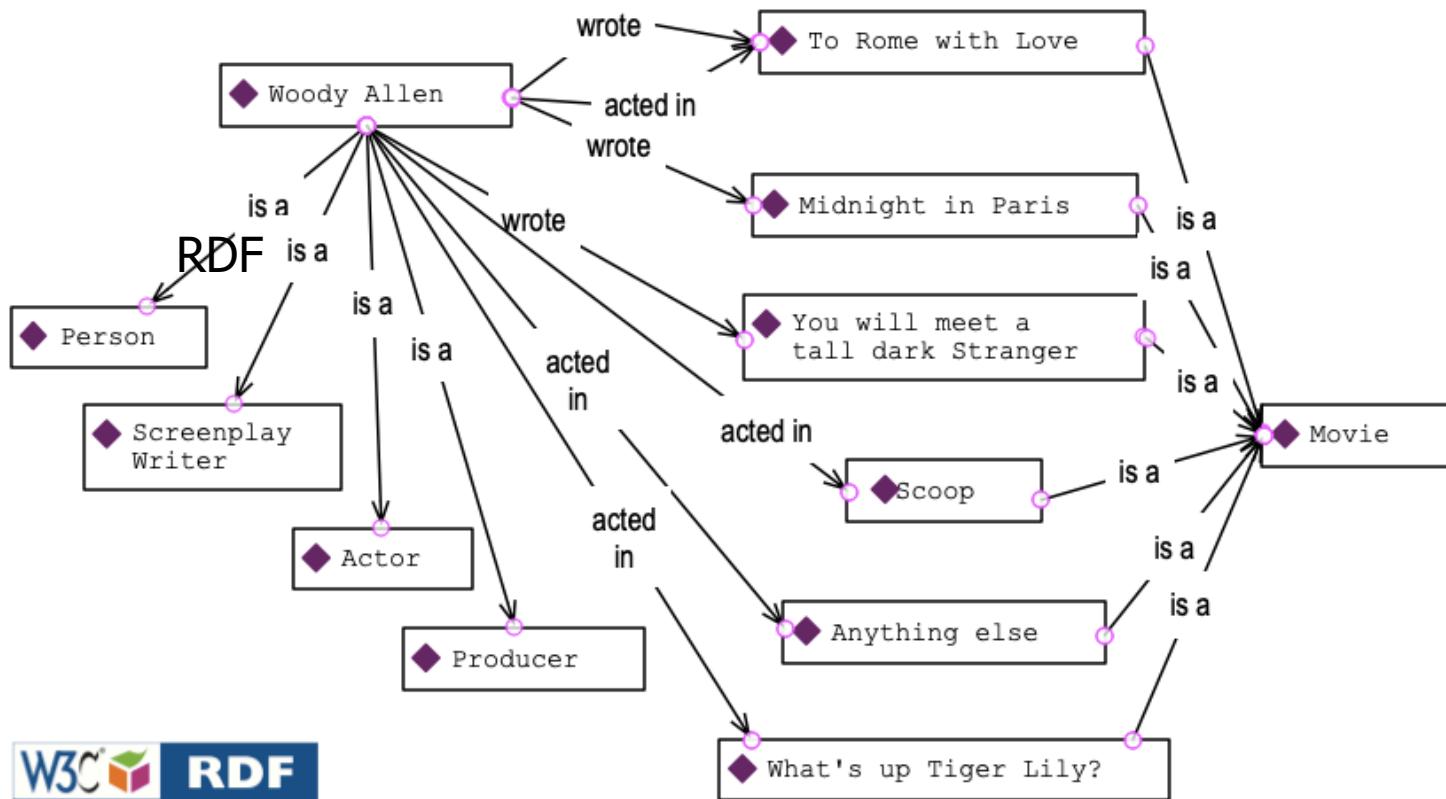
Read more about the [method](#) of calculating the scores.

include secondary database models 32 systems in ranking, June 2019

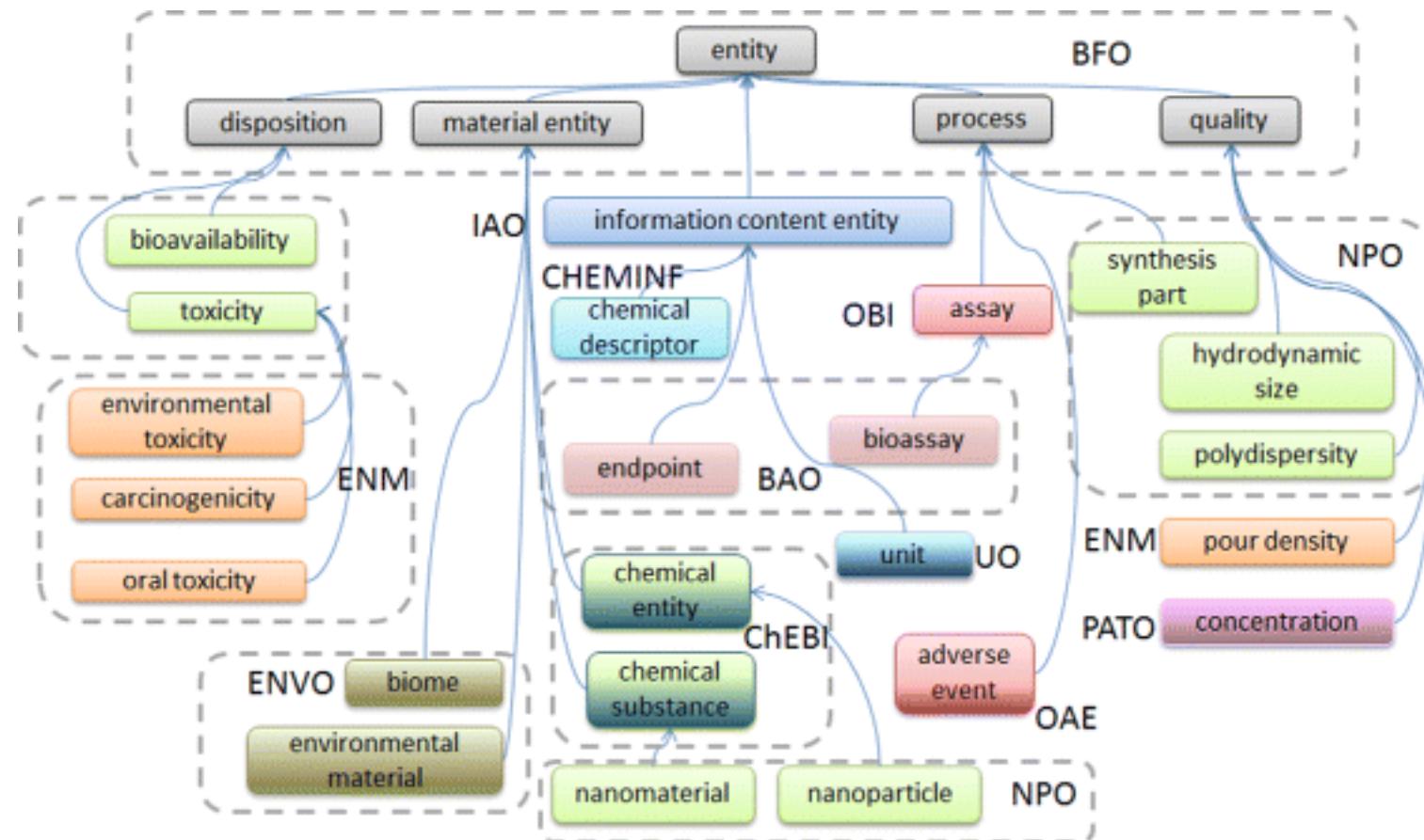
Rank	DBMS			Database Model	Score		
	Jun 2019	May 2019	Jun 2018		Jun 2019	May 2019	Jun 2018

Microsoft Azure Cosmos DB	Multi-model	Document store, Graph DBMS, Key-value store, Wide column store	05
OrientDB	Multi-model		25
ArangoDB	Multi-model		05
Virtuoso	Multi-model	Wide column store	33
JanusGraph	Graph	1.55 -0.07 +1.19	

SEMANTIC WEB – RDF: RESOURCE DESCRIPTION FRAMEWORK – OWL: WEB ONTOLOGY LANGUAGE - W3C STANDARDS



SEMANTIC WEB – RDF: RESOURCE DESCRIPTION FRAMEWORK – OWL: WEB ONTOLOGY LANGUAGE - W3C STANDARDS



PROPERTY GRAPHS

Apache (2015) Tinkerpop / Gremlin (2009):

Amazon Neptune,
Azure CosmosDB,
Datastax Graph,
JanusGraph,
Cypher for Gremlin,
OrientDB,
Stardog

Neo4J Cypher (2011):

openCypher (2017) – see list to the right:

SQL Style:

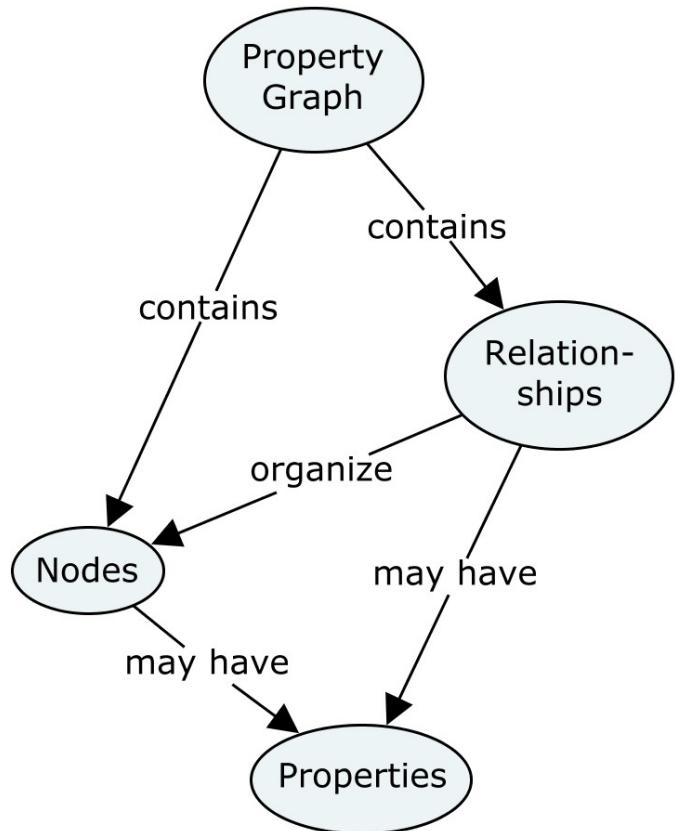
ISO SQL Property Graph Extensions (2017-2020?),
Oracle PGQL (2016),
SQL Server 2017 Graph,
TigerGraph GSQL (2018)

*(*Own not authoritative research mid 2019, errors and unintended omissions may exist. If so, I apologize)*

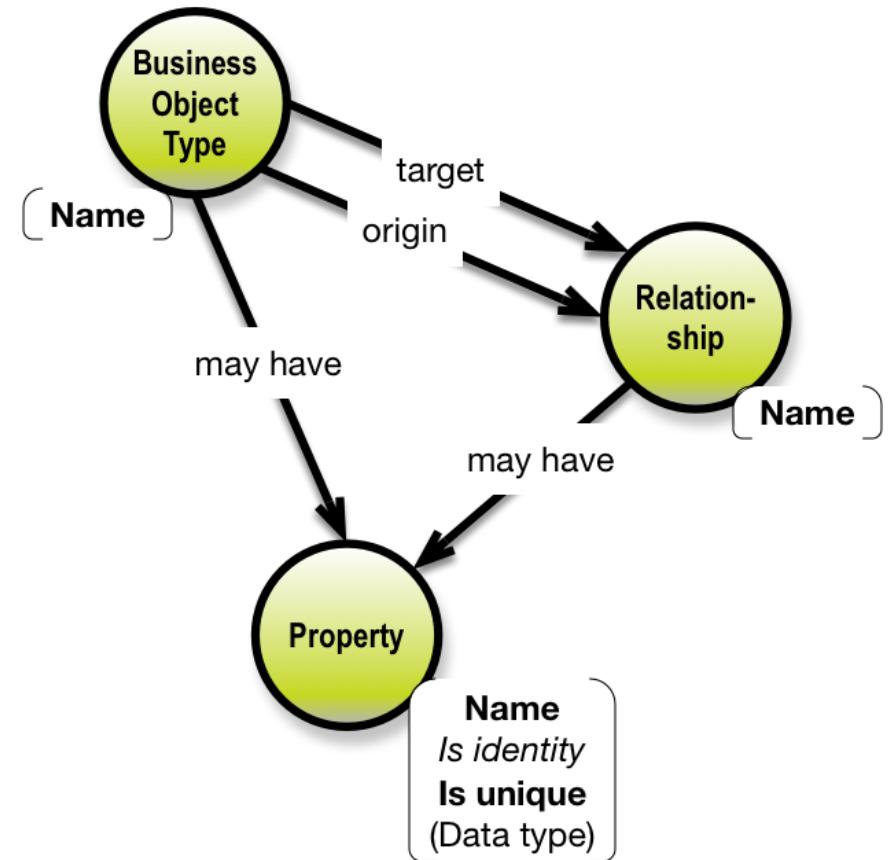
The screenshot shows a screenshot of a web browser displaying the openCypher project page at <https://www.opencypher.org/projects>. The page has a header with links for Home, About, Blog, Events, Cypher in use, Features, and Contributions. Below the header, a message states: "The implementations are ordered alphabetically within each category." There are two main sections: "Cypher implementations (vendor/industrial)" and "Cypher implementations (research)". The "Cypher implementations (vendor/industrial)" section lists: Agens Graph: A multi-model database; AnzoGraph: A native massively parallel (MPP) graph analytical database; CAPS: Cypher for Apache Spark; Cypher for Gremlin; Memgraph: An in-memory, transactional graph database; Neo4j: A native, transactional property graph database; RedisGraph: A graph module for Redis; and SAP HANA Graph. The "Cypher implementations (research)" section lists: Gradoop: Distributed graph analytics on Hadoop; Graphflow: An active graph database; ingraph: Graph query engine; and Ruruki: An in-memory graph database.

THE BASICS OF PROPERTY GRAPHS

Property Graph Concepts

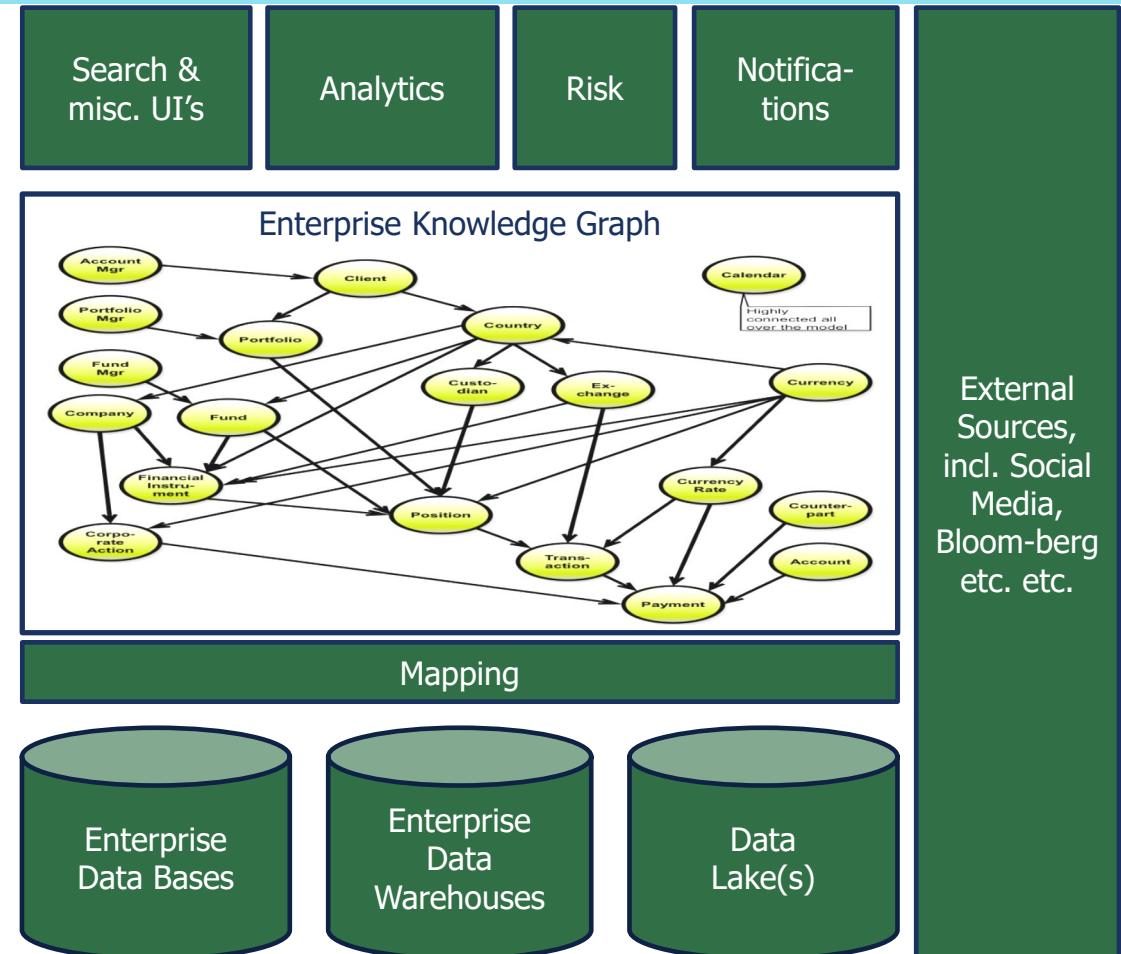


Property Graph Meta Model



WHAT ABOUT KNOWLEDGE GRAPHS?

