



NOSQL stores, graph databases and data analytics tools

Advances in Data Management, 2016

Petra Selmer,

Birkbeck College (PhD candidate) and Developer at Neo Technology

Agenda

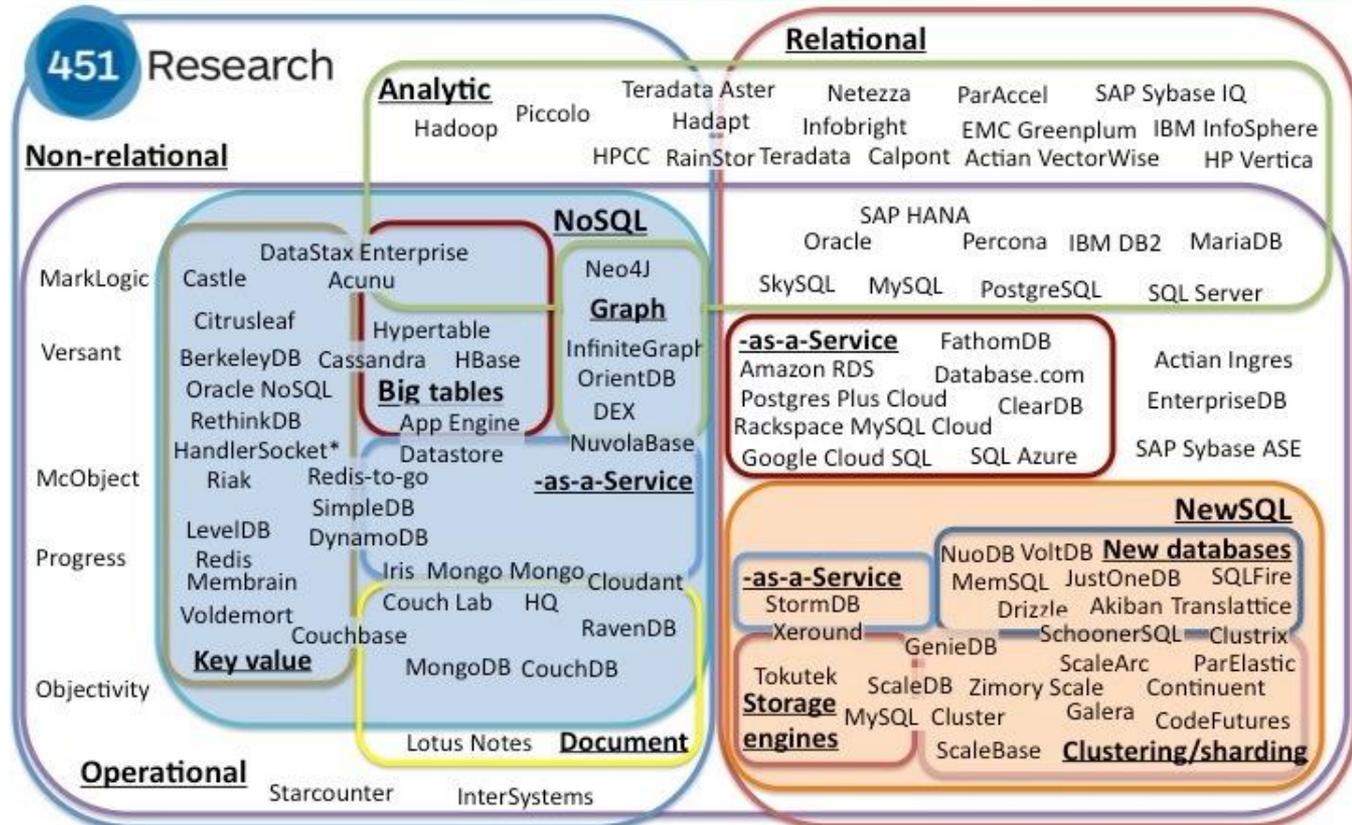
- × The wider landscape
- × NOSQL - i.e. non-relational stores
- × A closer look at Neo4j, a graph database
- × Data Analytics

Preamble

- × The area is HUGE
- × The area is ever-changing!

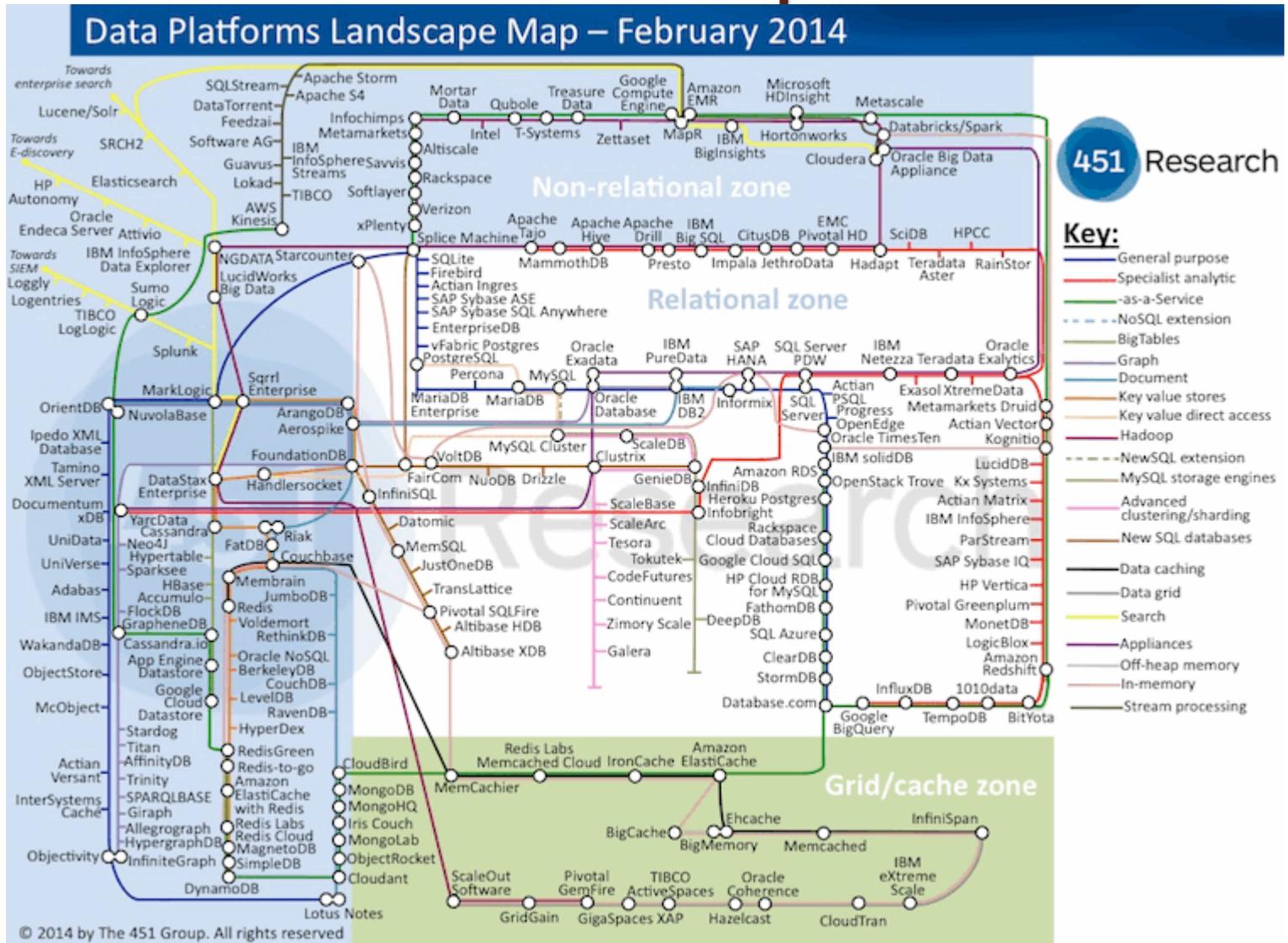
The wider landscape - 2012

The evolving database landscape



© 2012 by The 451 Group. All rights reserved

The wider landscape - 2014



The wider landscape

- × Several dimensions in one picture:
 - Relational vs. Non-relational
 - Analytic (batch, offline) vs. Operational (transactional, real-time)
- × Increasingly difficult to categorise these data stores:
 - Everyone is now trying fiercely to integrate features from databases found in other spaces.
- × The emergence of “multi-model” data stores:
 - ArangoDB and others.
 - One may start with one data model (e.g. Document model) and add other models (graph or key-value) as new requirements emerge.

NOSQL - non-relational

× NOSQL:

- “Not Only SQL”, not “No SQL”
- Basically means “not relational” – however this also doesn't quite apply, because graph data stores are very relational; they just track different forms of relationships than a traditional RDBMS. A more precise definition would be the union of different data management systems differing from Codd's classic relational model
- The name is not a really good one, because some of these support SQL and SQL is really orthogonal to the capabilities of these systems. However, tricky to find a suitable name.
- A good way to think of these is as “the rest of the databases that solve the rest of our problems”

× Scalability:

- Horizontal (scale out): the addition of more nodes (commodity servers) to a system (cluster) - simple NOSQL stores
- Vertical (scale up): the addition of more resources – CPU, memory – to a single machine

Non-relational vs. Relational

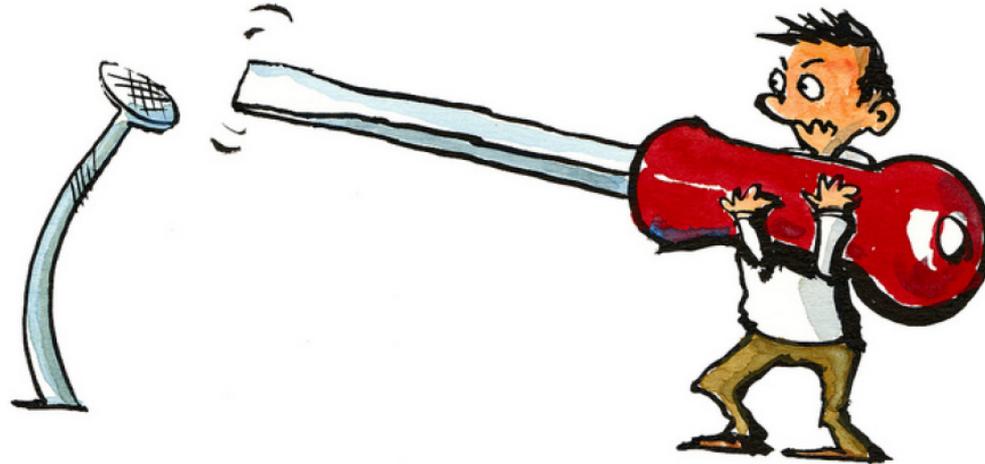
- × What's wrong with relational DBs? They're great!
 - ACID
 - Enforcement of referential integrity and constraints
 - SQL
 - Excellent support by many languages and technology stacks
 - Excellent tooling
 - Well-understood operational processes (DBAs): backups, recovery, tuning etc
 - Good security management (user access, groups etc)

NOSQL vs. Relational

- × BUT...there are problems:
 - Scaling with large and high-velocity data
 - × 'Big Data'
 - × Expensive / difficult / impossible to scale reads and writes vertically and horizontally
 - Complexity of data
 - × Impedance mismatch
 - × Performance issues (joins)
 - × Difficult to develop and maintain
 - Schema flexibility and evolution
 - × Not trivial
 - × Application downtime

Of hammers and nails...