- Example to the right is for the financial sector
- Navigation across silos
- Integration – even streaming
- Automation, e.g. of alerts
- *Pick your choice:*
 - *Strong semantic focus: RDF*
 - *Strong network focus: Property Graphs*
 - *Hybrid solutions*



SQL?

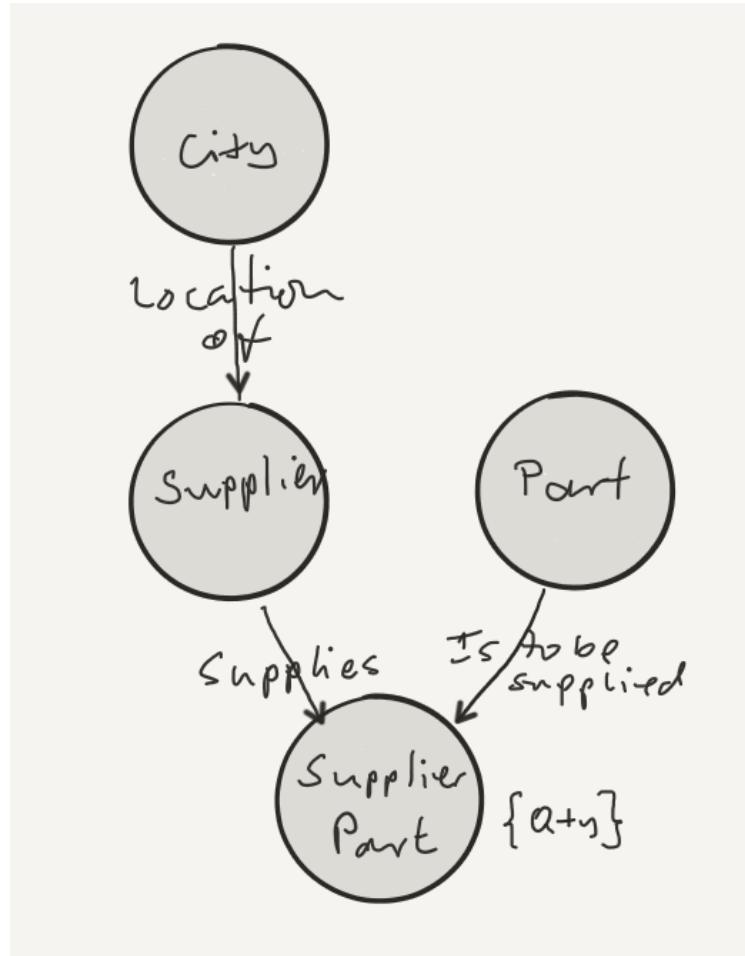
W.I.P.: PROPERTY GRAPH EXTENSIONS

- In 2017, INCITS in the US invited contributors to “Standardizing Graph Database Functionality”
- Ongoing work in the ISO/IEC SC32/WG3 (the information below is not authoritative)
- Graph Query Language
- Embed graph queries in the FROM clause of SELECT statements
- Visualization of graph paths as part of the syntax
- Describe property graphs (schema syntax)
- Map SQL tables to property graph
- Snapshot from 2018 of work in progress
 - <https://www.youtube.com/watch?v=hIhbUeoCPio>

“EVERYTHING LOOKS LIKE A GRAPH”

GRAPHS ARE INTUITIVE, COMMUNICATE WELL, AND ARE FUN TO DRAW – TOGETHER!

COGNITION AND PERCEPTION: GRAPHS JUST COME, NATURALLY, TO US ON THE WHITEBOARD



SPATIAL THINKING IS THE FOUNDATION OF ABSTRACT THOUGHT

Barbara Tversky, Prof. Emerita of Psychology, Stanford:

*Book: "Mind in Motion: How Action Shapes Thought",

Hachette UK, May 2019

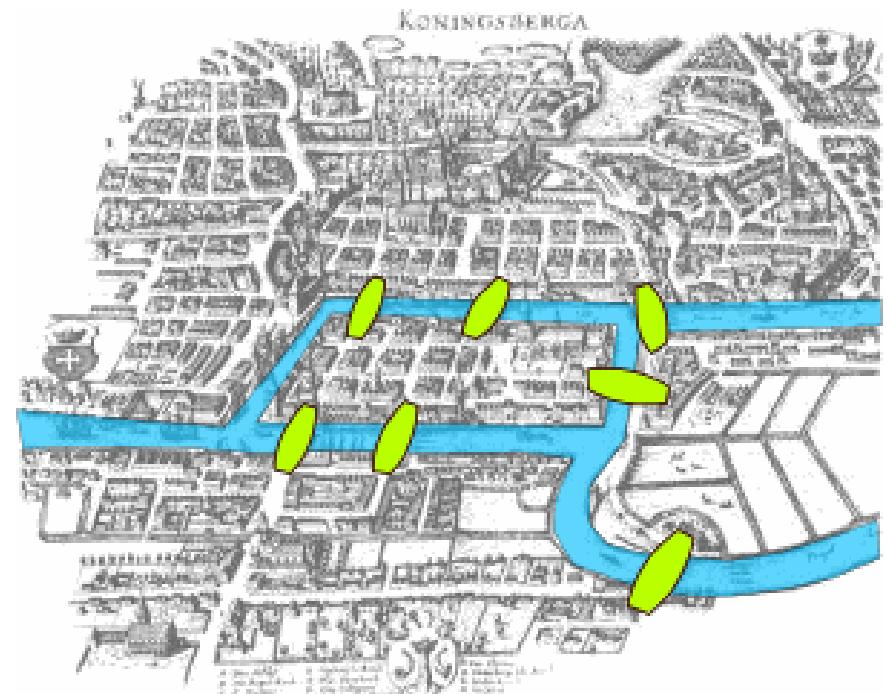
"Maps in Minds", "Space: Maps", "The World Is Diagrammed" ...

The Nine Laws of Cognition

1. There are no benefits without costs
2. Action molds perception
3. Feeling comes first
4. The mind can override perception
5. Cognition mirrors perception
6. Spatial thinking is the foundation of abstract thought
7. The mind fills in missing information
8. When thought overflows the mind, the mind puts it into the world
9. We organize the stuff in the world the way we organize the stuff in the mind

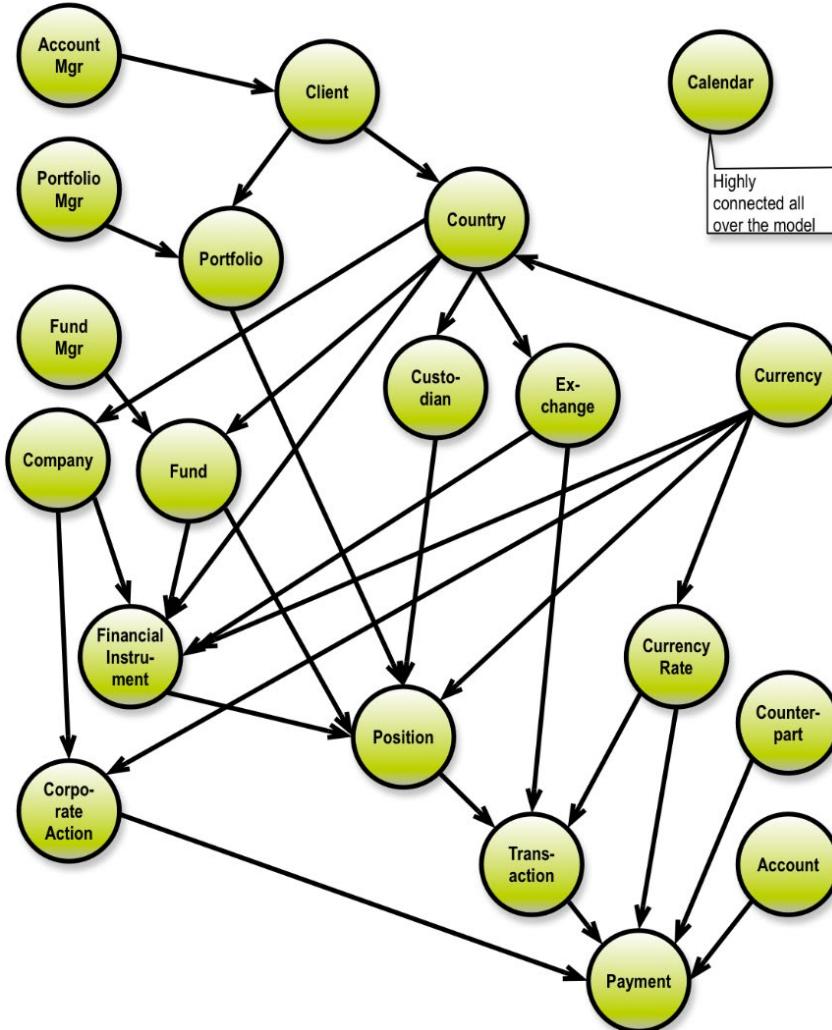
Leonhard Euler, 1736:

The Seven Bridges of Königsberg

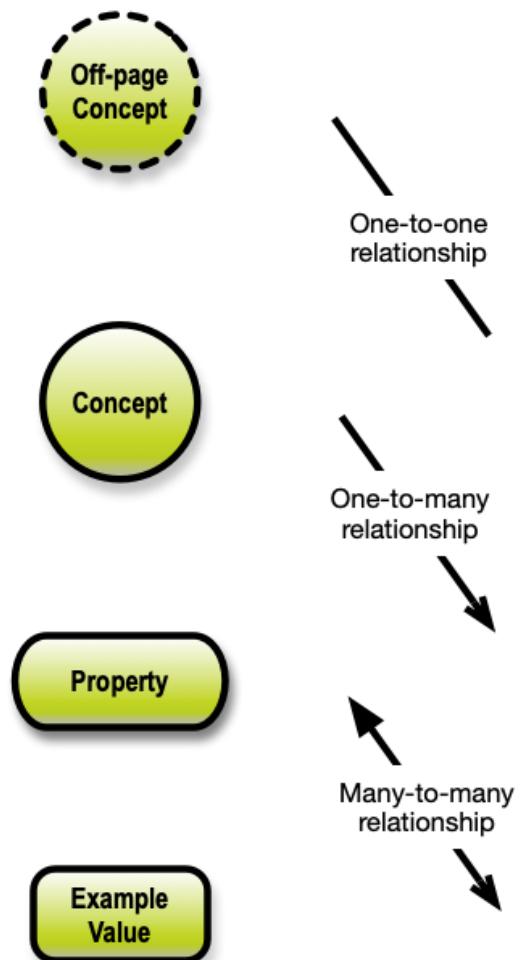


The birth of graph theory!

MAKE COMPLEX FINANCIALS SIMPLE

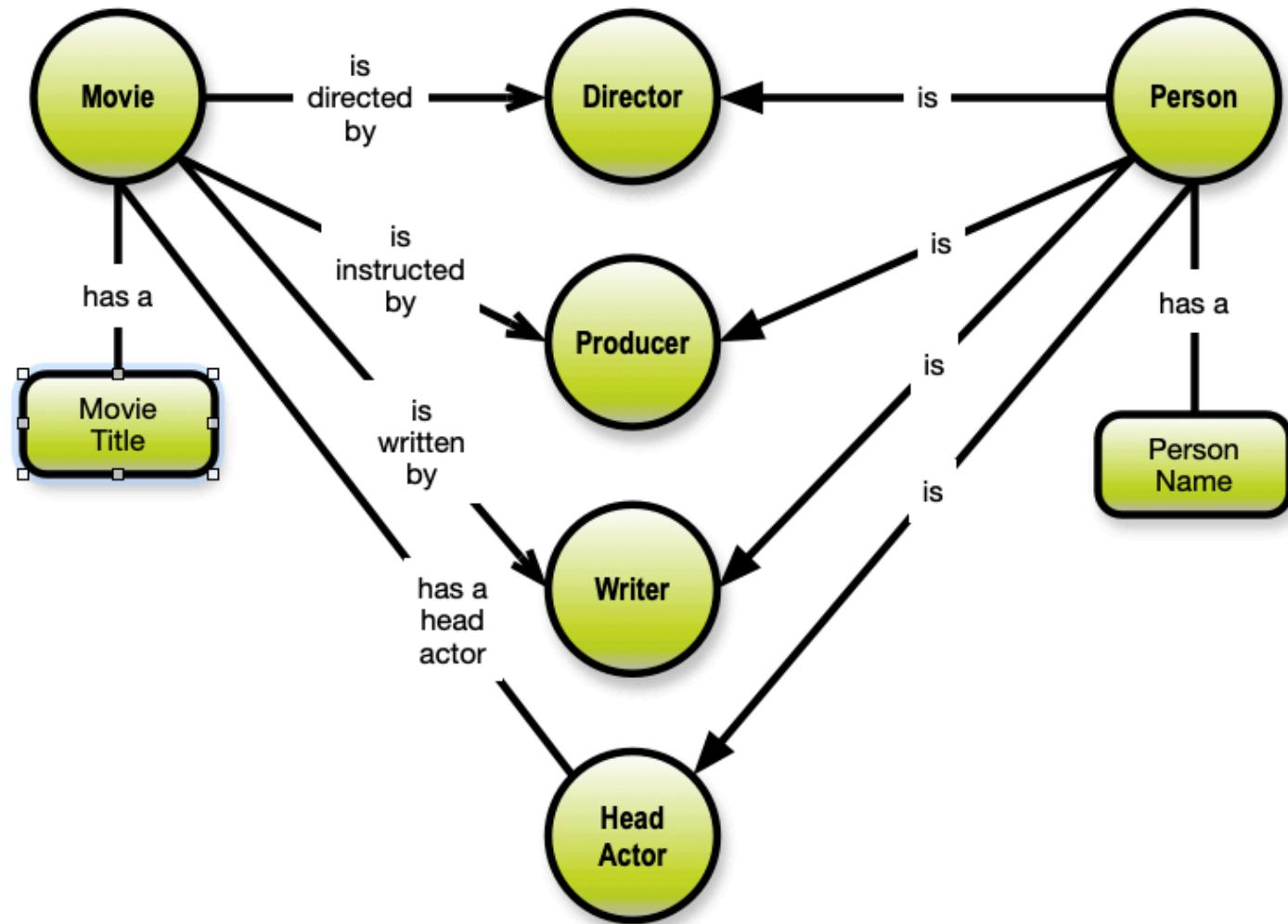


THE VISUAL SYNTAX (MY WAY)



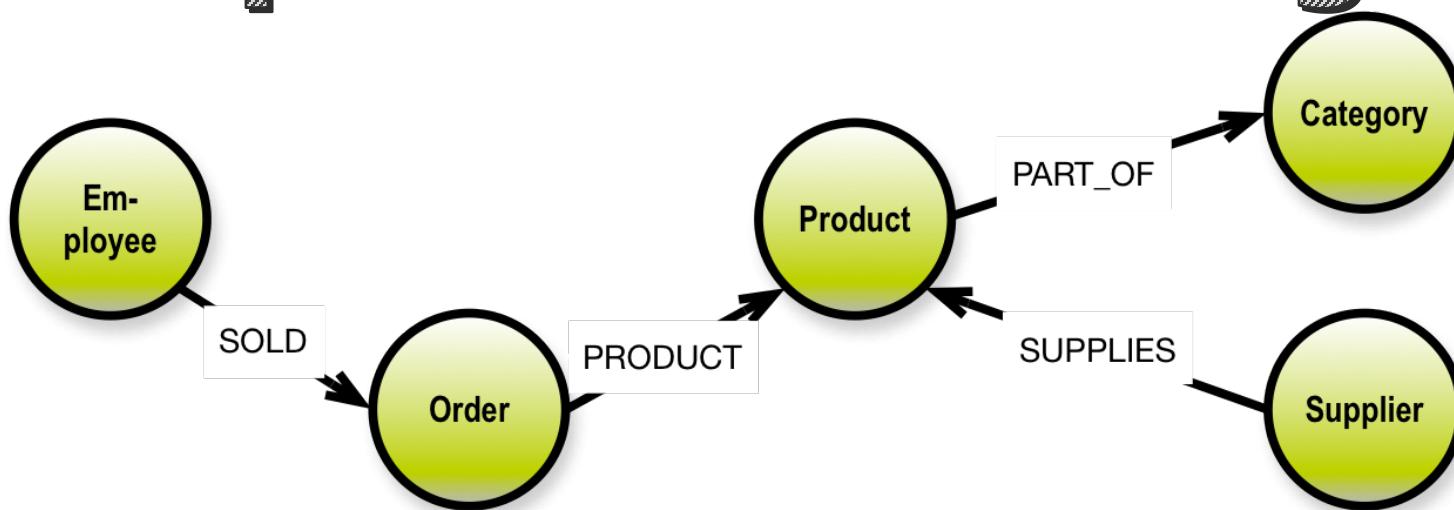
IMDB-INSPIRED DATA MODEL FOR MOVIES

- Read the Concept Model as little sentences!
- In CmapTools the texts on the arrow are called “Linking Phrases”
- In Relational Modeling they are called “Dependencies,” but they did not have names (i.e., no semantics!)
- In Graph Data Models they are called “Relationships” (or “Edges” in the math version)

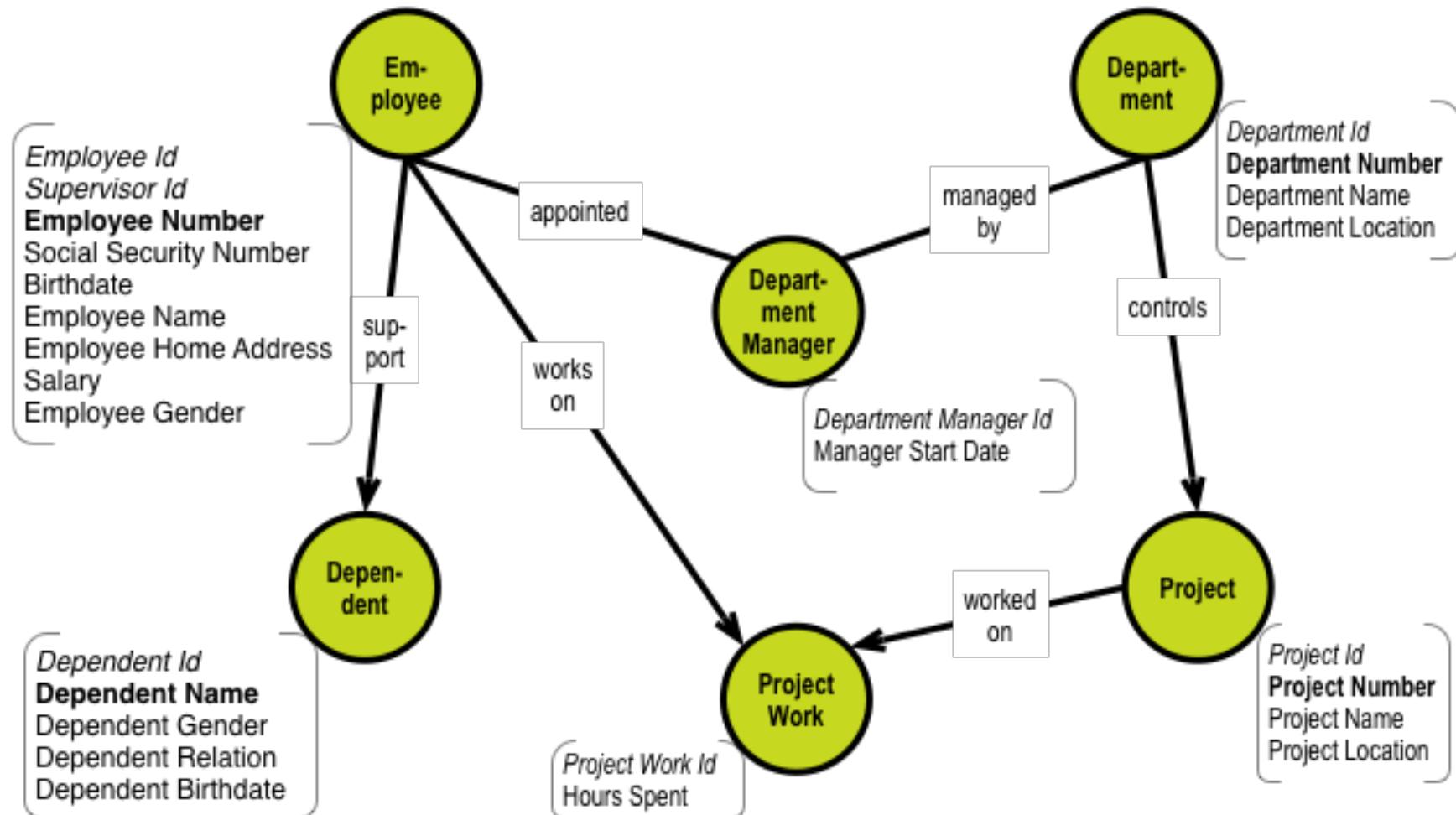


PRIMARY TAKEAWAY:

**Named Relationships is what
matters most in
Graph Data Modeling!!!**



STRUCTURE AND MEANING EXPRESSED AS A PROPERTY GRAPH

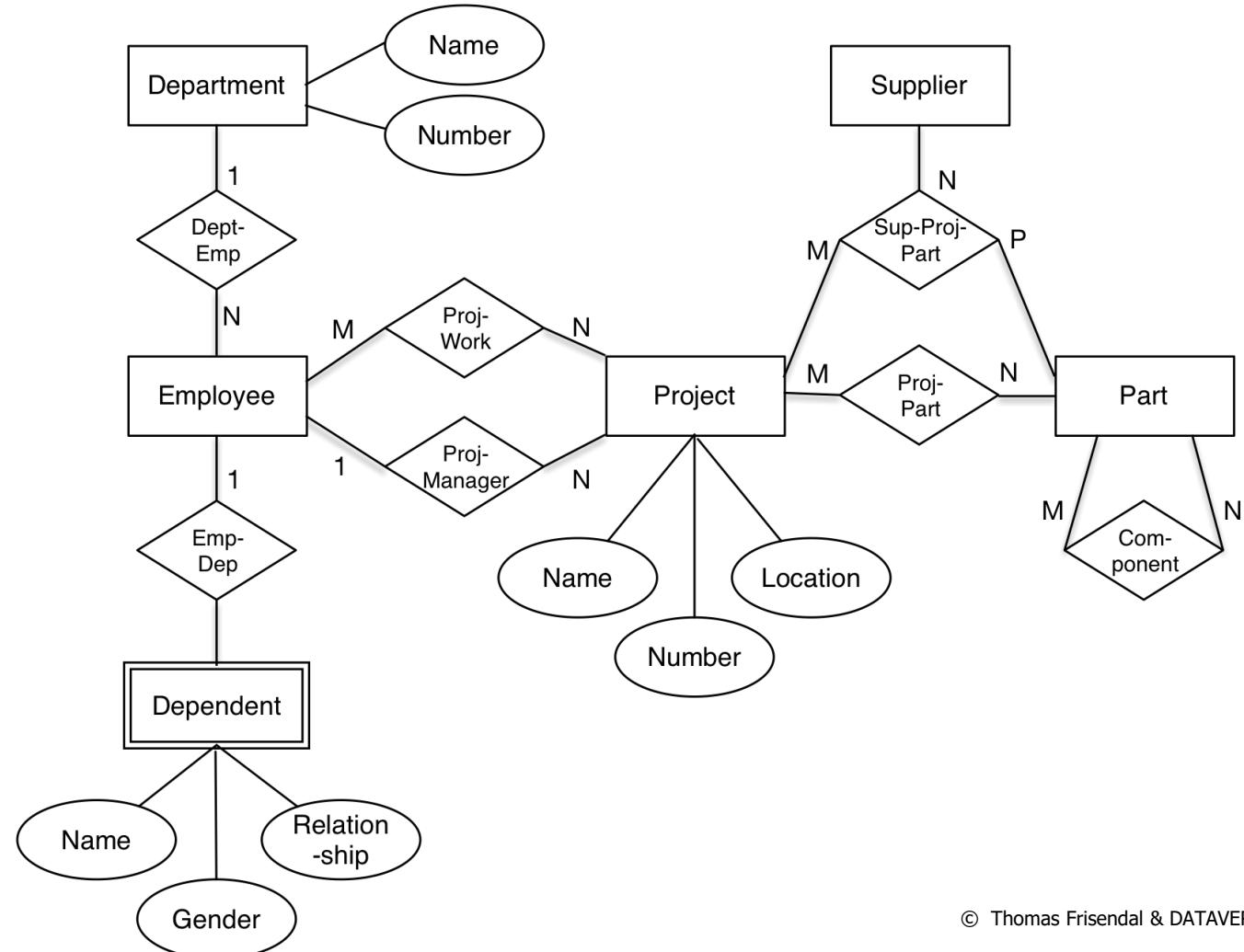


GRAPH DATA MODELING DIFFERENCES

- Graph Data Modeling compared to
- Classic Data Modeling

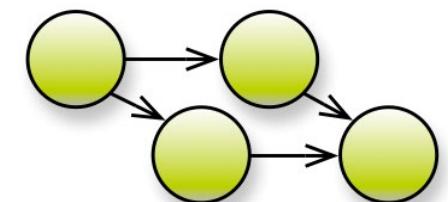
PETER CHEN – ENTITIES, ATTRIBUTES, RELATIONSHIPS

1976

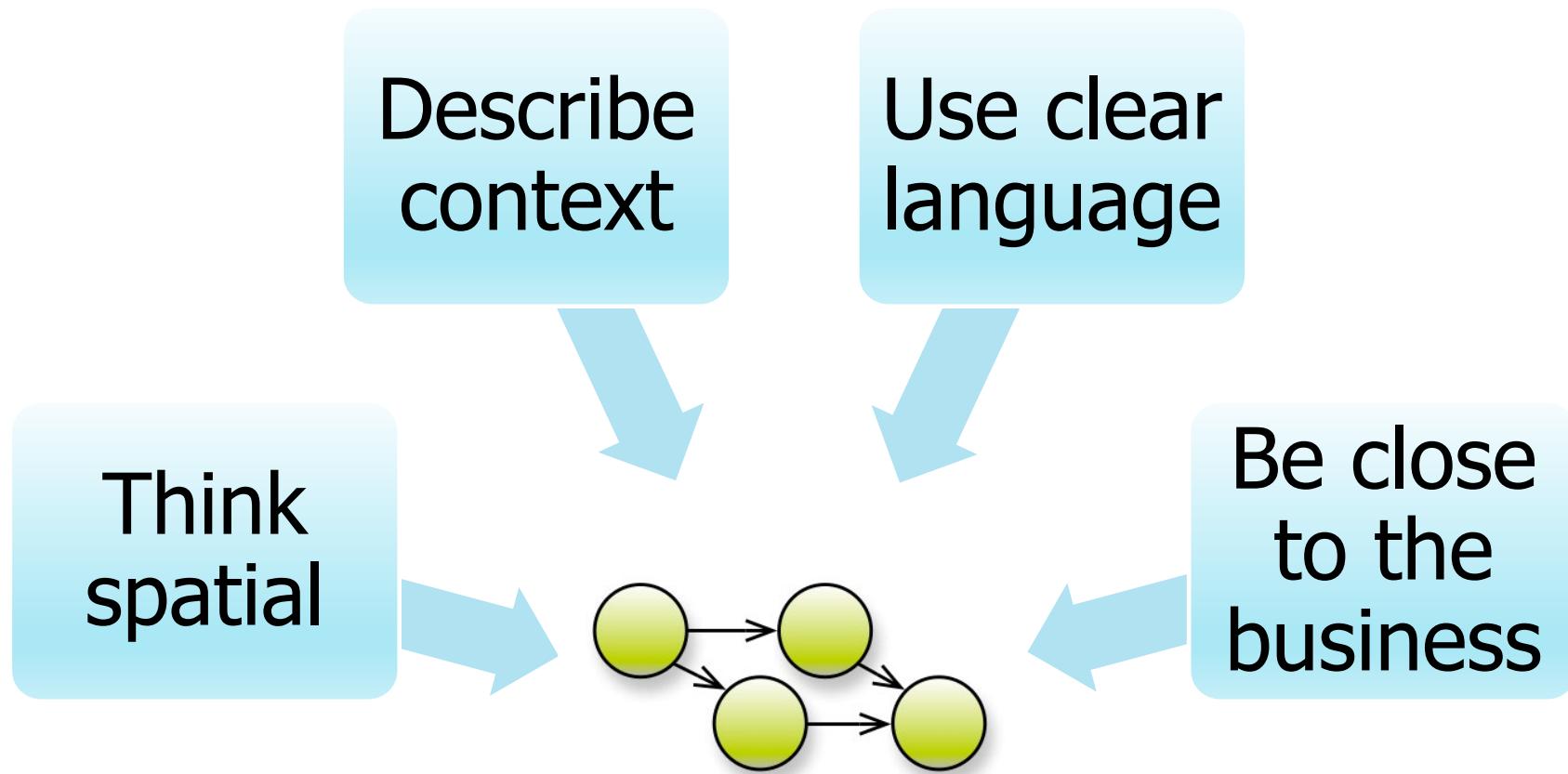


ISSUES IN MODELING DATA AS PROPERTY GRAPHS COMPARED TO CLASSIC DATA MODELING

- Visualization
- Highly Connected Data
- Normalization
- Many to Many
- Think in Paths
- Labels (and Types)
- The Role of the Schema (if any)
 - Schema First
 - Schema Last
- New Opportunities for the Data Modeler

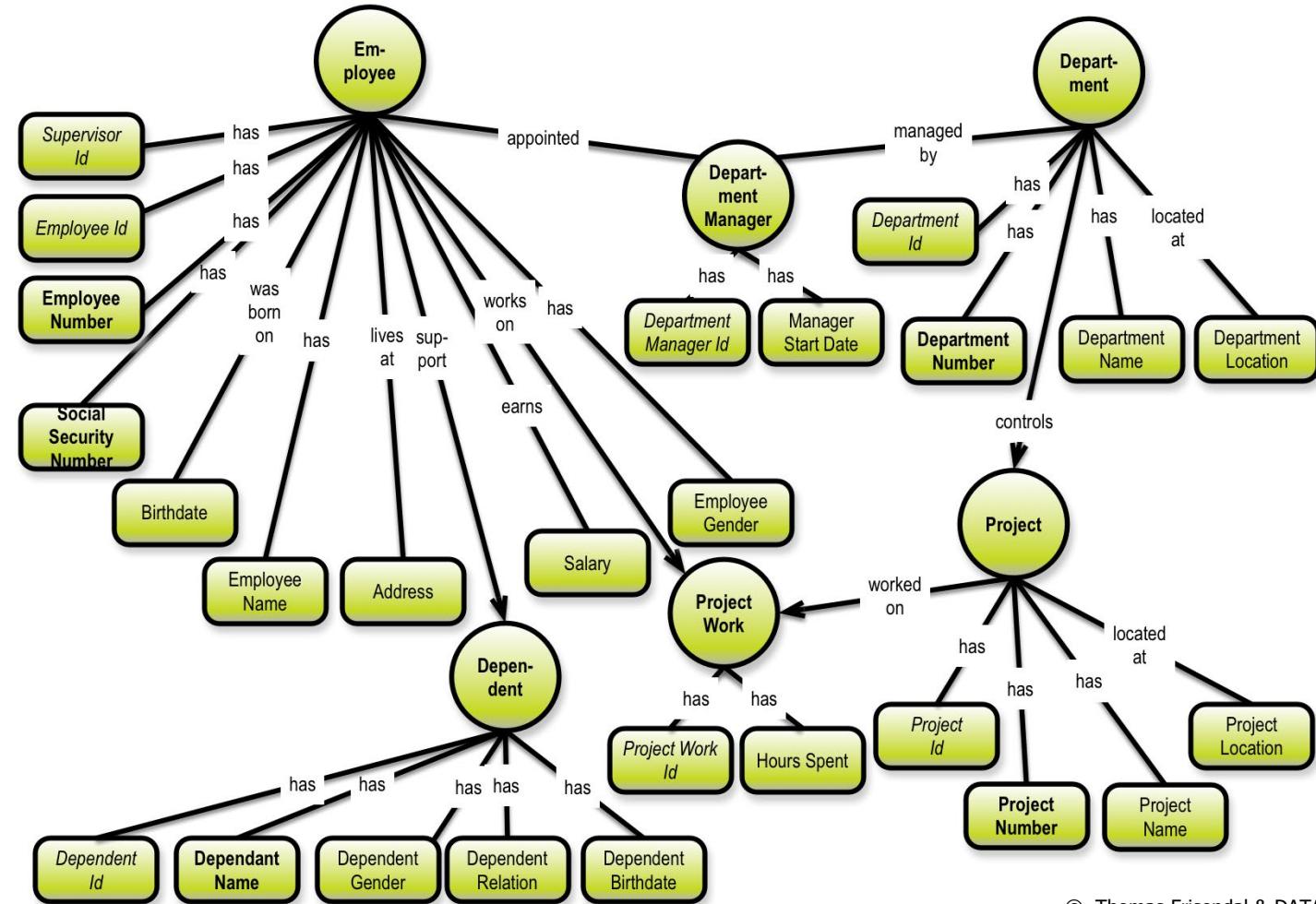


THE BASICS OF GRAPHING DATA

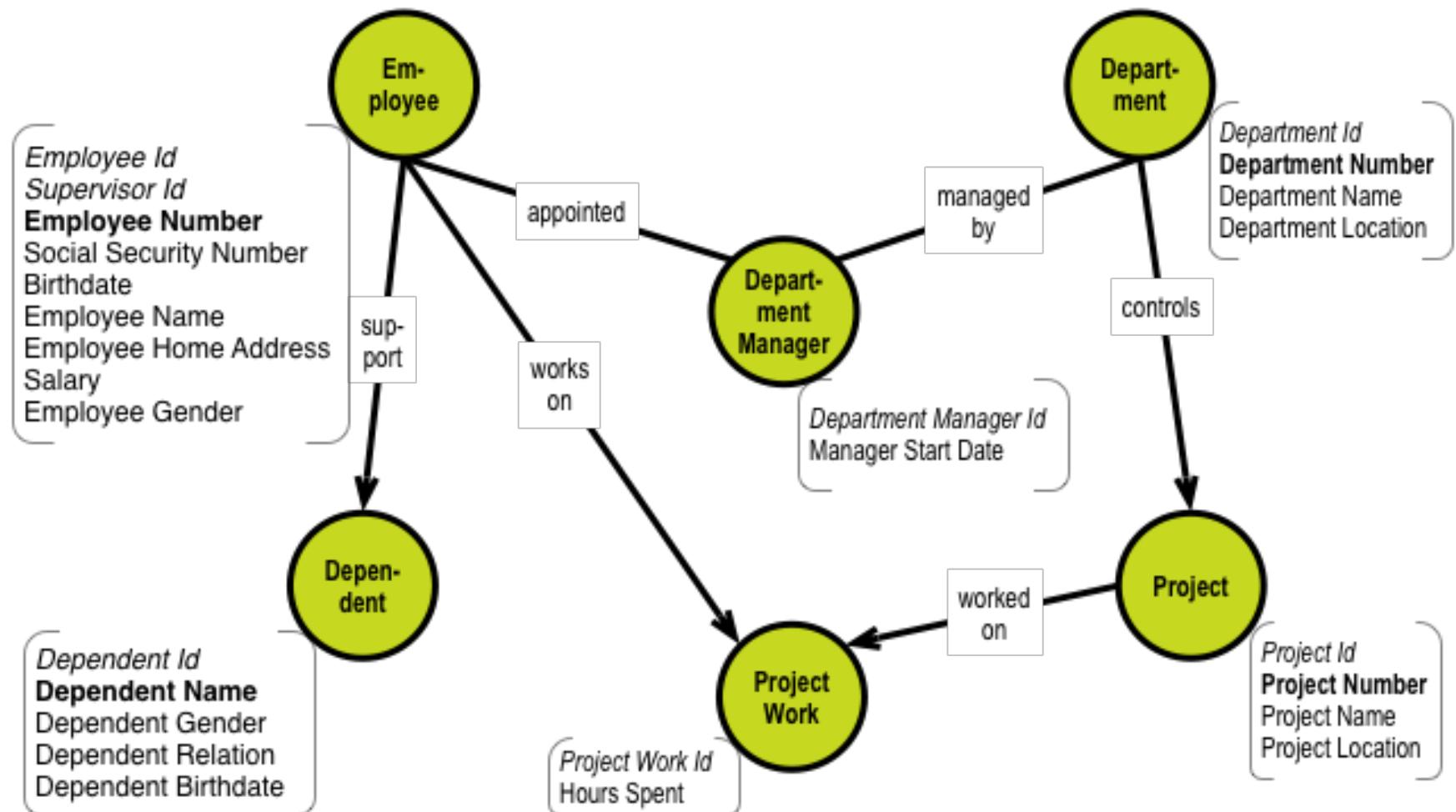


Graph it!