The Law of the Hammer

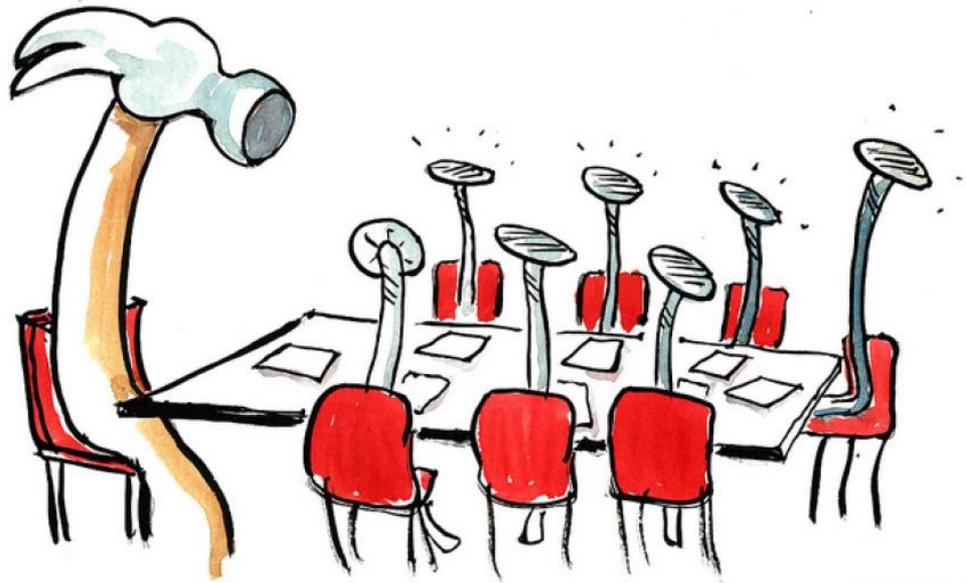


If the only tool you have is a hammer,
everything looks like a nail.

Abraham Maslow - The Psychology of Science - 1966

Of hammers and nails (cont)..

The Law of the Relational Database



By HikingArtist.com

If the only tool you have is a relational database,
everything looks like a table.

A Walk in Graph Databases - 2012

Non-relational vs. Relational

- × Why did these limitations not become problems before now?
 - In the past, there have been non-relational databases: Object Databases, XML Databases and proprietary formats; IBM's IMS, Lotus Notes...and Matlab
 - Types of systems – mostly operational and highly structured, and thus simpler:
 - × Payroll, inventory, stock management etc
 - × A simplified data model: employees previously only had one 'phone number, one title etc

NOSQL vs. Relational

- × Some datasets can be mapped easily to key-value pairs - flattening the data doesn't make it any less meaningful, and no reconstruction of its relationships is necessary.
- × For other datasets, the *relationship* to other items of data is as important as the items of data themselves.
- × Relational databases are based on **relational algebra (set theory)**:
 - Many datasets have relationships based on set theory, so an RDBMS is a good fit
 - However, for datasets where hierarchical or distance of relationships are required, set theory is not the best solution. In these cases, **graph theory** is a better match.
- × Summary:
 - RDBMS are too complex for data that can be effectively used as key-value pairs: we lose **scalability**
 - RDBMS are not complex enough for data that needs more context: we lose **performance**

Non-relational

- × Not intended as a replacement for RDBMS
- × One size doesn't fit all
- × Use the right tool for the job
- × Just as we shouldn't try to solve all of our problems with an RDBMS, we shouldn't try to solve all of our maths problems with set theory.
 - Today's data problems are getting complicated: the scalability, performance (low latency), and volume needs are greater.
 - In order to solve these problems, we're going to have to use an alternative data store or use more than one database technology.

NOSQL

- × Relational vs. Aggregate Data Model
- × Relational
 - Data are divided into rows (tuples) with pre-defined columns (attributes)
 - There is no nesting of tuples
 - There is no list of values
- × Aggregate
 - Think of this as a collection of related objects, which should be treated as a unit

Relational vs. Aggregate Data Model

Relational Instance

CUSTOMER	
ID	NAME
1	Guido

PRODUCT	
ID	NAME
1000	iPod Touch
1020	Monster Beat

BILLING_ADDRESS		
ID	CUSTOMER_ID	ADDRESS_ID
1	1	55

ADDRESS			
ID	STREET	CITY	POST_CODE
55	Chaumontweg	Spiegel	3095

ORDER		
ID	CUSTOMER_ID	SHIPPING_ADDRESS_ID
90	1	55

ORDER_ITEM			
ID	ORDER_ID	PRODUCT_ID	PRICE
1	90	1000	250.55
1	90	1020	199.55

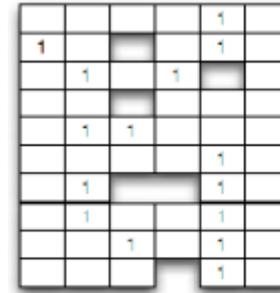
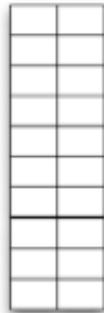
Aggregate Instance

```
{
  „id“:1,
  „name“:„Guido“,
  „billingAddress“:[{„street“:„Chaumontweg“,„city“:„Spiegel“,„postCode“:„3095“}]
}

{
  „id“:90,
  „customerId“:1,
  „orderItems“:[
    {
      „productId“:1000,„price“: 250.55, „produtName“: „iPod Touch“
    },
    {
      „productId“:1020,„price“: 199.55, „produtName“: „Monster Beat“
    }
  ],
  „shippingAddress“:[{„street“:„Chaumontweg“,„city“:„Spiegel“,„postCode“:„3095“}]
}
```

Non-relational Families

Key-Value

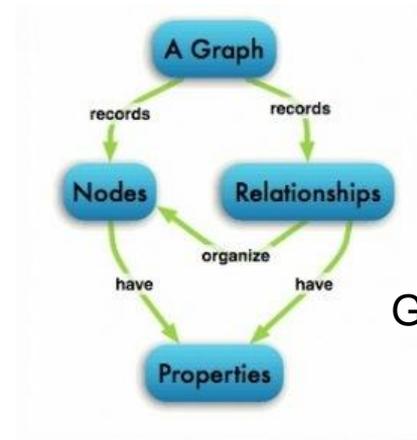


Column Store (also known as Big Table)



Document data model

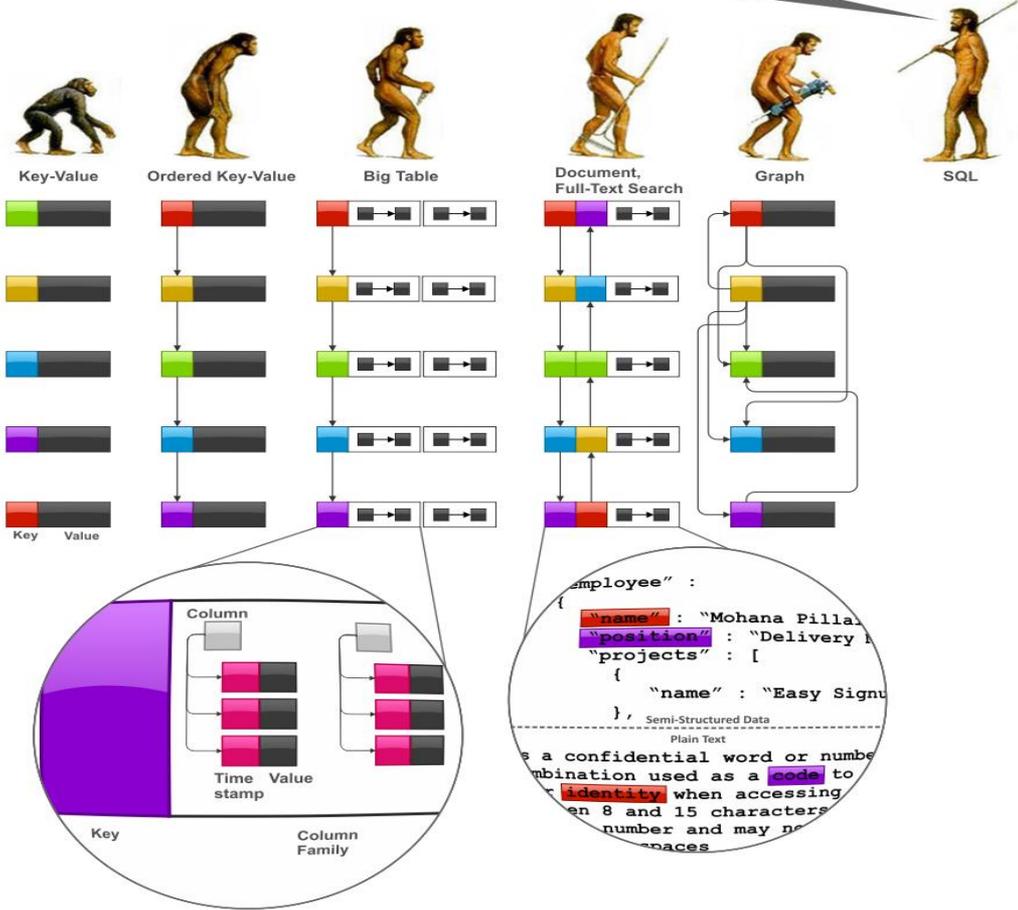
Collection of complex documents with arbitrary, nested data formats and varying "record" format.