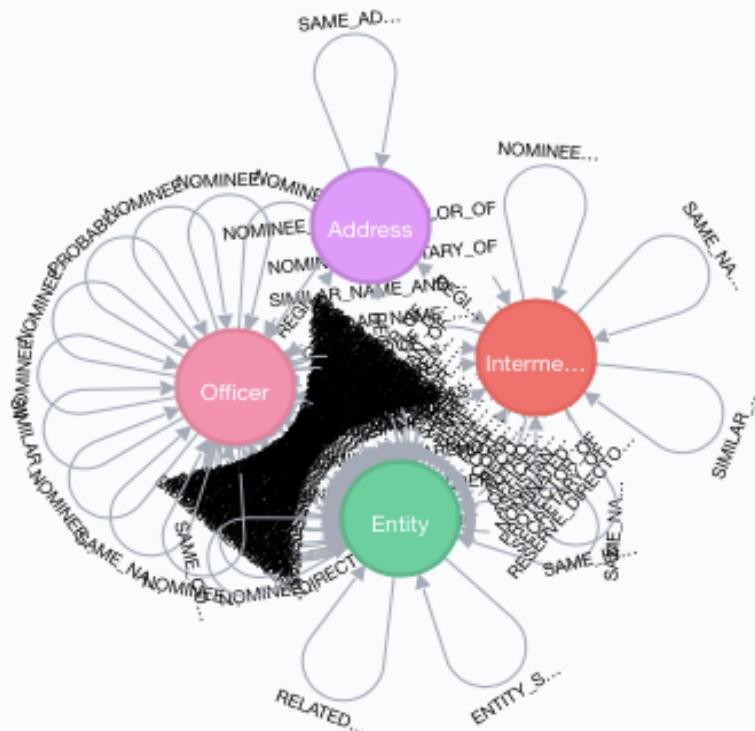
A CONCEPT MODEL (MAP) "A DATA-DRIVEN USER STORY"



PROPERTY GRAPH REPRESENTATION COMMUNICATES STRUCTURE AND MEANING!



HIGHLY CONNECTED DATA: THE PANAMA PAPERS



THE RELATIONAL MODEL

MAKING SENSE BY EXPLORING RELATIONSHIPS?

Famous text book example from several of Chris Date's books and presentations about the relational model.

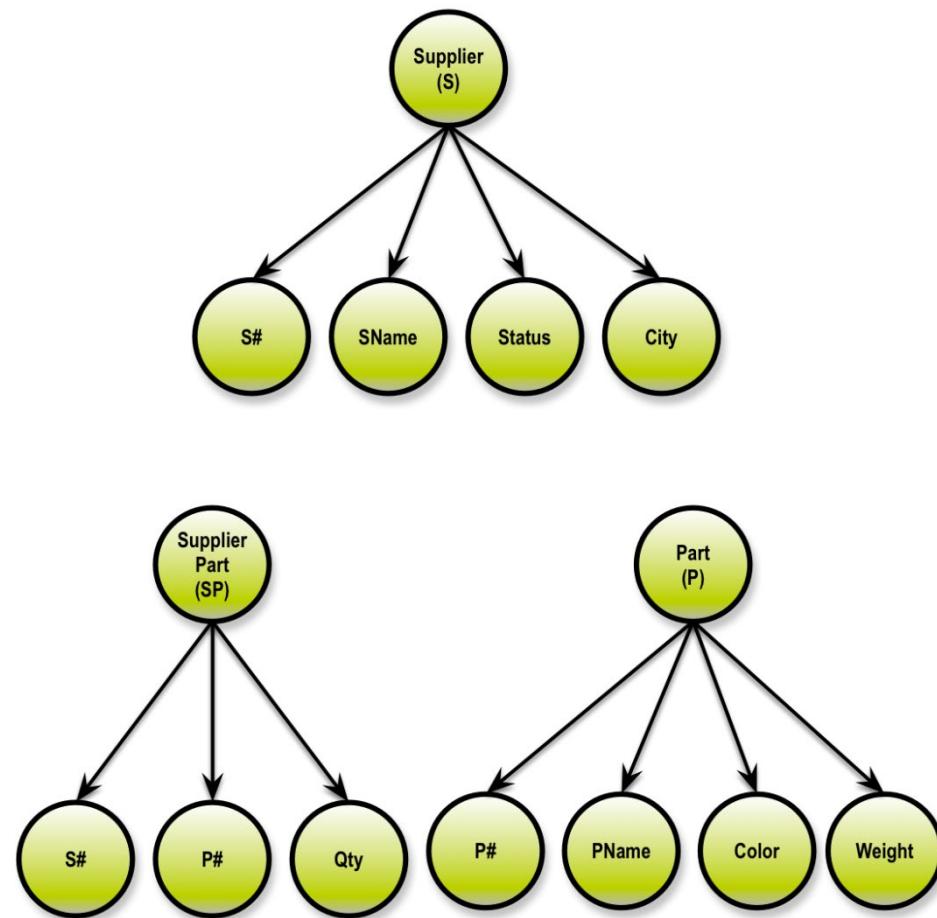
I believe the first appearance of this example was in "An Introduction to Database Systems, Volume 1", C.J. Date, Addison-Wesley (I have the Fourth Edition from 1986).

SNO	SNAME	STATUS	CITY
S1	Smith	20	London
S2	Jones	30	Paris
S3	Blake	30	Paris
S4	Clark	20	London
S5	Adams	30	Athens

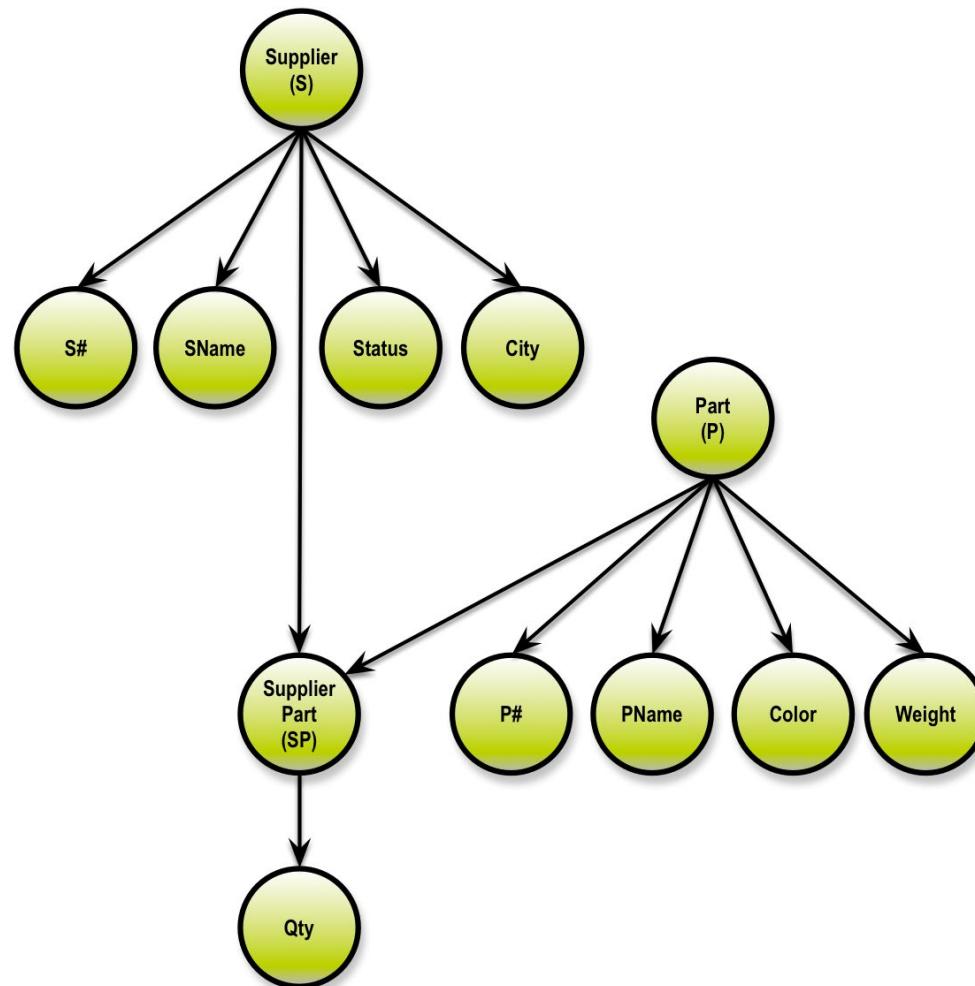
PNO	PNAME	COLOR	WEIGHT	CITY
P1	Nut	Red	12.0	London
P2	Bolt	Green	17.0	Paris
P3	Screw	Blue	17.0	Paris
P4	Screw	Red	14.0	London
P5	Can	Blue	12.0	Paris
P6	Cog	Red	19.0	London

SNO	PNO	QTY
S1	P1	300
S1	P2	200
S1	P3	400
S1	P4	200
S1	P5	100
S1	P6	100
S2	P1	300
S2	P2	400
S3	P2	200
S4	P2	200
S4	P4	300
S4	P5	400

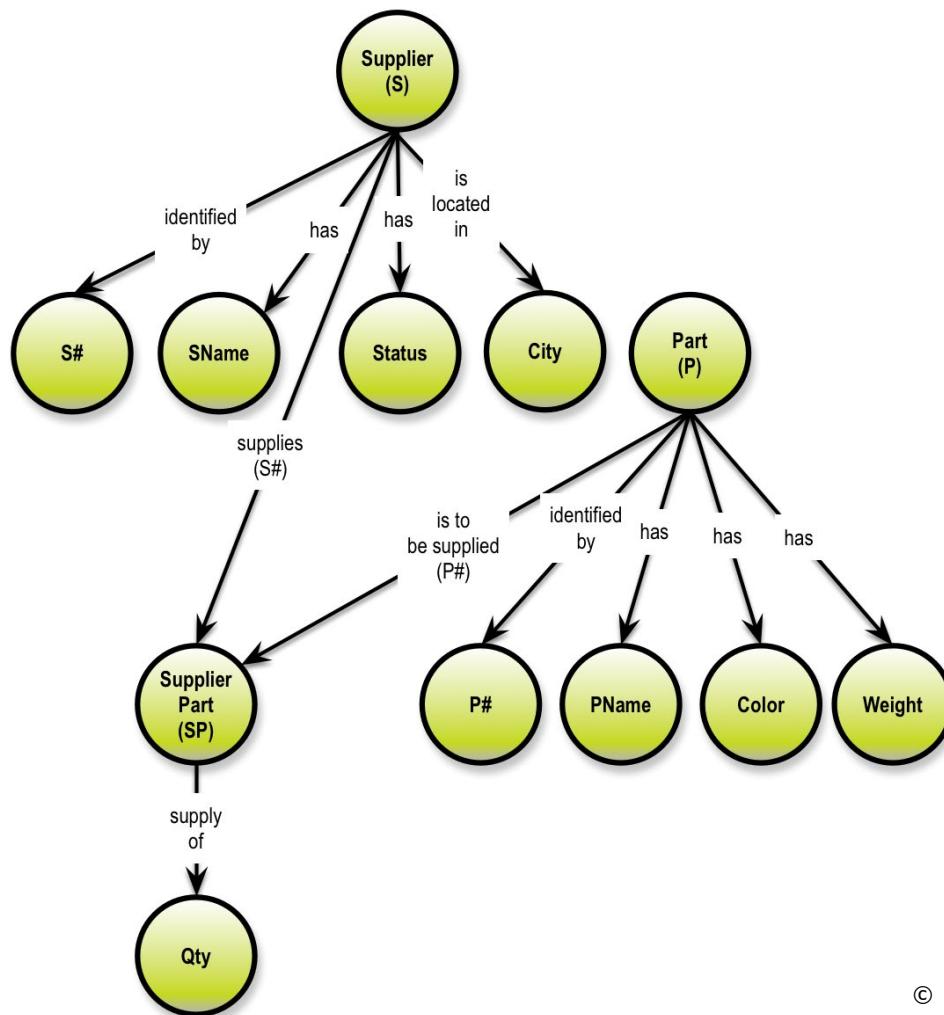
DRAW A NAIVE CONCEPT MAP



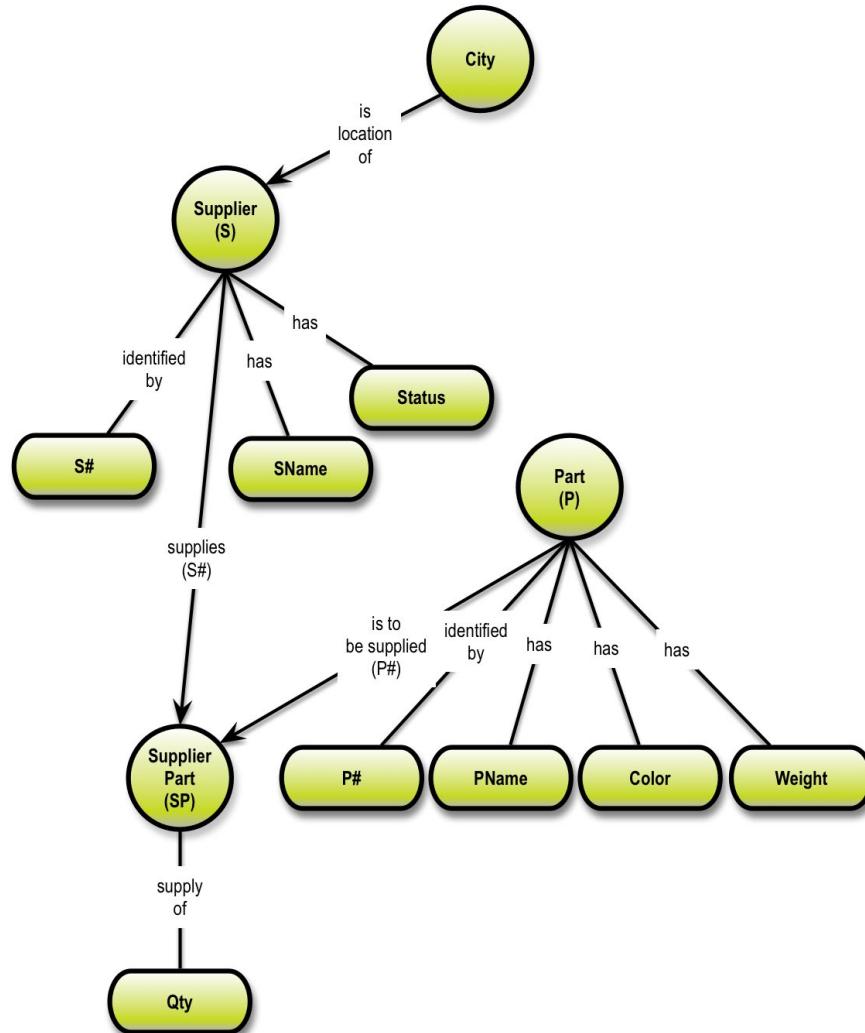
VISUALIZE IMPLICIT RELATIONSHIPS



NAME THE RELATIONSHIPS



VISUALIZE THE PROPERTIES AND FINALIZE THE CONCEPT DESIGN



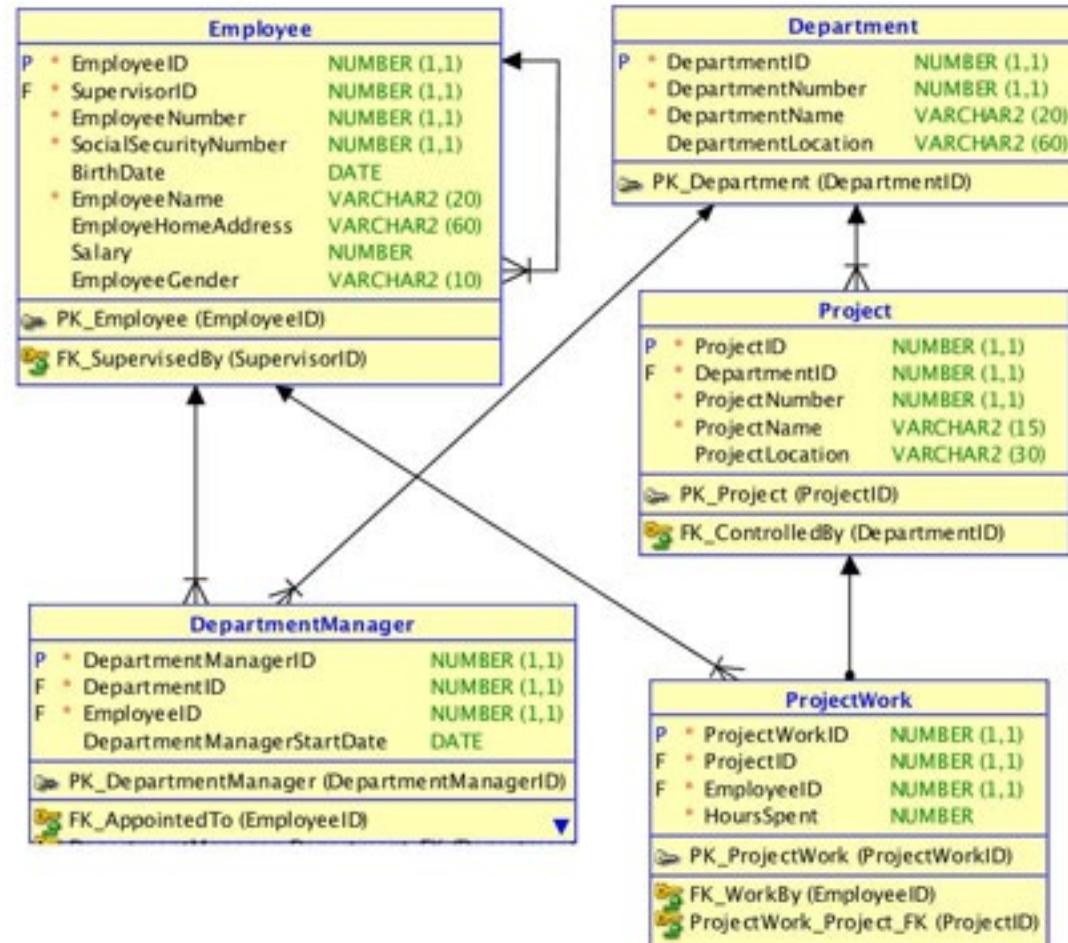
WHICH COMMUNICATES BEST? THIS?

<u>SNO</u>	<u>SNAME</u>	<u>STATUS</u>	<u>CITY</u>
S1	Smith	20	London
S2	Jones	30	Paris
S3	Blake	30	Paris
S4	Clark	20	London
S5	Adams	30	Athens

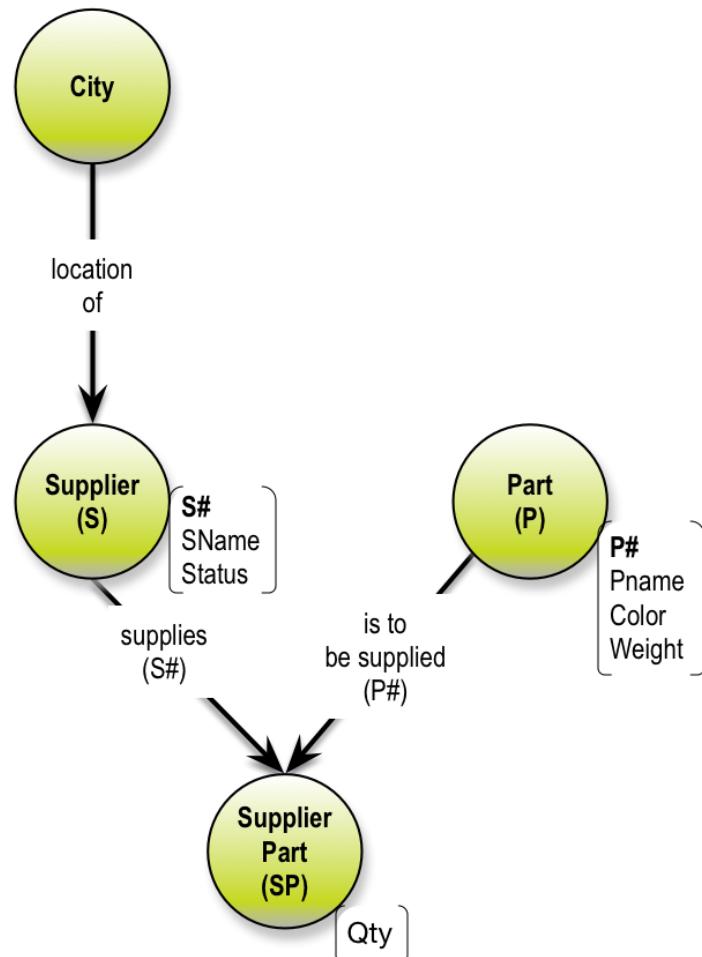
<u>PNO</u>	<u>PNAME</u>	<u>COLOR</u>	<u>WEIGHT</u>	<u>CITY</u>
P1	Nut	Red	12.0	London
P2	Bolt	Green	17.0	Paris
P3	Screw	Blue	17.0	Paris
P4	Screw	Red	14.0	London
P5	Can	Blue	12.0	Paris
P6	Cog	Red	19.0	London

<u>SNO</u>	<u>PNO</u>	<u>QTY</u>
S1	P1	300
S1	P2	200
S1	P3	400
S1	P4	200
S1	P5	100
S1	P6	100
S2	P1	300
S2	P2	400
S3	P2	200
S4	P2	200
S4	P4	300
S4	P5	400

OR THIS?



OR THIS?



THE CONTEMPORARY STYLE OF NORMALIZATION

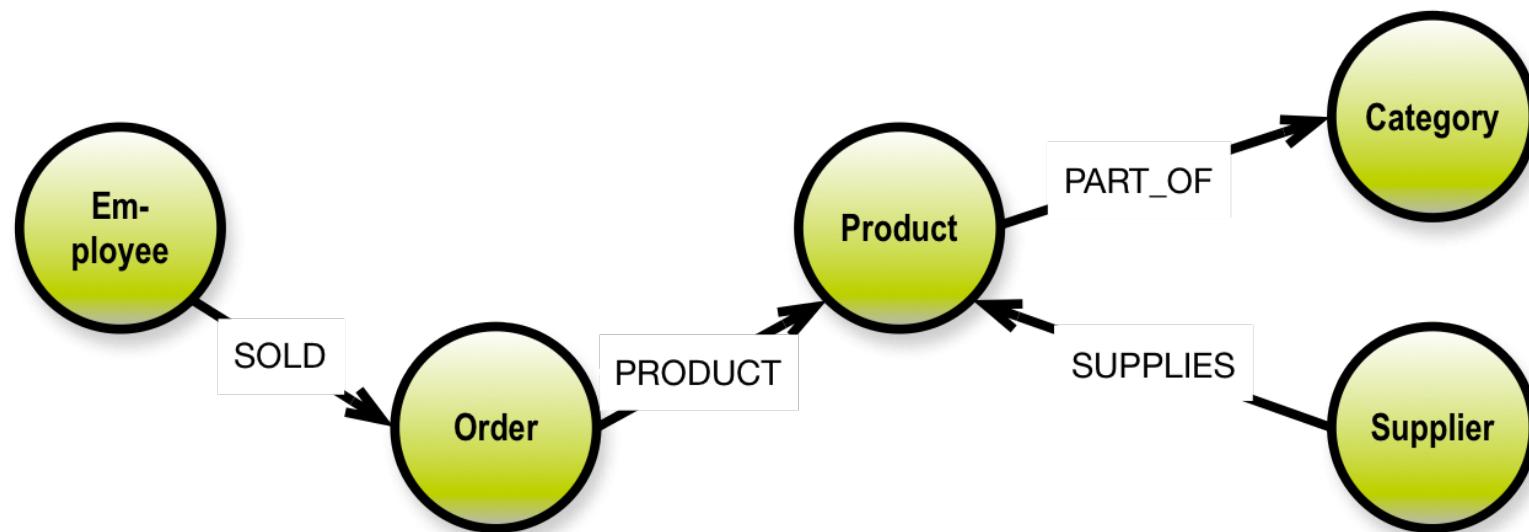
- Visualizing the dependencies as a directed graph gives a deep understanding of the functional dependency structure
- Named relationships between all concepts explain the semantics really well
- Identity and uniqueness contribute hugely to setting the context in precise ways

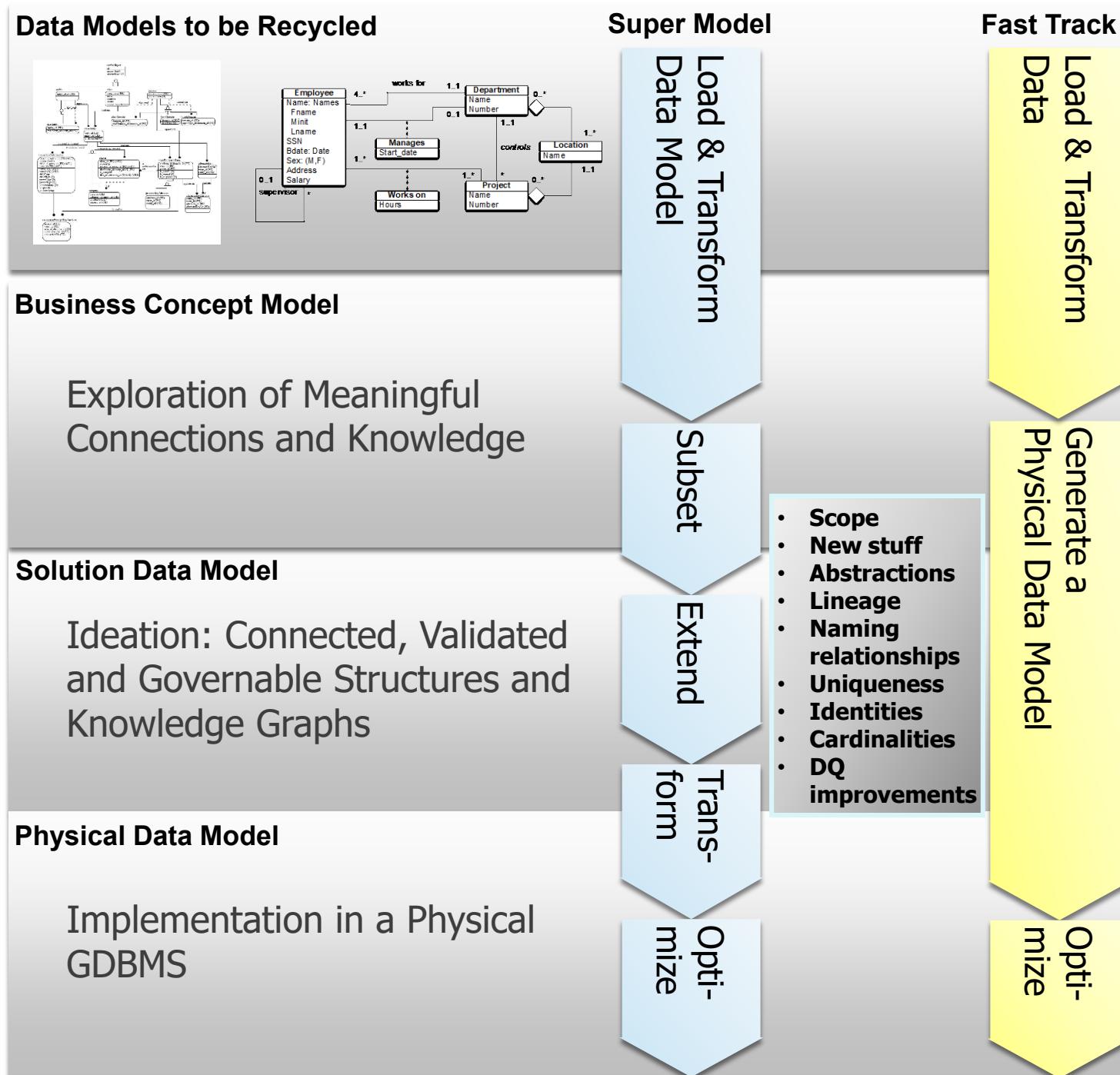
GRAPH DATA MODELING ARTISANSHIP

- Data is data, but you need structure and meaning to understand the data!
- Data models can be non-existent, bad, so-so or good!

PRIMARY TAKEAWAY:

**Named Relationships is what
matters most in
Graph Data Modeling!!!**

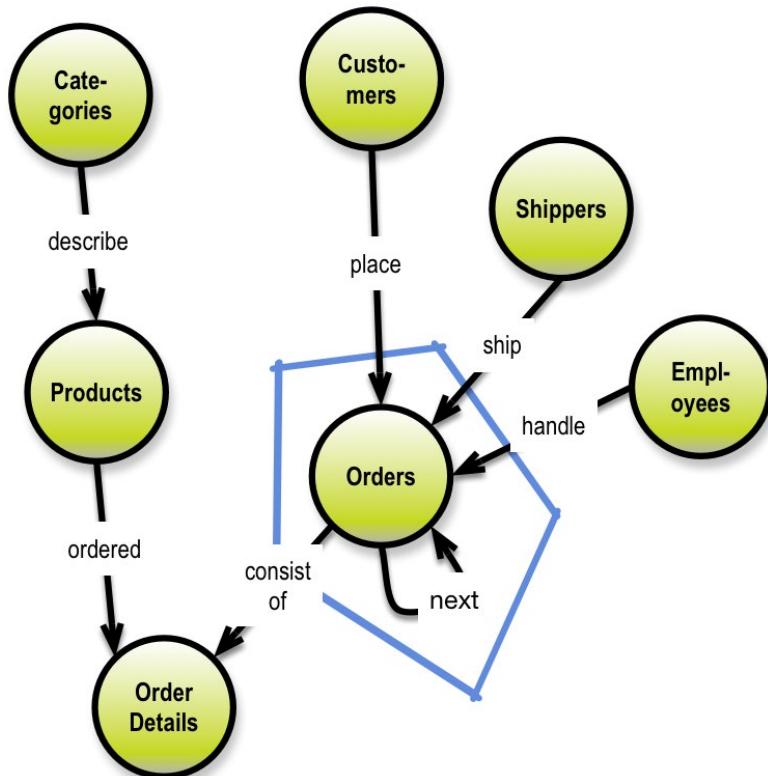




“THE ATOMS AND MOLECULES” OF DATA MODELS

- Concepts, which materialize as either
 - Object types (Customer, Order etc.), or
 - Properties (CustomerNumber, ProductName etc., some of which are mandatory)
- Dependencies, which structure the model as either:
 - Functional dependencies (internal within the object type and determining the primary key), or
 - Intra-object relationships (like “business address of” and so forth), which have
 - Cardinalities
 - An associated key, which can be either a single identity or a combined set of uniqueness criteria
 - a type (name, not always present)
 - a direction (from/to, not always present)
- Uniqueness criteria (object type level),
 - either an identifying property (key) or
 - (frequently) a list of concatenated properties
- Identity (object type level),
 - either a single identifying property (key) like CustomerNumber, or
 - a system generated surrogate key
- Data types

PROPERTY GRAPH DATABASES - THE UBIQUITOUS POINTER!