Graph Store

NOSQL Families

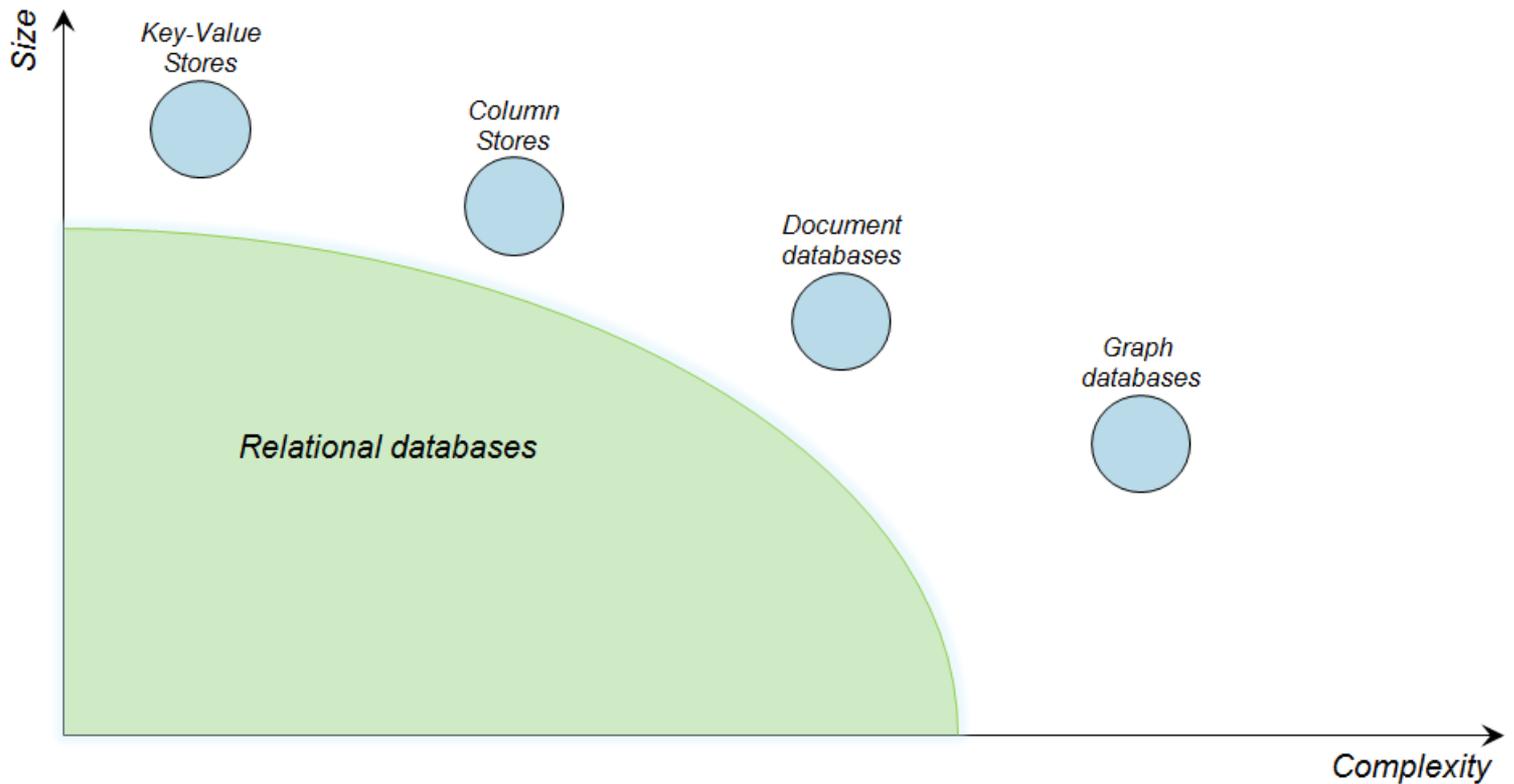
Stop following me, you freaks!



Non-relational Families

	Key/Value Store	Column Store	Document Store	Graph Store
Design	Key/Value pairs; indexed by Key	Columns and Column Families. Directly accesses the column values	Multiple Key/Value pairs form a document. Values may be nested documents or lists as well as scalar values	Focus on the <u>connections</u> between data and fast navigation through these connections
Scalability / Performance	+++	+++	++	++
Aggregate-Oriented	Yes	Yes	Yes	No
Complexity	+	++	++	+++
Inspiration / Relation	Berkley DB, Memcached, Distributed Hashmaps	SAP Sybase IQ, Google BigTable	Lotus Notes	Graph Theory
Products	Voldemort Redis Riak	HBase Cassandra Hypertable	MongoDB CouchDB Couchbase	Neo4j Sparksee InfiniteGraph [Triple and Quad Stores]

Non-relational Families





1) Key-Value Stores

- × History – Amazon decided that they always wanted the shopping basket to be available, but couldn't take a chance on RDBMS. So they built their own...
 - “Dynamo: Amazon's Highly Available Key-Value Store” (2007)
- × A key-value store is a simple hash table
- × Generally used when all access to the data is via a primary key
- × **Simplest** non-relational data store
- × Value is a BLOB → data store does not care or necessarily know what is 'inside'
- × Aggregate-oriented
- × Accessing and writing the data: PUT, GET, DELETE (matches REST)
- × Data model:
 - Global key-value mapping
 - Big scalable HashMap
 - Highly fault tolerant (typically)
- × Examples:
 - Riak, Redis



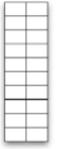
1) Key-Value Stores

× Strengths

- Simple data model
- Great at scaling out horizontally for reads and writes
 - × Scalable
 - × Available
 - × No database maintenance required when adding / removing columns

× Weaknesses:

- Simplistic data model – moves a lot of the complexity of the application **into** the application layer itself
- Poor for complex data
- Querying is simply by a given key: more complex querying not supported



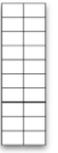
1) Key-Value Stores

× Suitable Use Cases

- Storing Session Information
- User Profiles, Preferences
- Shopping Cart Data
- Sensor data, log data, serving ads

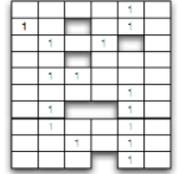
1) Key-Value Stores - Riak

Key-Value



- × Developed by Basho
- × Riak is a distributed database architected for:
 - Availability: replication of data means it is available for read and write operations, even in failure conditions;
 - Fault-tolerance: loss of access to many nodes owing to network partition or hardware failure does not mean a loss of data;
 - Operational simplicity: new machines can be added to the Riak cluster easily without incurring a larger operational burden;
 - Scalability: Riak automatically distributes data around the cluster and yields a near-linear performance increase as you add capacity.
- × See <http://basho.com/about/customers/> for a list of customers

2) Column Stores

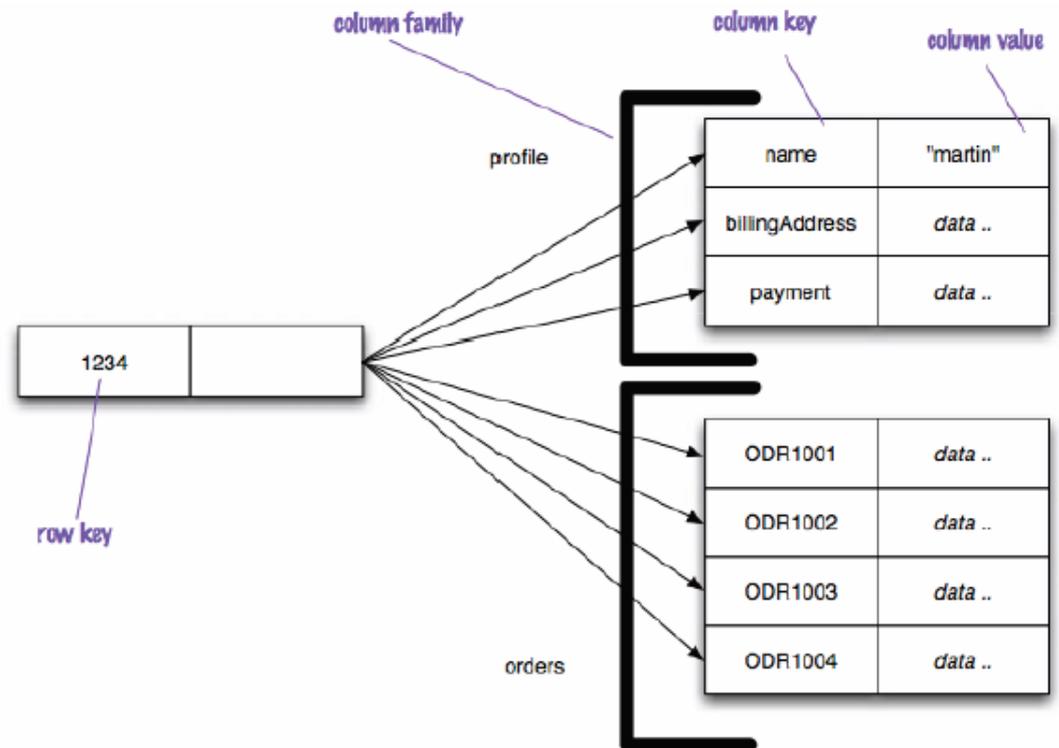


Column Store

- × Google's "Bigtable: A Distributed Storage System for Structured Data" (2006). Sometimes this family is called *Big Table*, *Wide Column* etc
- × Data model:
 - Rows are split across multiple nodes through sharding on the primary key
 - A big table, with *column families*. Column families are groups of related data, often accessed together
 - New columns may be added within the column family on a per-record basis, when needed. Lists of values may be stored in the column.
 - MapReduce for querying/processing
 - The records may be partitioned horizontally (sharded) across multiple servers, or parts of a SINGLE record may be stored on multiple servers (vertical partitioning)
- × Examples:
 - HBase, Cassandra
- × Aggregate-oriented

2) Column Stores - Example

- × One row for Customer 1234
- × Customer table partitioned into 2 column families: profile and orders
- × Each column family has columns (e.g. name and payment) and supercolumns (have a name and an arbitrary number of associated columns)
- × Each column family may be treated as a separate table in terms of sharding:
 - Profile for Customer 1234 may be on Node 1
 - Orders for Customer 1234 may be on Node 2



Source: NoSQL Distilled

2) Column Store

								1	
1									1
	1						1		
		1							
			1						
				1					
					1				
						1			
							1		
								1	

Column Store

- × Many key-value stores offer some form of grouping for columns and can be considered "column" stores as well.
- × Some databases – like HBase - were designed as column stores from the beginning:
 - This is a more advanced form of a key-value pair database. Essentially, the keys and values become composite.
 - Think of this as a hash map crossed with a multidimensional array. Essentially each column contains a row of data.