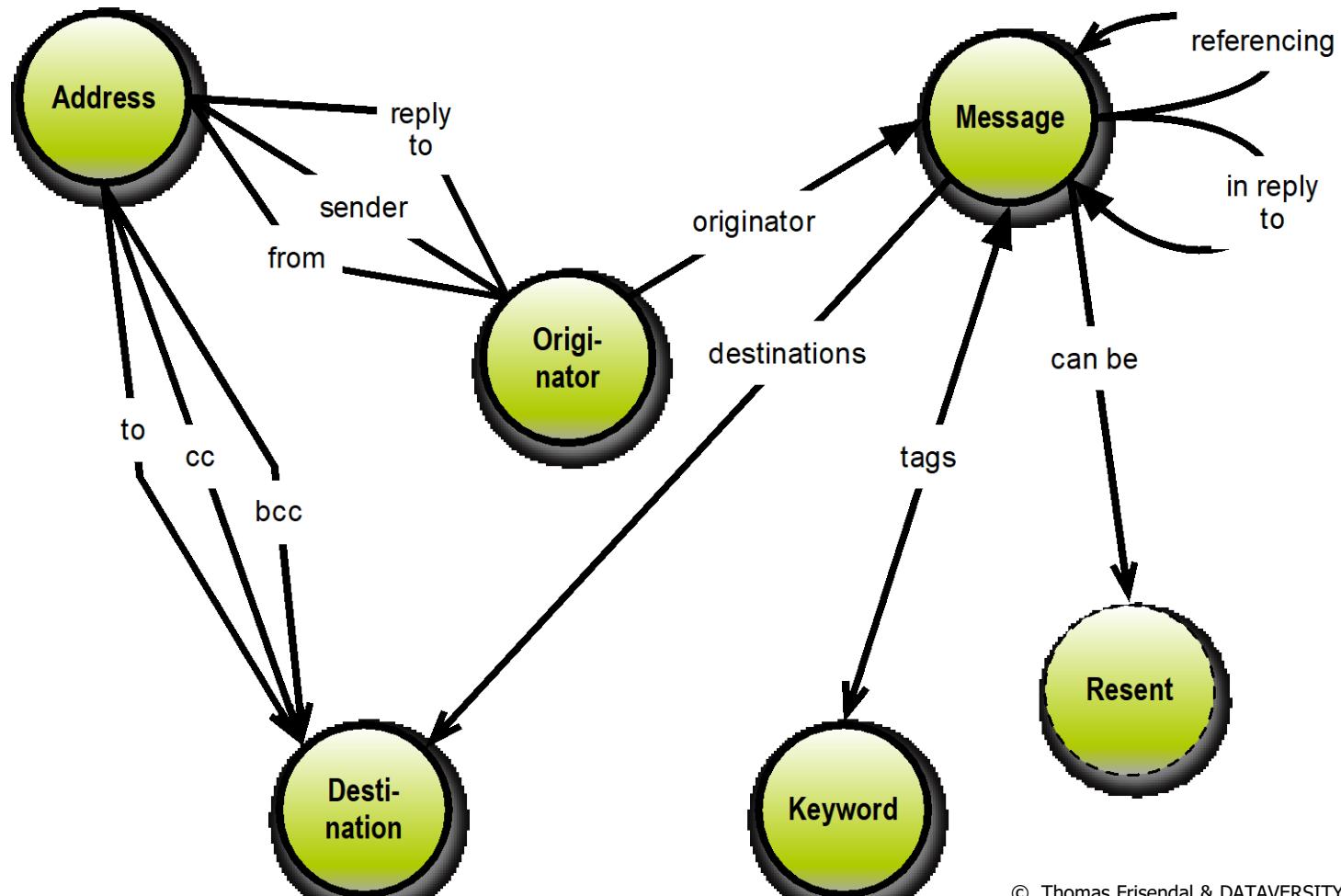


Next, prior, first, last, time-series

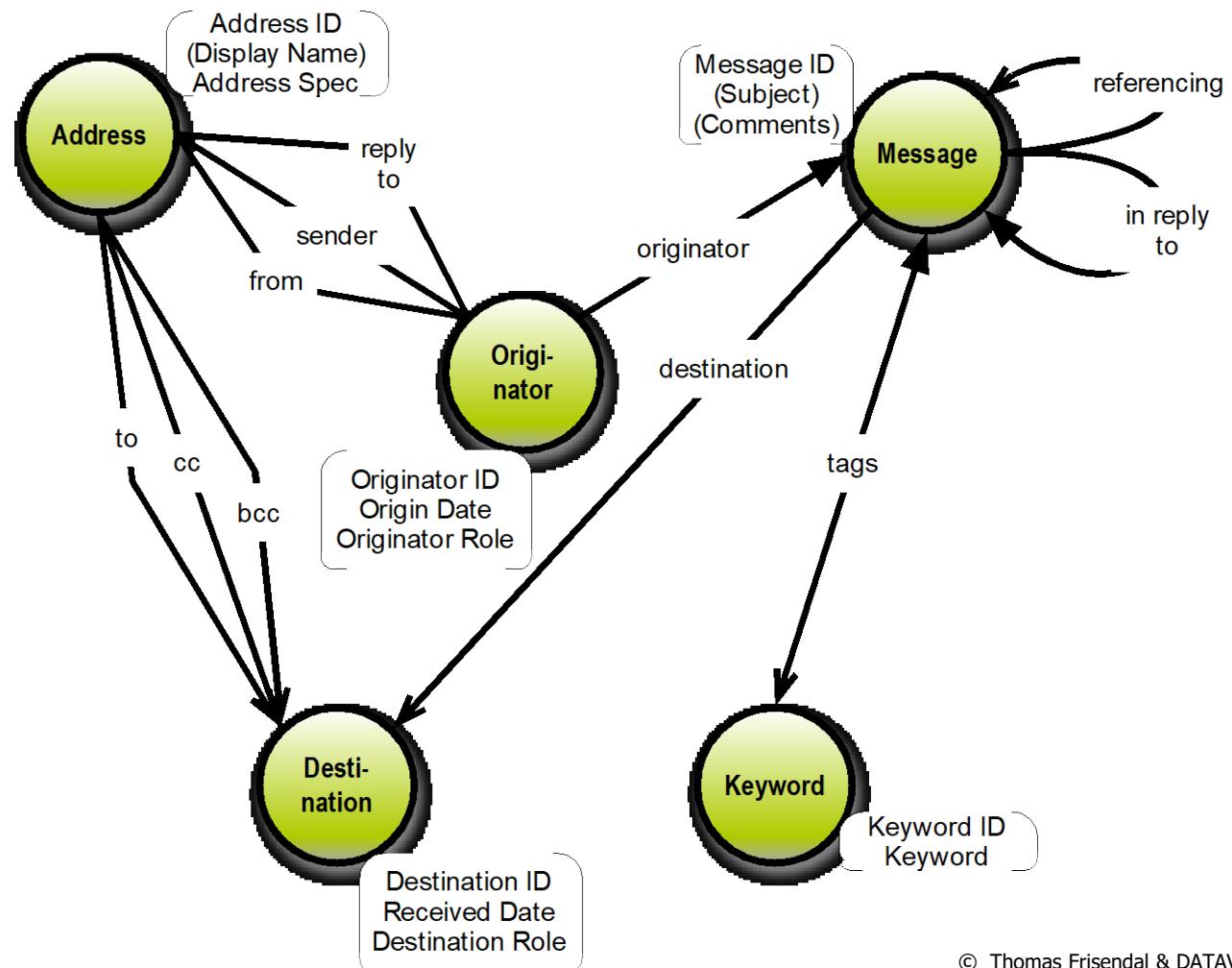


© Thomas Frisendal & DATAVERSITY Education, LLC | All Rights Reserved

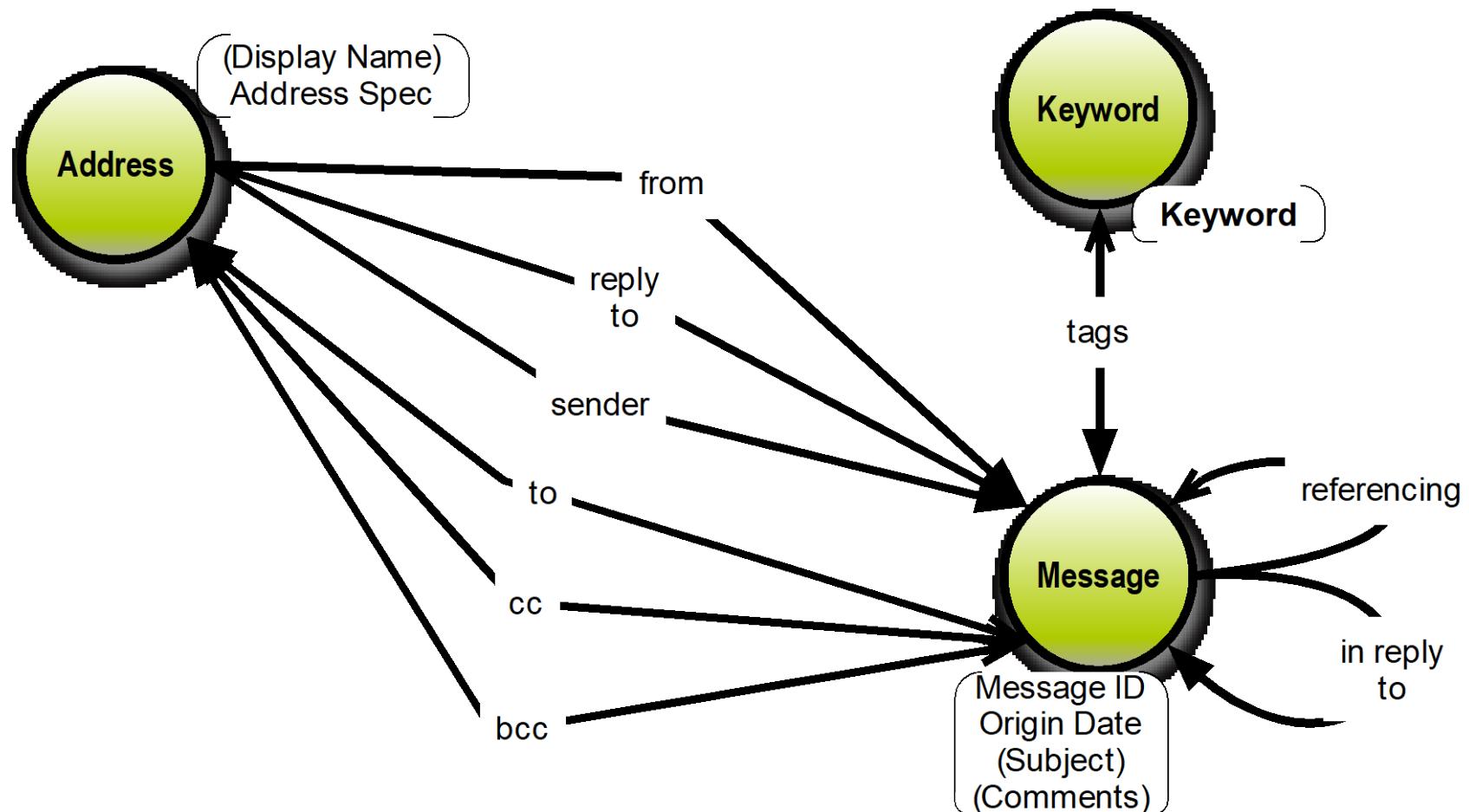
EMAIL CONCEPTS



NAIVE EMAIL PROPERTY GRAPH



A "GRAPHISH" EMAIL PROPERTY GRAPH



WHAT ABOUT M:M?

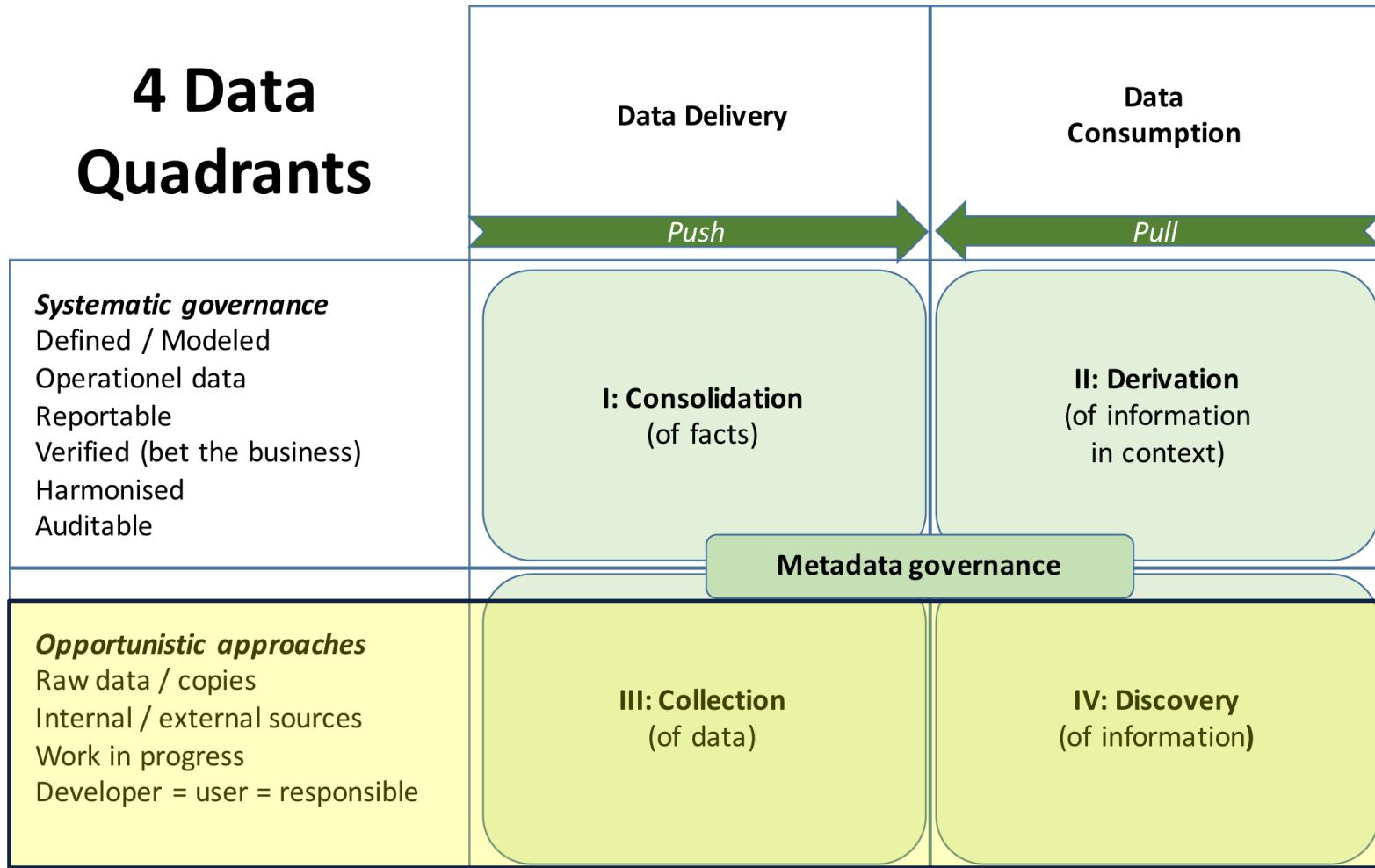
- **Relational:**
 - Relationship table (bridge)
 - No(?) properties?
 - "Factless facts" in star schemas
 - Non-information bearing M:M relationships are relatively seldom:
 - At closer inspection they most often carry information
 - A quantity, an amount, a percentage, something related to the relatedness
- **Graph**
 - Many to many works fine without additional modelling
 - But: No object-owned properties on relationships, except for simple things
 - Objects / events are nodes
 - Objects / events may be role-playing (a reply is also an email)
 - Objects / events may be state-changing (Sold, Invoiced, Paid for)

GRAPH DATA MODELING AGILITY

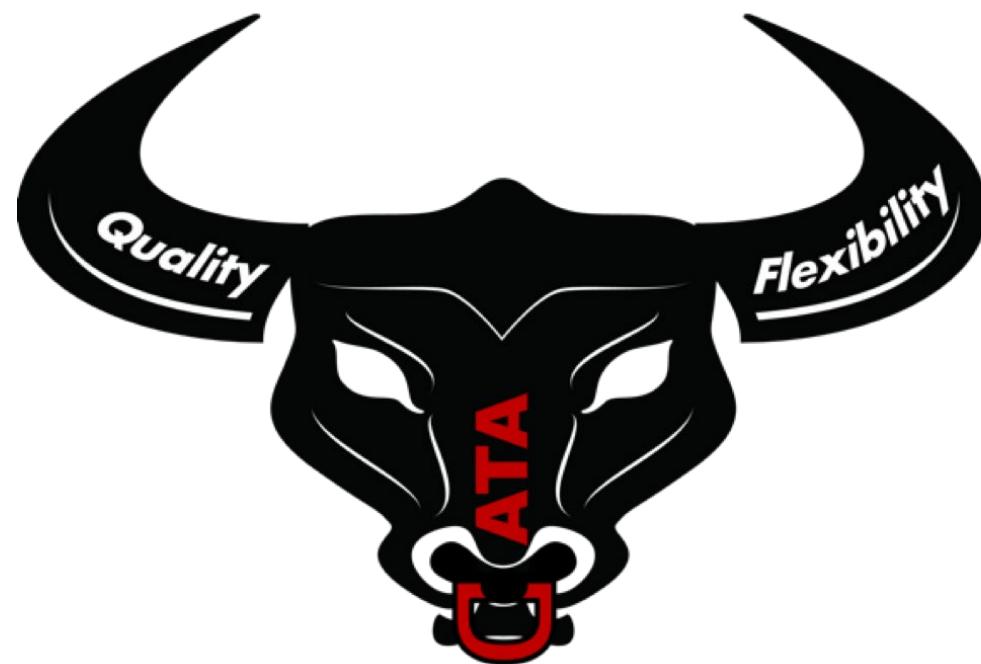
- Schema first
- Schema less
- Schema eventually