2) Column Stores

								1	
1								1	
	1							1	
		1	1						
									1
	1							1	

Column Store

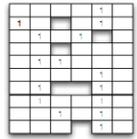
× Strengths

- Data model supports (sparse) semi-structured data
- Naturally indexed (columns)
- Good at scaling out horizontally
- MapReduce is very often used on these, so they can be good analytical stores for semi-structured data
- Can see results of queries in real time

× Weaknesses:

- Unsuitable for interconnected data: if the relationships between the data are as important as the data itself (such as distance or path calculations), then **don't** use a column store
- Unsuitable for complex data reads and querying
- Require maintenance – when adding / removing columns and grouping them
- Queries need to be pre-written; no ad-hoc queries defined “on the fly” : NOT for use for non real-time, unknown queries

2) Column Store

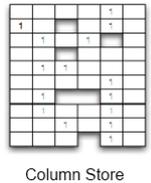


Column Store

× Use cases

- Netflix use it for logging and customer analytics, among others
- Ebay use it for search optimisation
- Adobe use it for structured data processing and Business Intelligence (BI)
- Used for 'firehose' data for TV shows such as BGT, The X Factor etc (audience and viewer voting): high amount of writes, and fast real-time basic analytics (Cassandra)
- Event Logging
- Counters
- Smart meters and monitoring
- Sensor data
- Mobile information

2) Column Stores - Cassandra



- × Apache project; also distributed by third-parties such as Datastax (who provide extra functionality on top of the base technology)
- × **Model:**
 - Column oriented, key value.
 - The values are split into columns which are pre-indexed before the information can be retrieved.
 - Eventually consistent.
 - This makes it better for highly distributed use cases or ones where the data is spread over an unreliable networks – lends itself well to geographically-distributed networks, in particular.
- × Robin Schumacher (VP of products for DataStax): "A popular use case for Cassandra is time series data, which can come from devices, sensors, websites (e.g., Web logs), financial tick data, etc. The data typically comes in at a high rate of speed, can come from multiple locations at once, adds up quickly, and requires fast write capabilities as well as high-performance reads over time slices."
- × Real-time query examples at which Cassandra excels:
 - Give me X
 - How many Y?
 - What is the top K?
 - How many distinct P in Q?

3) Document Stores



Document data model
Collection of complex documents with arbitrary, nested data formats and varying "record" format.

- × Documents are the main concept
- × Data model
 - Collections of documents
 - A document is a key-value collection
 - Index-centric: primary as well as secondary
- × Stores and retrieves documents, which can be XML, JSON, BSON..
- × Documents are self-describing, hierarchical tree data structures which can consist of maps, collections and scalar values, as well as nested documents
- × Documents stored are similar to each other but do not have to be exactly the same
- × Aggregate-oriented
- × Examples
 - MongoDB, Couchbase

```
{
  person: {
    first_name: "Peter",
    last_name: "Peterson",
    addresses: [
      {street: "123 Peter St"},
      {street: "504 Not Peter St"}
    ],
  }
}
```

3) Document Stores



Document data model
Collection of complex documents with arbitrary, nested data formats and varying "record" format.

× Strengths

- Simple but powerful data model – able to express nested structures
- Good scaling (especially if sharding supported)
- No database maintenance required to add / remove 'columns'
- Powerful query expressivity (especially with nested structures) – able to pose fairly sophisticated queries

× Weaknesses:

- Unsuitable for interconnected data
- Query model limited to keys (and indexes)
 - × MapReduce for larger queries (thus, might be slow)

3) Document Stores



Document data model
Collection of complex documents with
arbitrary, nested data formats and
varying "record" format.

- × Both key-value stores and document stores talk about “key-value” pairs – what is the difference?
 - From clustering to accessing data, document stores and key-value stores are exactly the same, except in a document store, the store understands the documents in the data store - the values are JSON, and the elements inside the JSON document can be indexed for better querying and search.
 - Because of this, the querying semantics within the document store will be much richer
- × When to consider using a document store (from a data model point of view)?
 - Your row schema is changing quickly over time and hence becomes too complex to model in a relational database.
 - If you don't have interconnected data

3) Document Stores - MongoDB

- Examples of enterprise uses of MongoDB at <http://www.mongodb.com/who-uses-mongodb>
- Financial Services use cases (<http://www.mongodb.com/presentations/webinar-how-and-why-leading-investment-organisations-are-moving-mongodb>):
 - High Volume Data Feeds
 - Tick Data capture
 - Risk Analytics & Reporting
 - Product Catalogs & Trade Capture
 - Portfolio and Position Reporting
 - Reference Data Management
 - Portfolio Management
 - Quantitative Analysis
 - Automated Trading
- CERN's Large Hadron Collider – DAS (Data Aggregation System):
 - <http://www.computerweekly.com/news/2240166469/Cornell-Cern-project-plumps-for-NoSQL-DBMS>
 - <http://www.slideshare.net/vkuznet/mongodb-at-the-energy-frontier>

3) Document Stores – use case



Document data model
Collection of complex documents with
arbitrary, nested data formats and
varying "record" format.

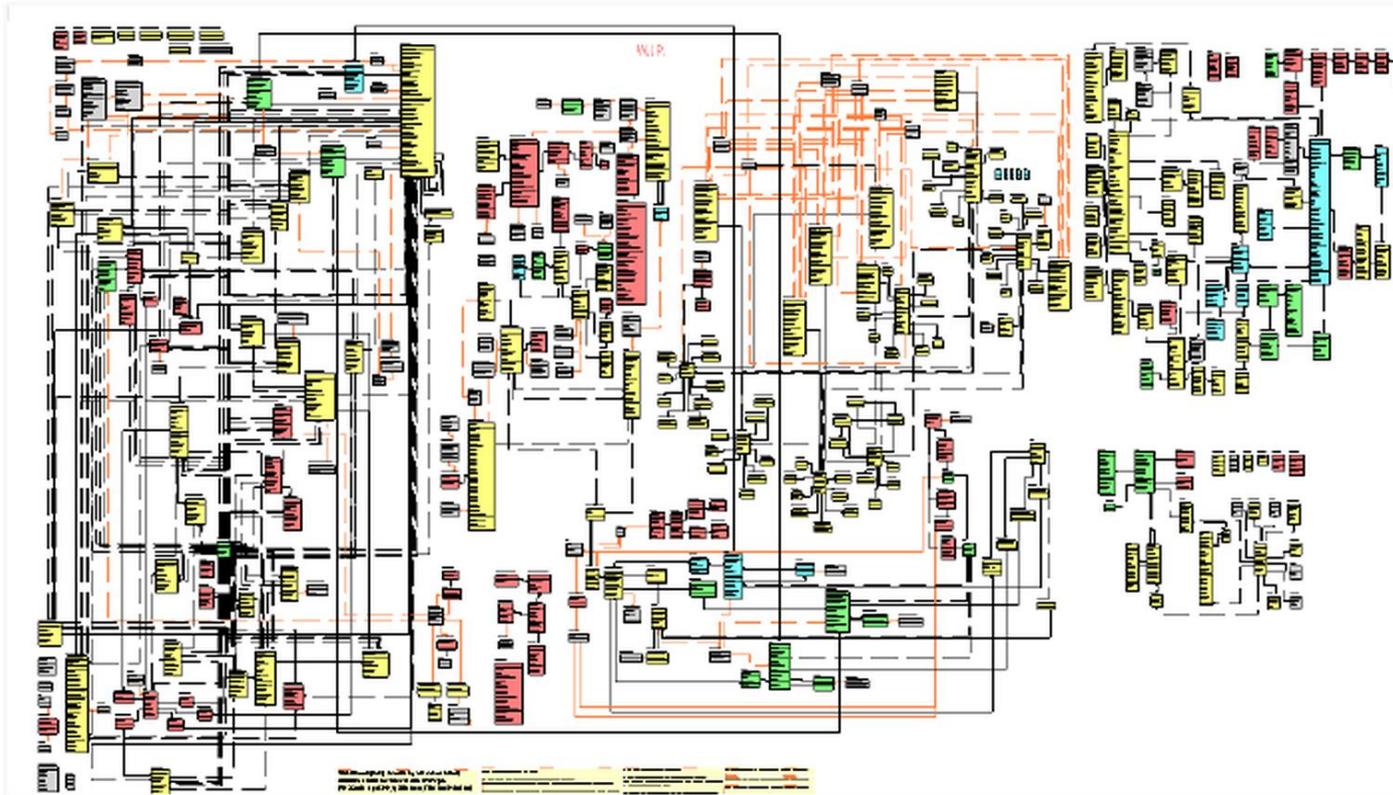
- × National Archives project:
 - <http://www.slideshare.net/AleksDrozdov/from-sql-server-to-mongo-d-bv10>
 - The National Archives is one of the world's largest records repositories, holding more than 11 million records spanning the Magna Carta to modern government papers – all of which is available to the public
 - Very interesting case of how they ran into massive problems with an RDBMS, and had to change the architecture to cope with the volume and heterogeneity of the data
- × Their systems encompassed:
 - Metadata
 - Digital images
 - e-Commerce; inventory; orders
 - User accounts; history
 - User participation
- × They have many searches through the catalogue
- × They had 2 000 tables and 56 000 attributes in SQL Server!
- × <http://discovery.nationalarchives.gov.uk>

3) Document Stores



Document data model
Collection of complex documents with arbitrary, nested data formats and varying "record" format.

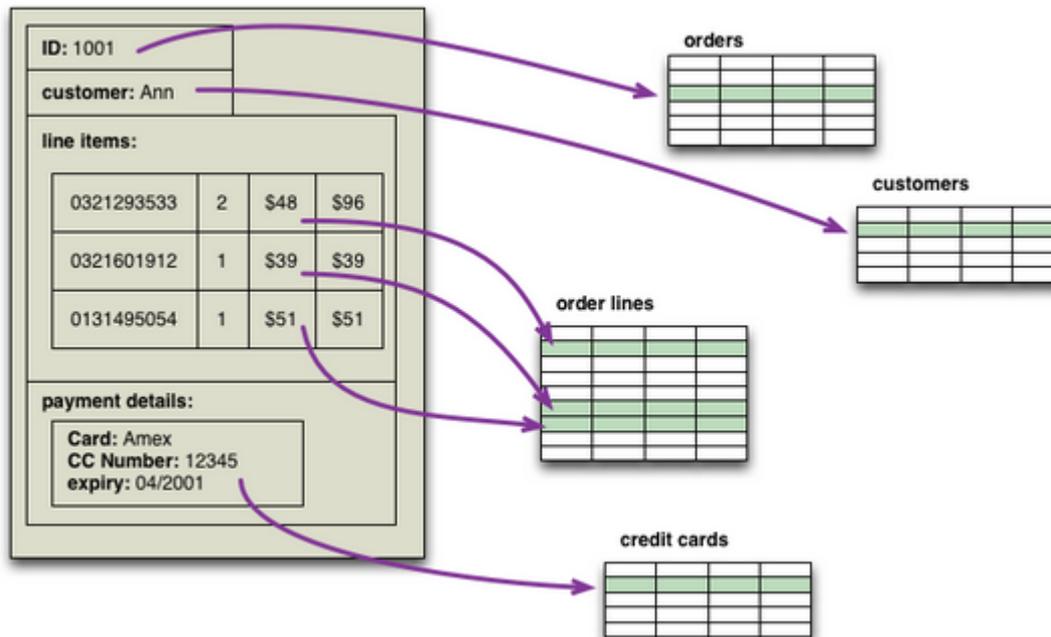
RELATIONAL MODEL



Aggregate-Oriented Databases

× <http://martinfowler.com/bliki/AggregateOrientedDatabase.html>

× “There's a big similarity between [key-value stores, column stores and document stores] - all have a fundamental unit of storage which is a rich structure of closely related data: for key-value stores it's the value, for document stores it's the document, and for column-family stores it's the column family. This group of data is an aggregate. An aggregate makes a lot of sense to an application programmer. If you're capturing a screenful of information and storing it in a relational database, you have to decompose that information into rows before storing it away. An aggregate makes for a much simpler mapping - which is why many early adopters of NoSQL databases report that it's an easier programming model.”





DENORMALISE

Aggregate data into documents

RICHER MODEL

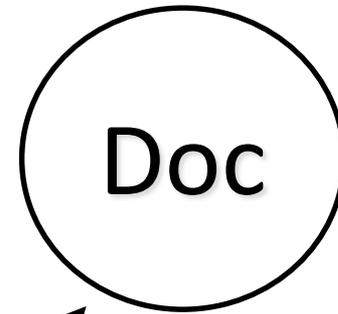
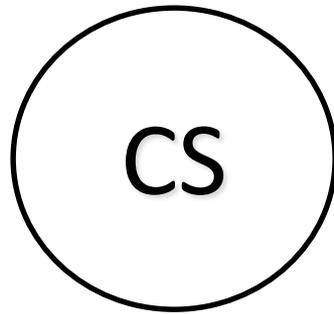
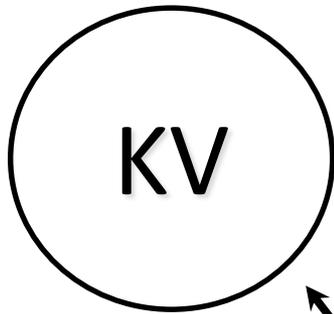
Connected structured data



*Simple data model
Map-reduce friendly*

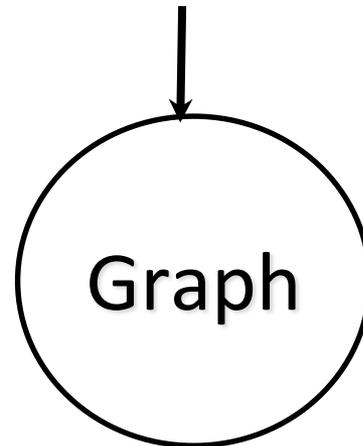
*Expressive power
Fast graph traversals*





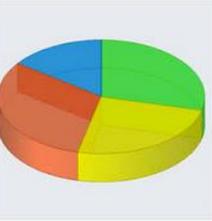
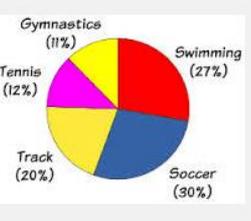
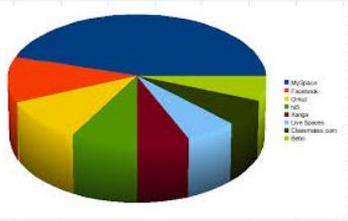
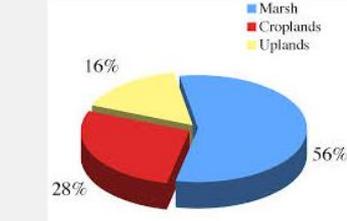
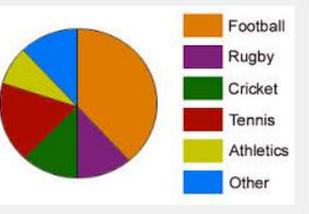
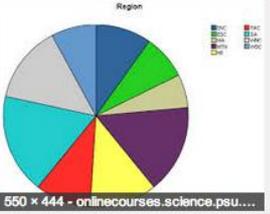
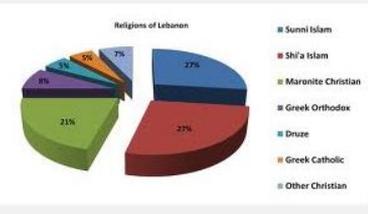
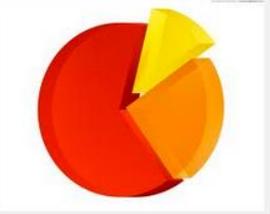
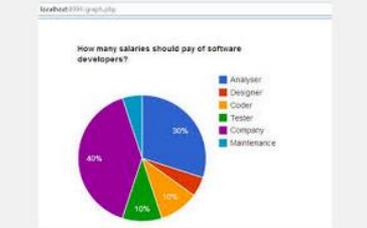
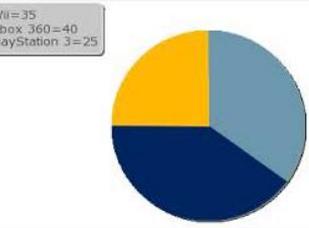
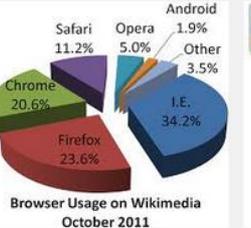
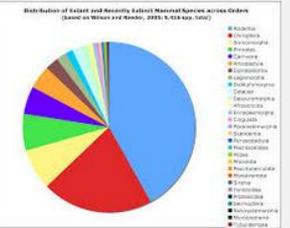
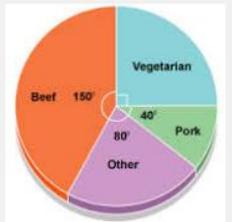
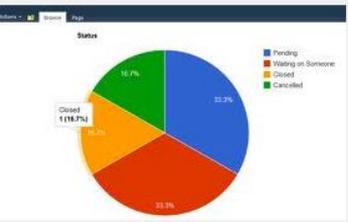
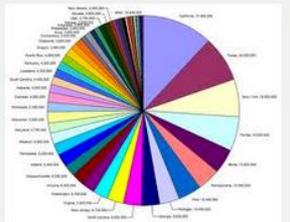
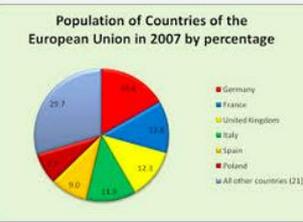
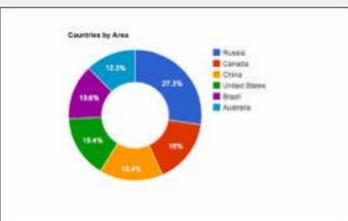
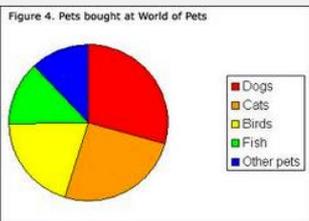
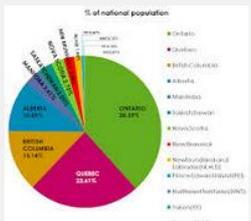
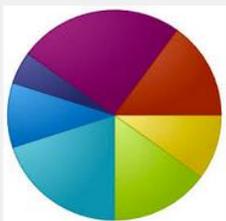
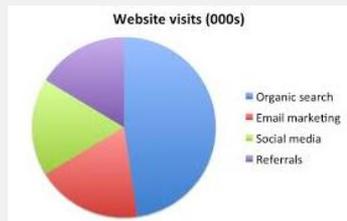
Denormalise

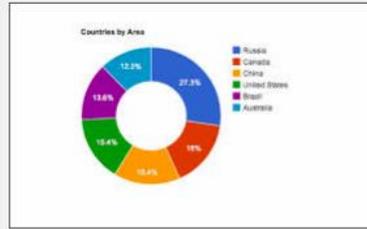
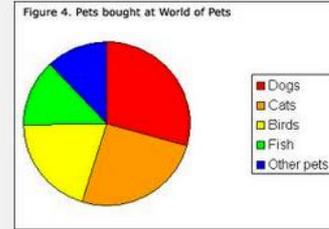
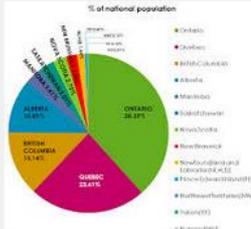
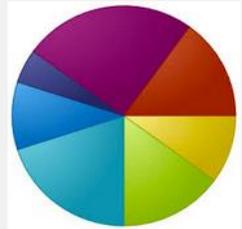
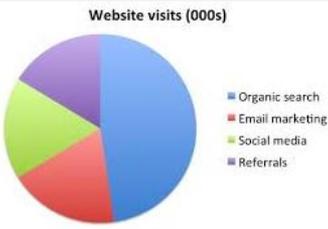
Four NOSQL Categories
arising from the "relational crossroads"



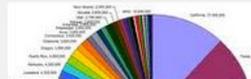
Normalise

Let's talk about graphs

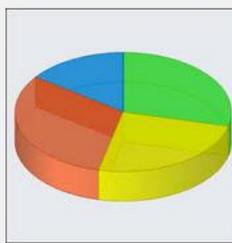
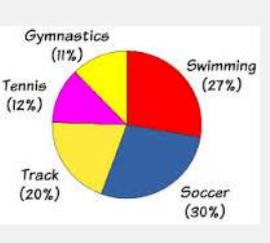
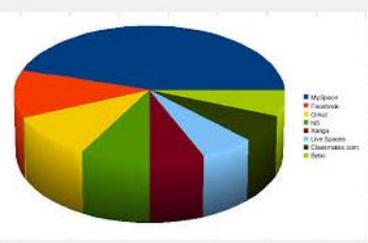
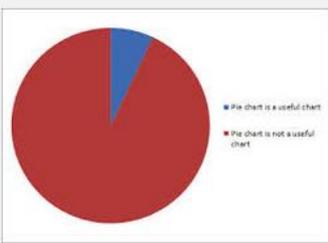
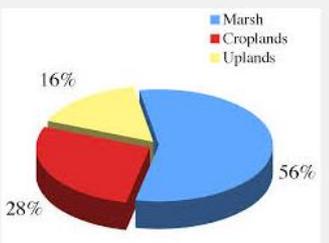
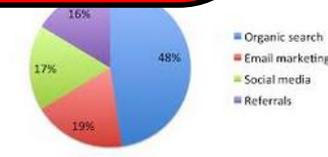
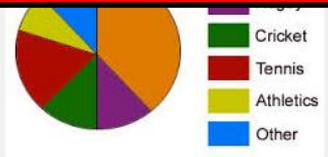
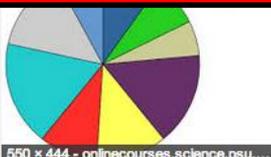
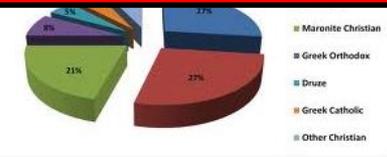