TWO-BY-TWO DIMENSIONS: MANAGE THE CONCERNS, WHICH ARE COMPETING

4 Data Quadrants

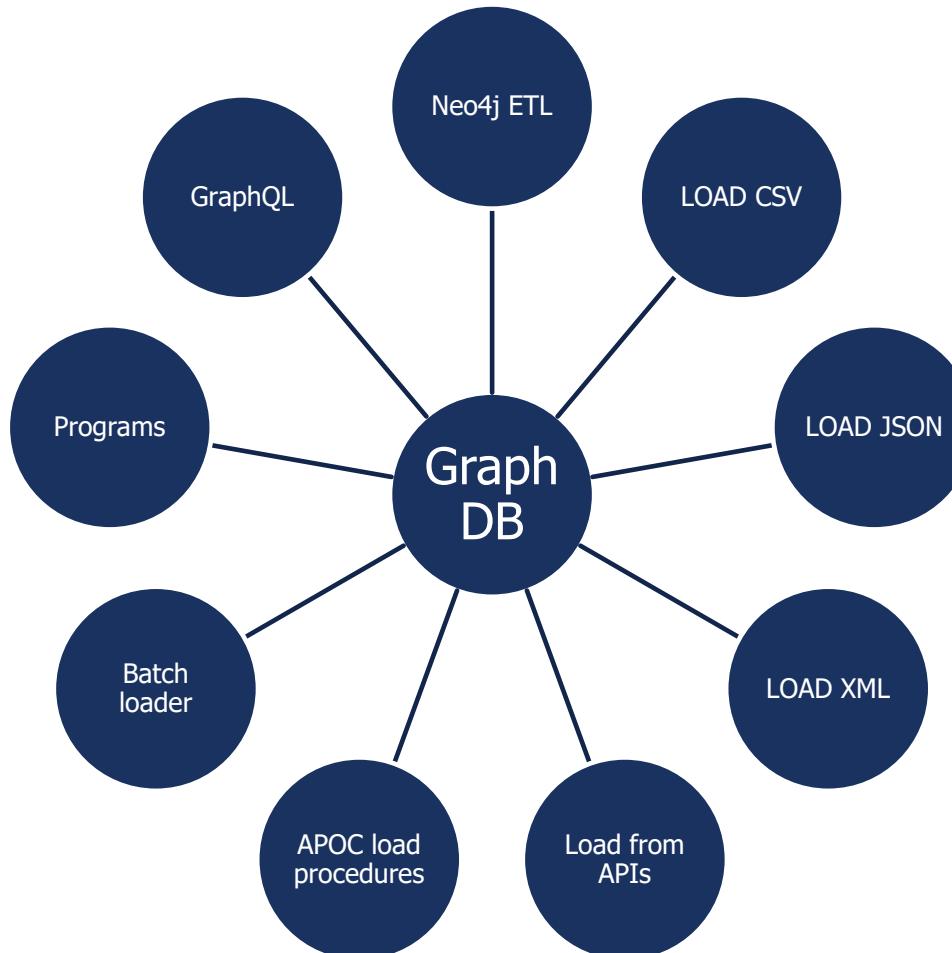


THE BULL IS COMING AT YOU!



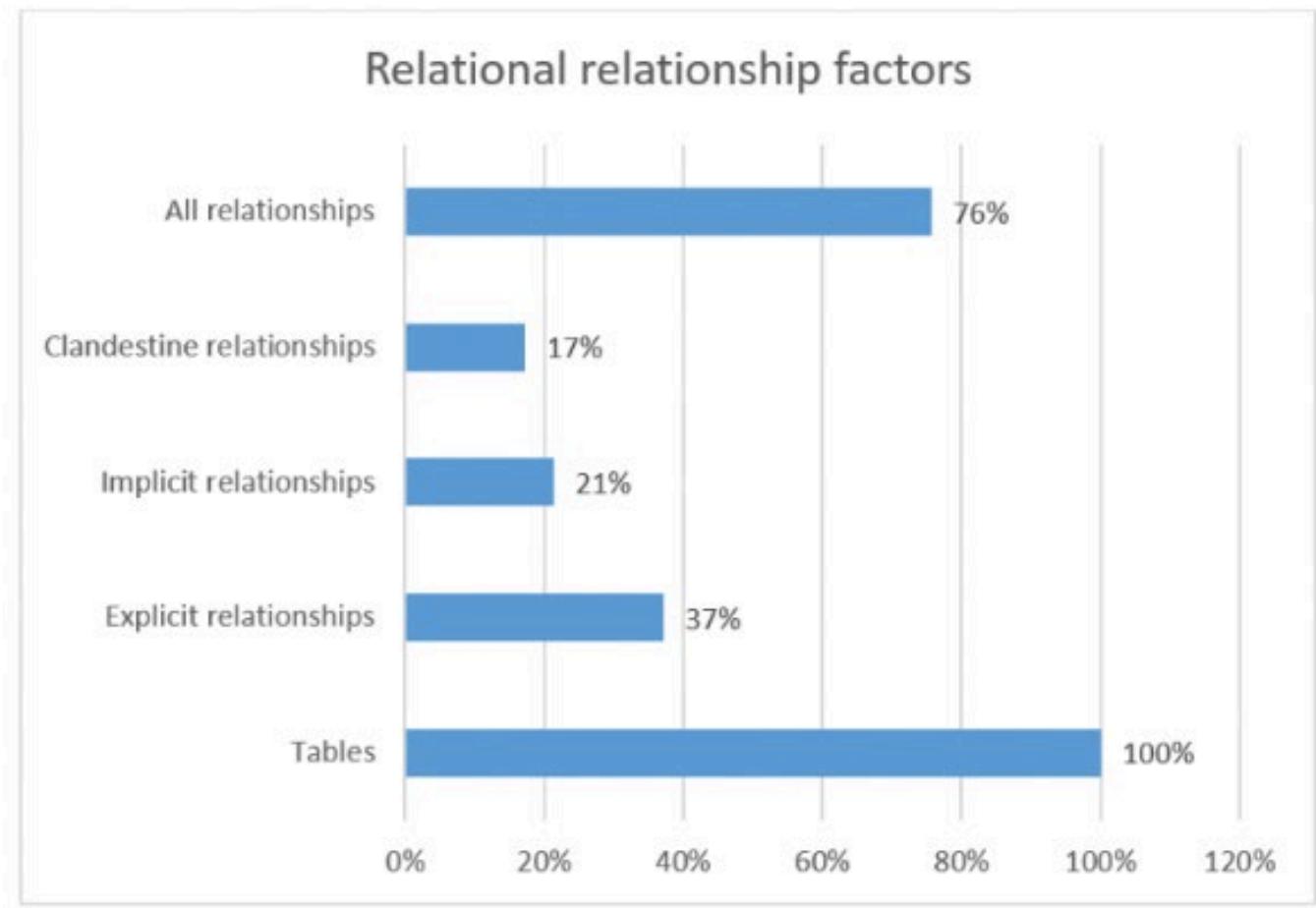
More on Facebook: <https://www.facebook.com/FullScaleDataArchitects/>,
and on Meetup: <https://www.meetup.com/Full-Scale-Data-Architects/>

WAYS TO LOAD DATA INTO A GRAPH DATABASE (NEO4J)



TWO MAJOR CHALLENGES

- What is the quality of the table and column names in the source?
- What is the quality of the relationship names?
 - Not all relationships declared (only 37 %)
 - Foreign key names?
 - FROM-TABLE_TO-TABLE?
 - PRIMARYKEY_FOREIGNKEY?
 - Physical primary keys are surrogate keys in many databases



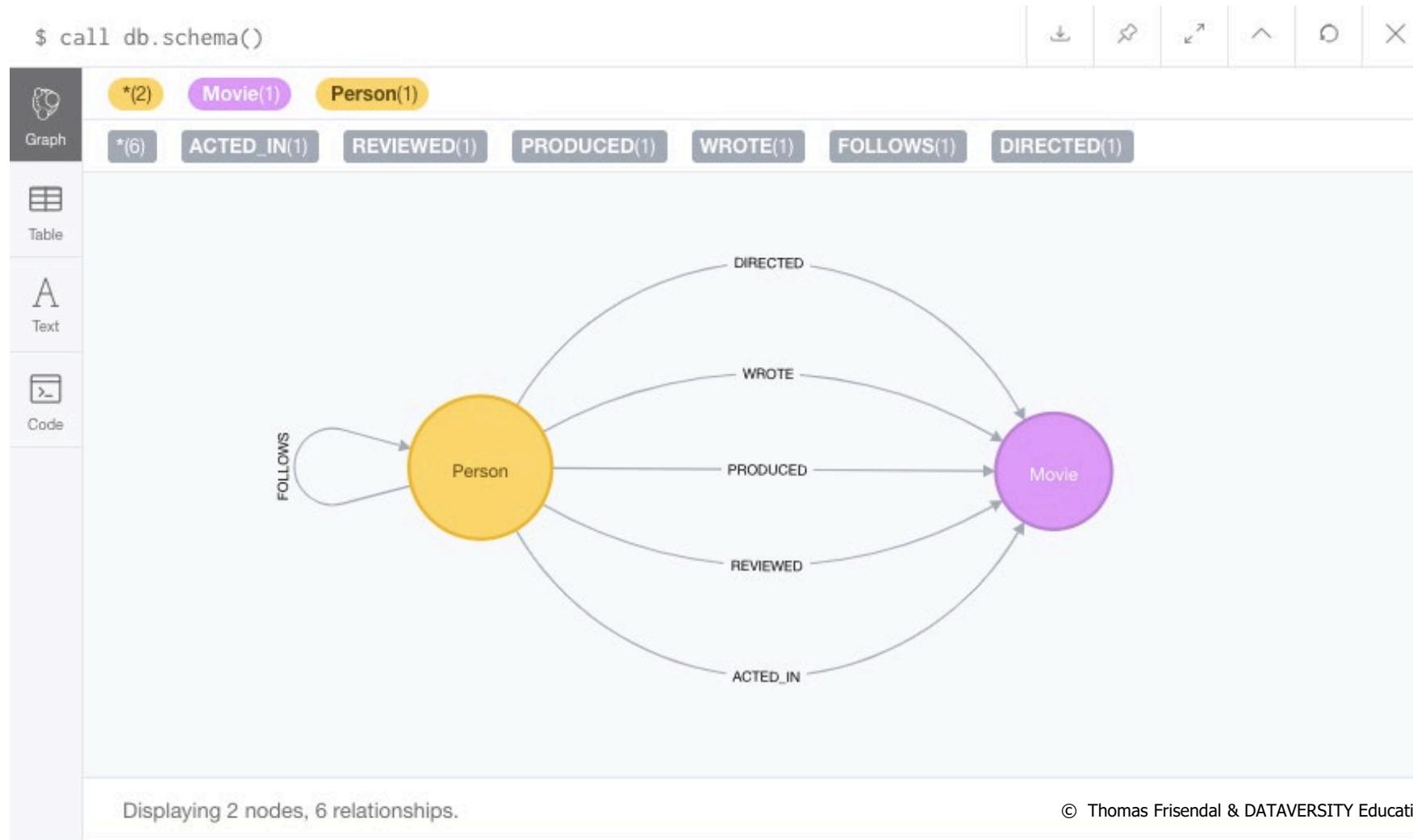
THE HUMAN FACTORS

- If data models are to be derived from the data, we need to master the Data Quality aspects rather completely
- A given data set may be interpreted in more than one way, and there are bound to be situations, where human judgement is necessary
- We are still better than machines when it comes to understanding a given context. Why is that? Because humans set the rules of the context.
- In our very human ways we leave room for oversights, simplifications, peculiar terminology, and laissez-faire “known errors.” The rightful owners of the business terminology are the business people, and they are more human than they are androids.
- The path from data values to a context description (another way to describe what a data model is) is long and works upwards in multiple layers.
- The bottom layers is traditional stuff from Data Quality but the upper layers really do not have a good name, yet. “Smart Data Discovery” is more a marketing term and there are no clearly sovereign technologies.

PROFILING YOUR GRAPH DATA

WHAT DID I LOAD?

INFERRING THE SCHEMA FROM THE DATA (NEO4J)



HOW MANY NODE TYPES AND RELATIONSHIPS?

```
MATCH (n)  
RETURN DISTINCT labels(n),  
count(*) AS SampleSize,  
avg(size(keys(n))) as Avg_PropertyCount,  
min(size(keys(n))) as Min_PropertyCount,  
max(size(keys(n))) as Max_PropertyCount,  
avg(size( (n)-[]-() ) ) as Avg_RelationshipCount,  
min(size( (n)-[]-() ) ) as Min_RelationshipCount,  
max(size( (n)-[]-() ) ) as Max_RelationshipCount
```

labels(n)	OccCount	Avg_PropertyCount	Min_PropertyCount	Max_PropertyCount	Avg_RelationshipCount	Min_RelationshipCount	Max_RelationshipCount
["Movie"]	38	2.973684210526316	2	3	6.578947368421054	2	14
["Person"]	133	1.9624060150375942	1	2	1.924812030075188	1	13

EVALUATE THE UNIQUENESS OF A PROPERTY WITHIN A LABEL

```
MATCH (m:Movie)  
RETURN count(DISTINCT m.title) AS DistinctTitle,  
       count(m.title) AS TotalTitle,  
       100*count(DISTINCT m.title)/count(m.title) AS Uniqueness
```



The screenshot shows a Neo4j browser window with a query results table. The table has three columns: 'DistinctTitle', 'TotalTitle', and 'Uniqueness'. The data row shows values 38, 38, and 100 respectively. The table icon in the sidebar is highlighted.

	DistinctTitle	TotalTitle	Uniqueness
38	38	100	

WHICH PROPERTIES ON RELATIONSHIPS?

```
$ call db.schema.relTypeProperties()
```



Table

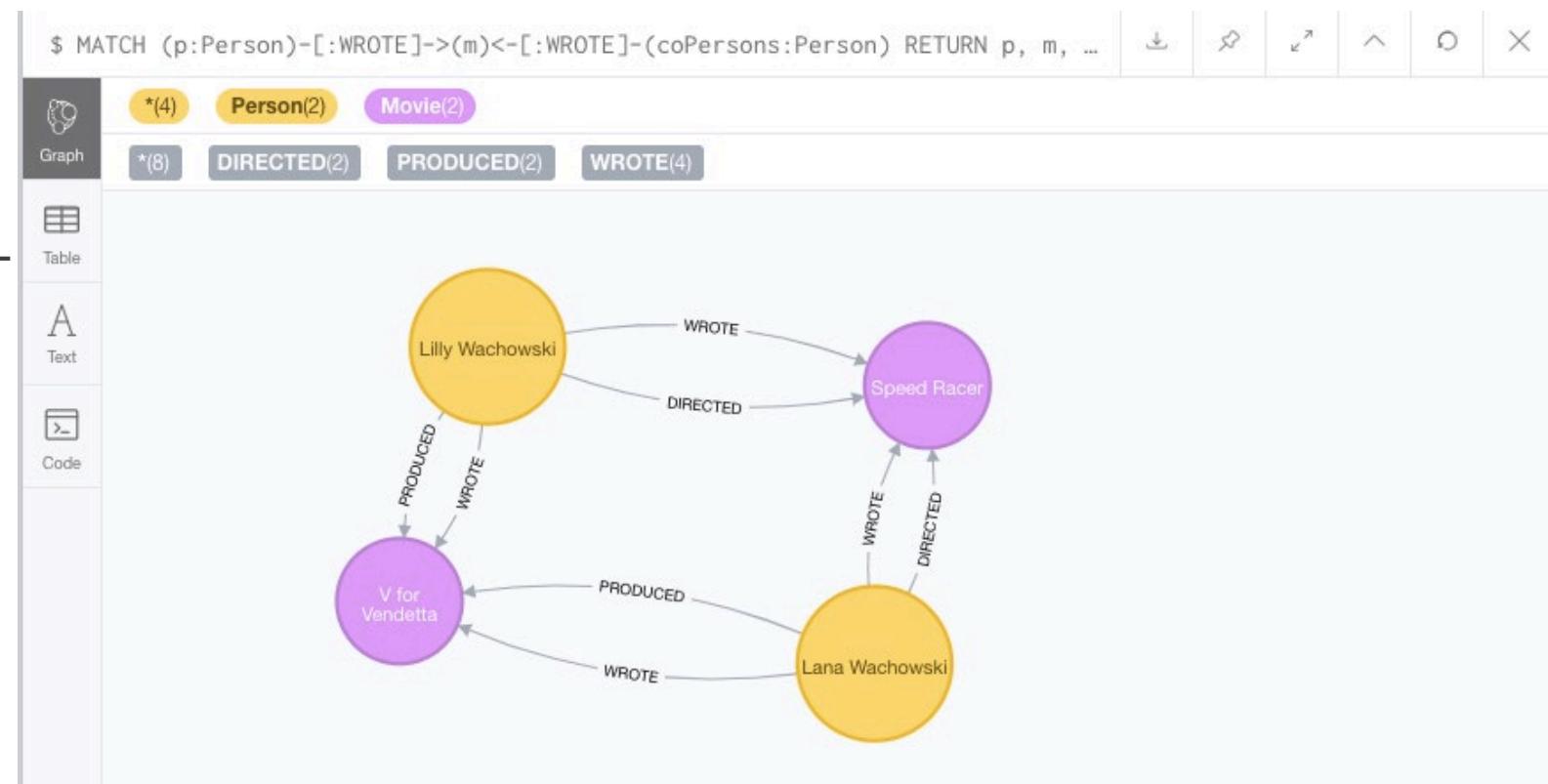
A
Text

Code

relType	propertyName	propertyTypes	mandatory
"':ACTED_IN'"	"roles"	["StringArray"]	true
"':DIRECTED'"	null	null	false
"':PRODUCED'"	null	null	false
"':WROTE'"	null	null	false
"':FOLLOWS'"	null	null	false
"':REVIEWED'"	"summary"	["String"]	true
"':REVIEWED'"	"rating"	["Long"]	true

FINDING M:M CARDINALITIES ACROSS RELATIONSHIPS

```
MATCH (p:Person)-[:WROTE]->(m)<-
[:WROTE]-(coPersons:Person)
RETURN p, m, coPersons
```



IDENTITY, UNIQUENESS, AND KEYS

- Surrogate keys
- Customer -> Order -> Orderline -> Product
- The “zero records”:
- Consider GUID’s

CustomerID	FirstName	LastName	FullName
0	Not specified	Not specified	Not specified
291	Gustavo	Achong	Gustavo Achong
293	Catherine	Abel	Catherine Abel
295	Kim	Abercrombie	Kim Abercrombie
297	Humberto	Acevedo	Humberto Acevedo
299	Pilar	Ackerman	Pilar Ackerman
301	Frances	Adams	Frances Adams
303	Margaret	Smith	Margaret Smith
305	Carla	Adams	Carla Adams
307	Jay	Adams	Jay Adams
309	Ronald	Adina	Ronald Adina
311	Samuel	Agcaoili	Samuel Agcaoili
313	James	Aguilar	James Aguilar
315	Robert	Ahlering	Robert Ahlering

GRAPH REFACTORING – PART OF THE APOC LIBRARY

- Just to give you a sense of what refactoring looks like:

```
MATCH (a1:Person{name:'John'}), (a2:Person {name:'Tom'})
```

```
WITH head(collect([a1,a2])) as nodes
```

```
CALL apoc.refactor.mergeNodes(nodes,{properties:"combine", mergeRels:true}) yield node
```

```
MATCH (n)-[r:WORKS_FOR]->(c) return *
```

- Learn more: <https://www.youtube.com/watch?v=wUMKM-uLNfg>

- You may, of course, also DELETE and (re-)load parts of the graph to make structural changes

SUMMARY

THE GRAPH DATA MODELING APPROACH

THE LIST OF CONCEPTS

Region

Region Name

Country

Segment

Country Code

Country Name

Company

Parent Company

Segment Name

Company Number

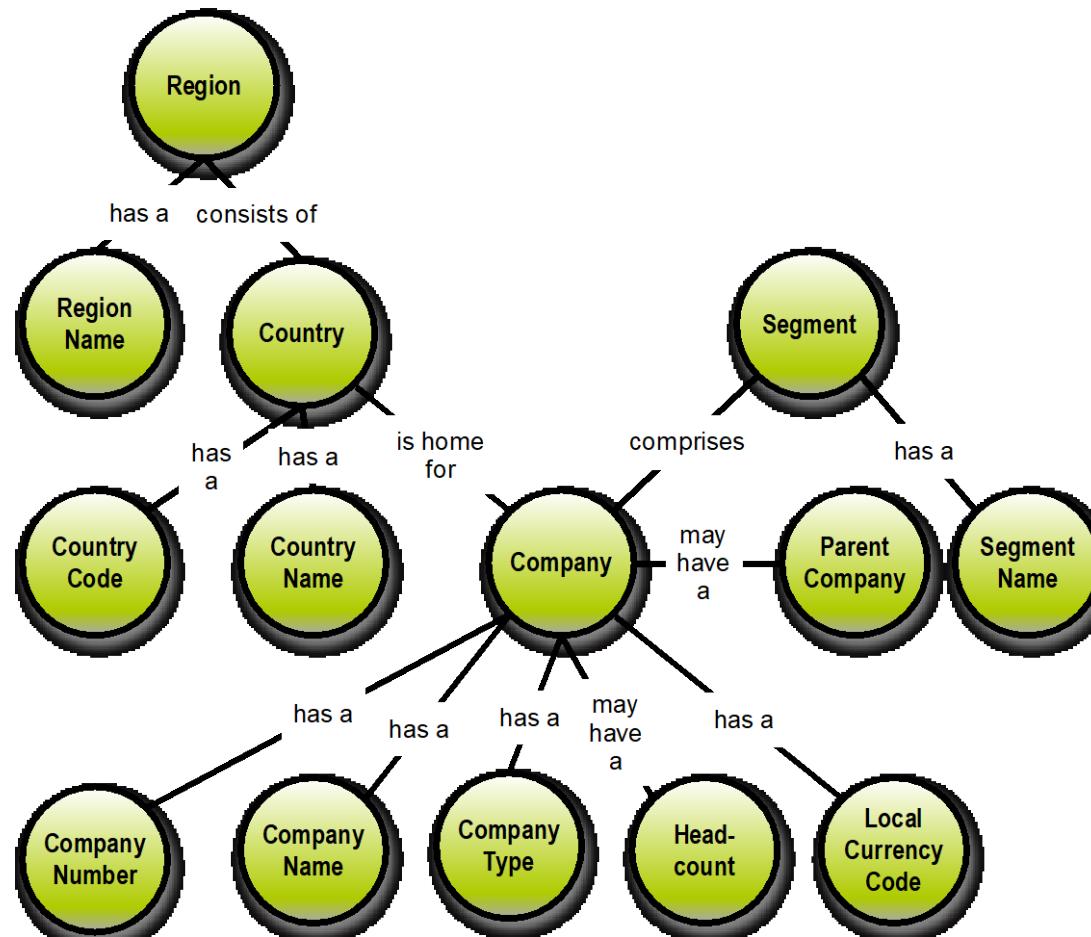
Company Name

Company Type

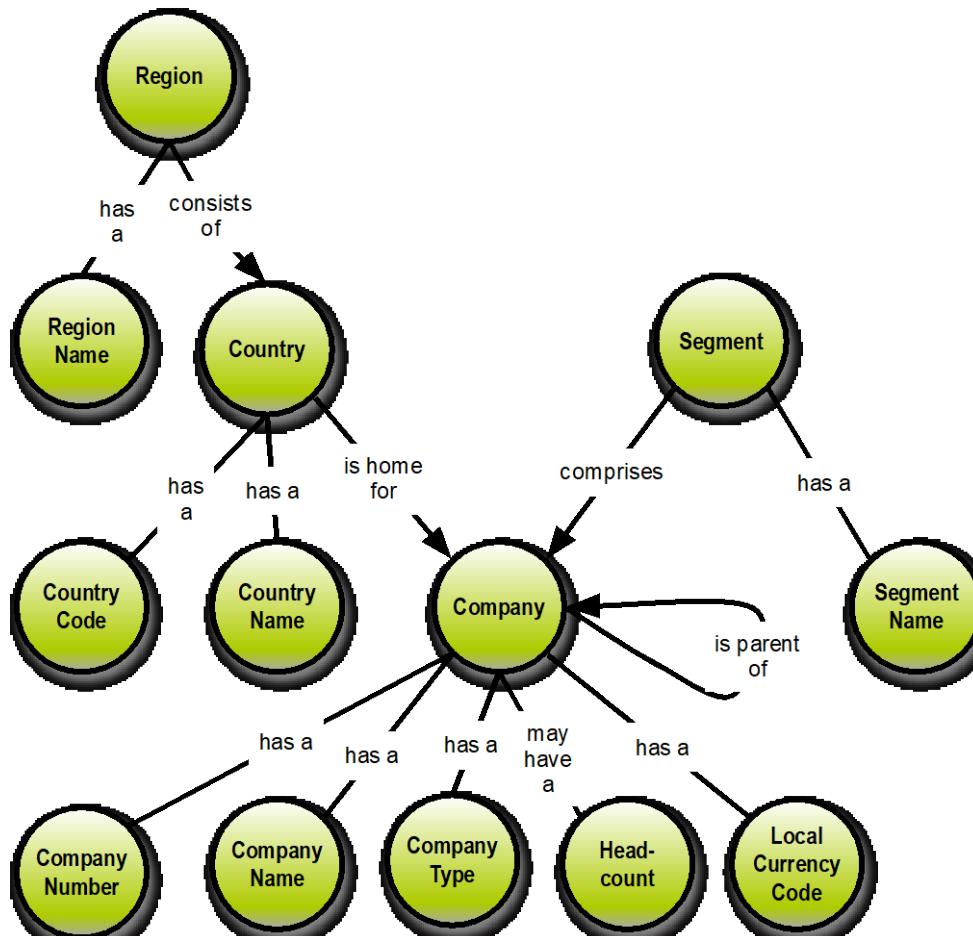
Head-count

Local Currency Code

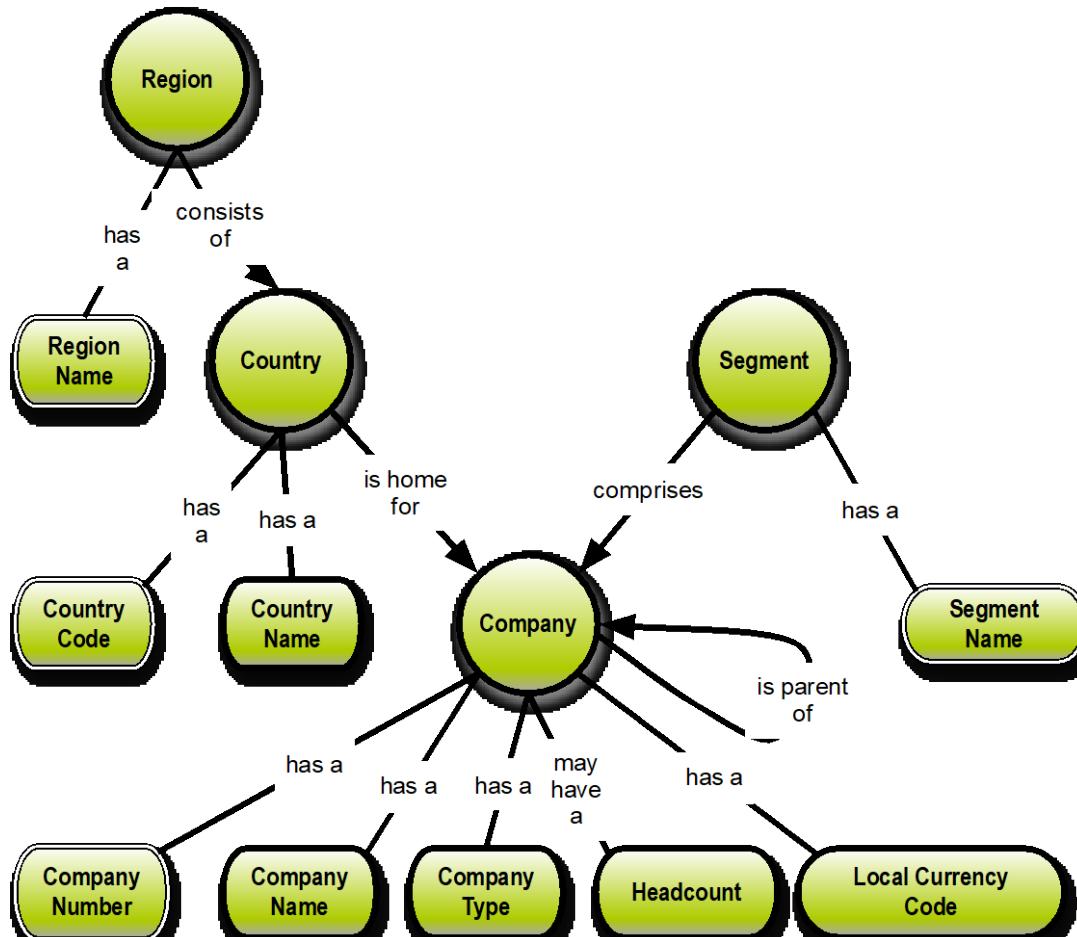
THE CONNECTED DATA MODEL



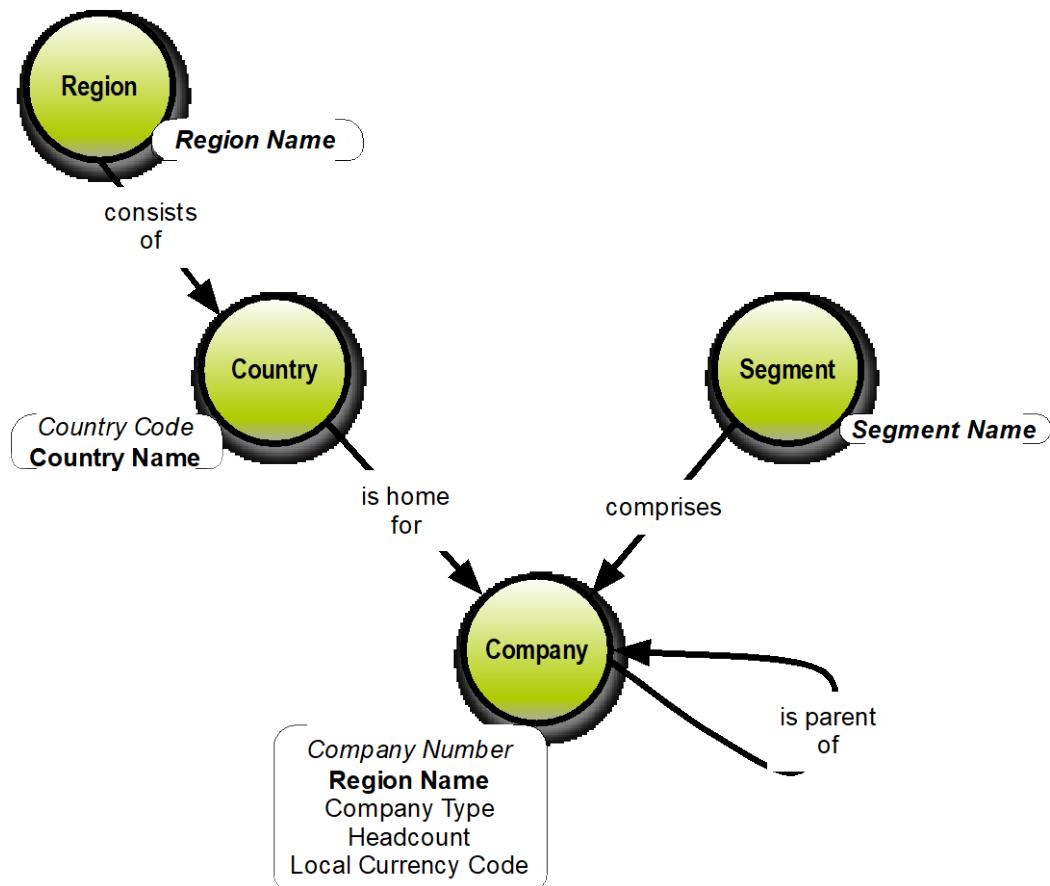
THE STRUCTURED DATA MODEL



THE GOVERNABLE DATA MODEL



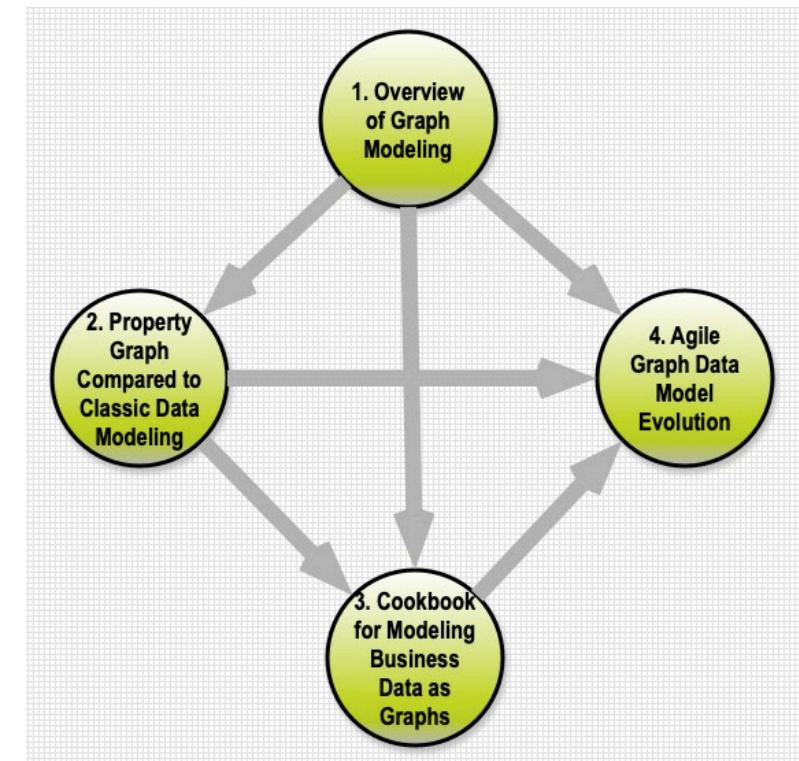
THE LABELED PROPERTY GRAPH VERSION OF THE GOVERNABLE DATA MODEL



MODELING DATA AS GRAPHS – COURSES

Modeling Data as Graphs Learning Plan at training.dataversity.net

- Course 1 – Overview of Graph Modeling
- Course 2 – Property Graph Compared to Classic Data Modeling
- Course 3 – Cookbook for Modeling Business Data as Graphs
- Course 4 – Agile Graph Data Model Evolution



THANK YOU!

- Thomas Frisendal
- Copenhagen (close to the airport)
- thomasf@tf-informatik.dk
- Member of the ISO IEC/JTC1/SC32/WG3 Database Languages Standards Committee
- @VizDataModeler
- www.graphdatamodeling.com
- <https://www.linkedin.com/in/thomas-frisendal-19a56a>