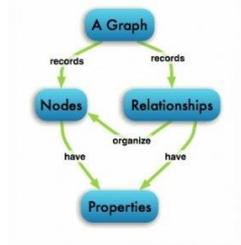
Population of Countries of the European Union in 2007 by percentage



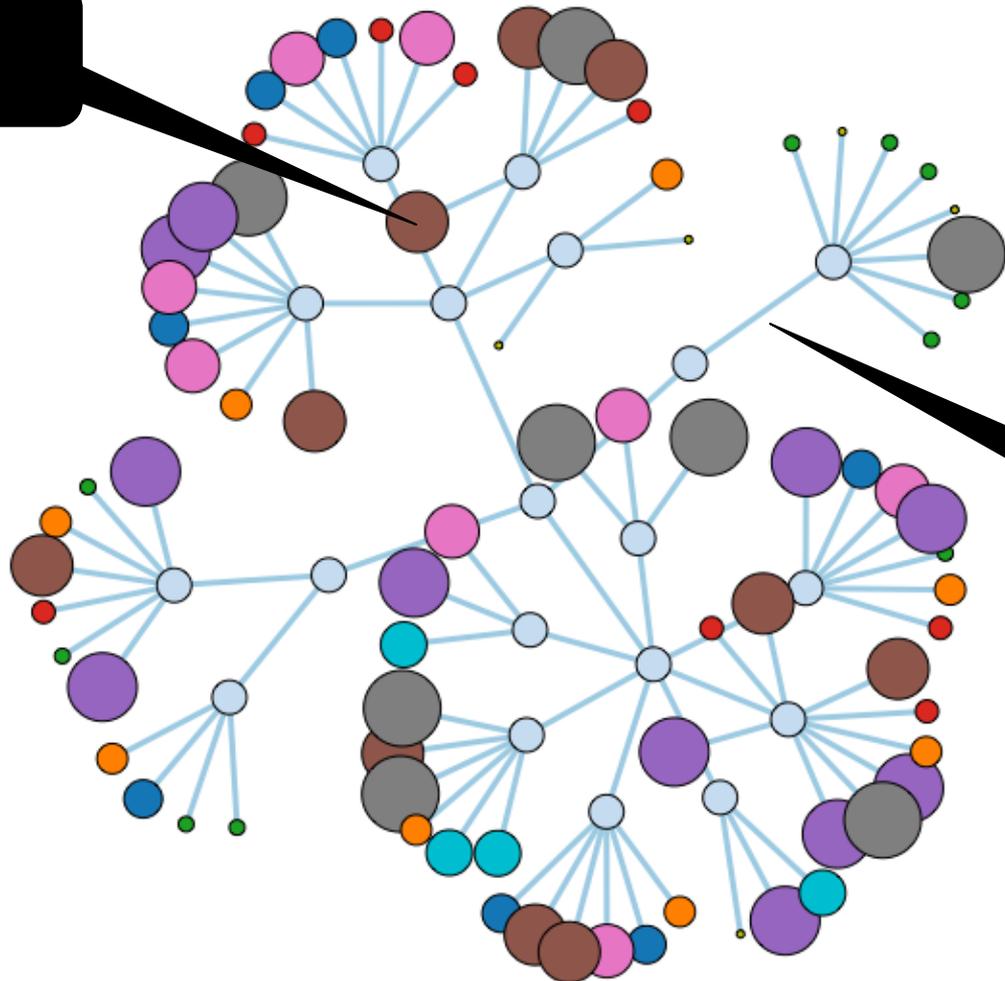
Not the graphs we're going to talk about!



4) Graph Stores

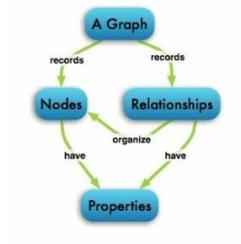


Node



Relationship

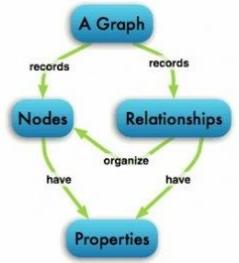
4) Graph Stores



- × “Odd man out” in the non-relational group
- × Designed for **COMPLEX** data – richer data, a lot of expressive power
- × Data model – nodes and edges:
 - Nodes (with properties)
 - Edges are named relationships between nodes (with properties)
- × A query on the graph is also known as traversing the graph
- × Traversing the relationships is very fast
- × Not aggregate-oriented
- × Examples:
 - Neo4j, OrientDB, InfiniteGraph, AllegroGraph
- × Graph theory:
 - People talk about Codd’s relational model being mature because it was proposed in 1969: 43 years old.
 - Euler’s graph theory was proposed in 1736: 276 years old!
- × Semantic Web technologies: RDF, ontologies, triple stores and SPARQL



4) Graph Stores



× Strengths

$complexity = f(size, \textit{variable structure},$

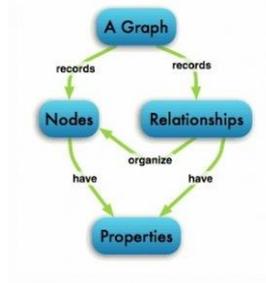
$\textit{connectedness})$

- Powerful data model
 - Fast
 - × For connected data, can be many orders of magnitude faster than RDBMS
 - Good, well-established querying models: Cypher, SPARQL and Gremlin
- ## × Weaknesses:
- Putting some thought into formulating the right data model to make the most of your queries
 - If the data has no / few connections, there is not much benefit in using a graph database

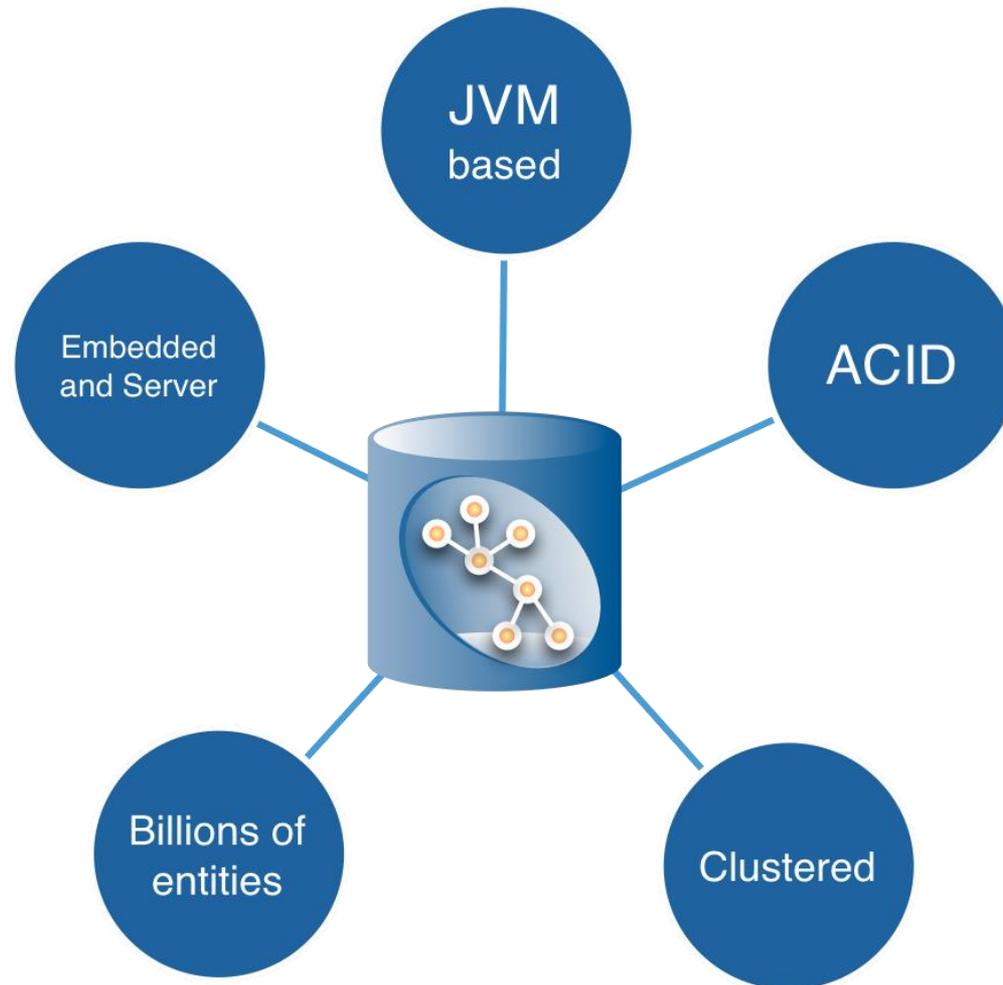
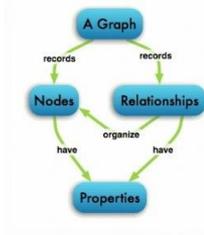
4) Graph Stores

× Suitable Use Cases:

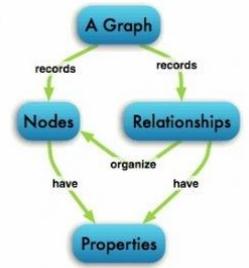
- Recommendation engines
- Business intelligence
- Network impact analysis
- Social computing
- Geospatial
- Systems management
- Web of things / Internet of things
- Genealogy
- Product catalogue
- Access Control
- Life Sciences and scientific computing (especially bioinformatics)
- *Connected* data
- *Hierarchical* data
- Routing, Dispatch, Logistics and Location-Based Services
- Financial services – finance chain, dependencies, risk management, fraud detection etc. For example, if you want to find out how vulnerable a company is to a bit of "bad news" for another company, the directness of the relationship can be a critical calculation. Querying this in several SQL statements takes a lot of code and won't be fast, but a graph store excels at this task.



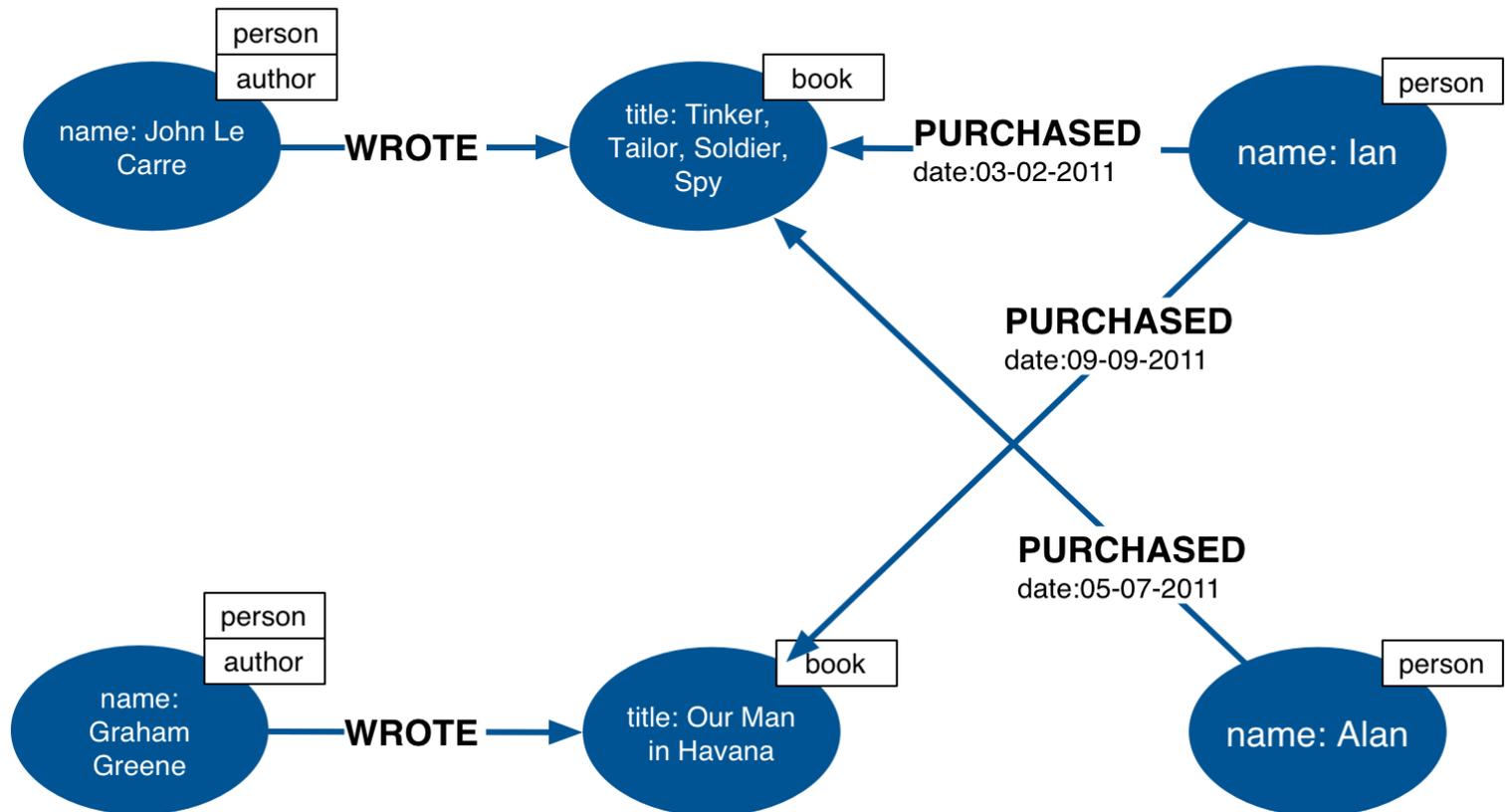
Neo4j is a Graph Database



4) Graph Stores - Neo4j



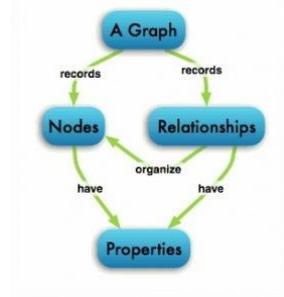
Labeled Property Graph Data Model



4) Graph Stores - Neo4j

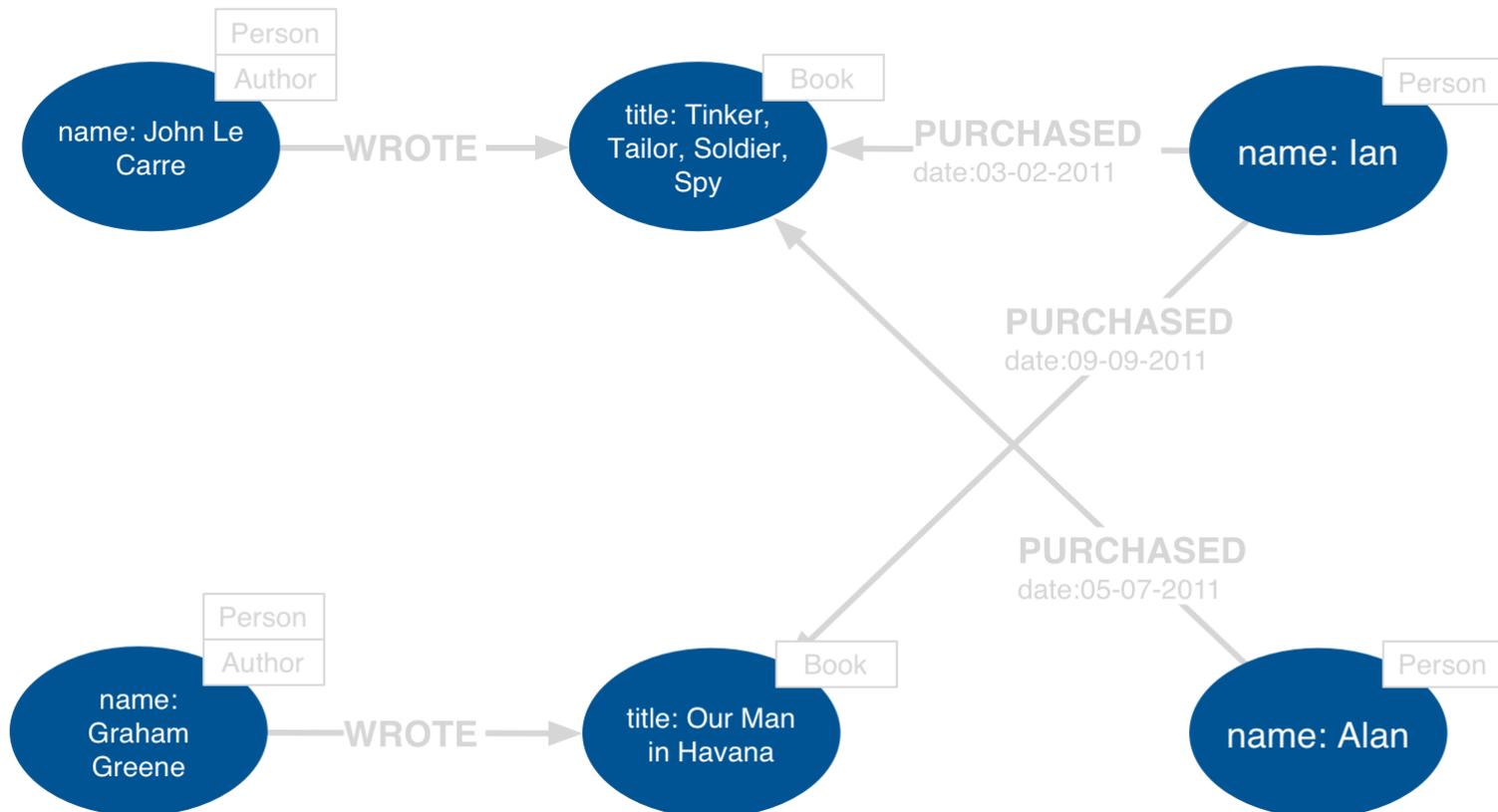
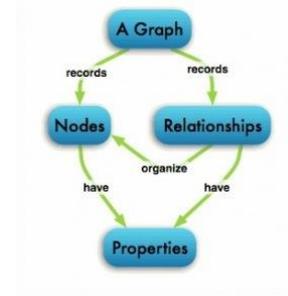
× Four building blocks:

- Nodes
- Relationships
- Properties
- Labels



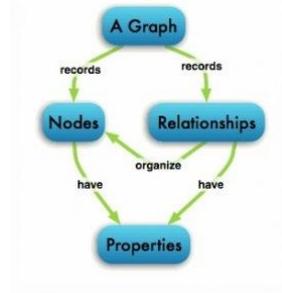
4) Graph Stores - Neo4j

Nodes



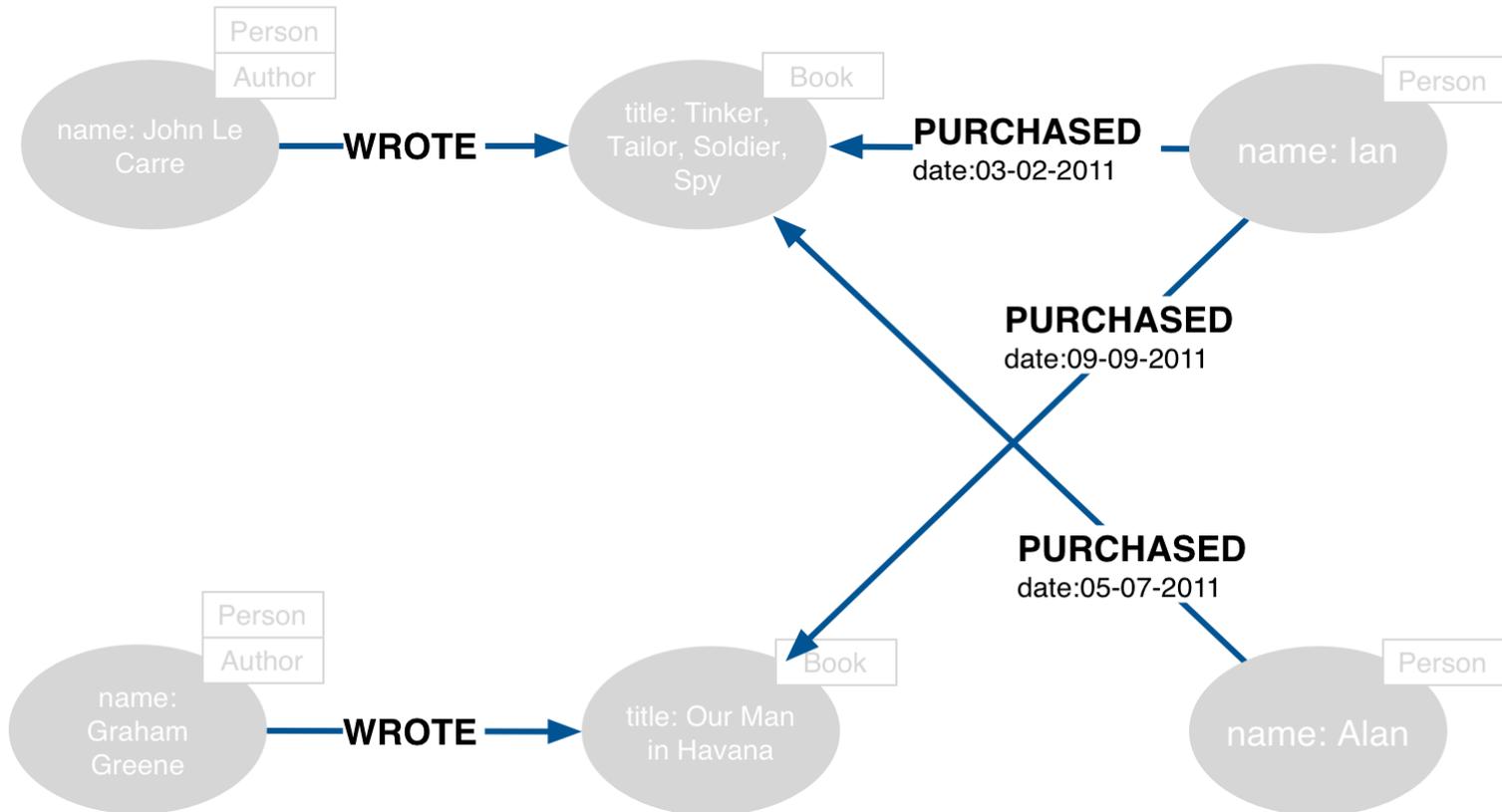
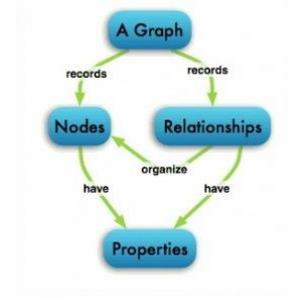
4) Graph Stores - Neo4j

Nodes

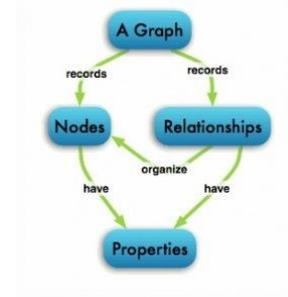


- Used to represent *entities* and *complex value types* in your domain
- Can contain properties
- Nodes of the same type can have different properties

4) Graph Stores - Neo4j Relationships



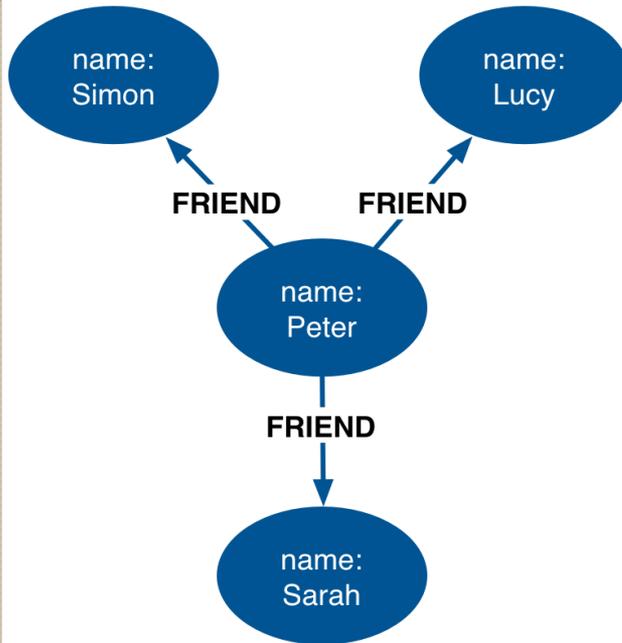
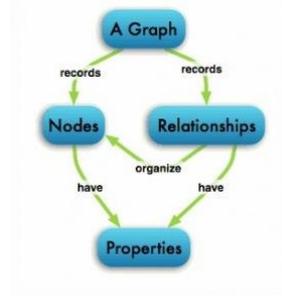
4) Graph Stores - Neo4j Relationships



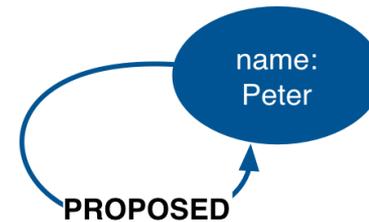
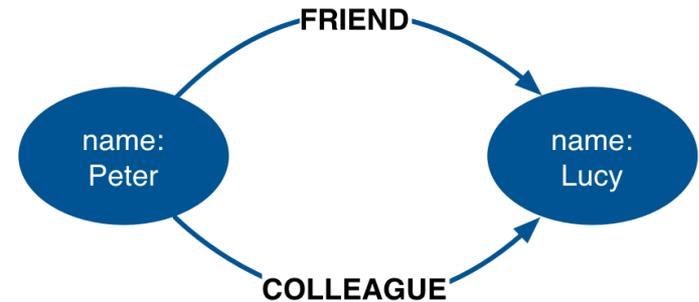
- Every relationship has a *name* and a *direction*
 - Add structure to the graph
 - Provide semantic context for nodes
- Can contain properties
 - Used to represent *quality* or *weight* of relationship or *metadata*
- Every relationship must have a *start node* and *end node*
 - No dangling relationships

4) Graph Stores - Neo4j

Relationships (continued)

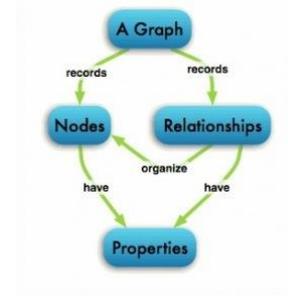


Nodes can have more than one relationship



Self relationships are allowed

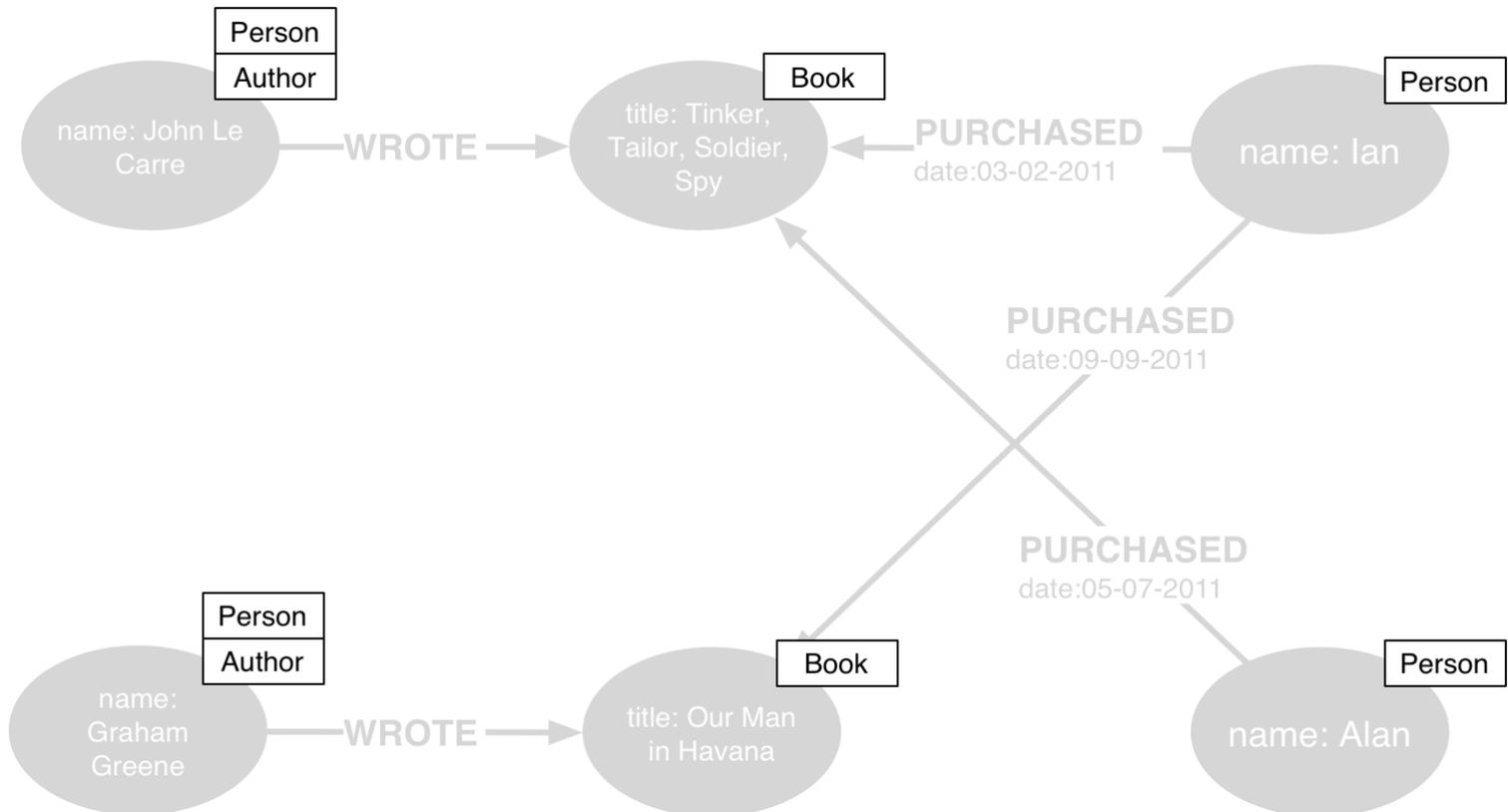
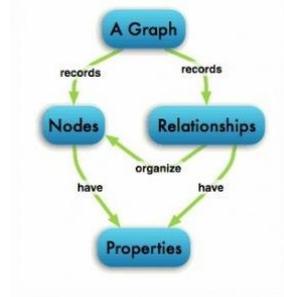
4) Graph Stores - Neo4j Variable Structure



- Relationships are defined with regard to node *instances*, not *classes* of nodes
 - Two nodes representing the same kind of “thing” can be connected in very different ways
 - Allows for structural variation in the domain
 - Contrast with relational schemas, where foreign key relationships apply to all rows in a table
 - No need to use *null* to represent the absence of a connection

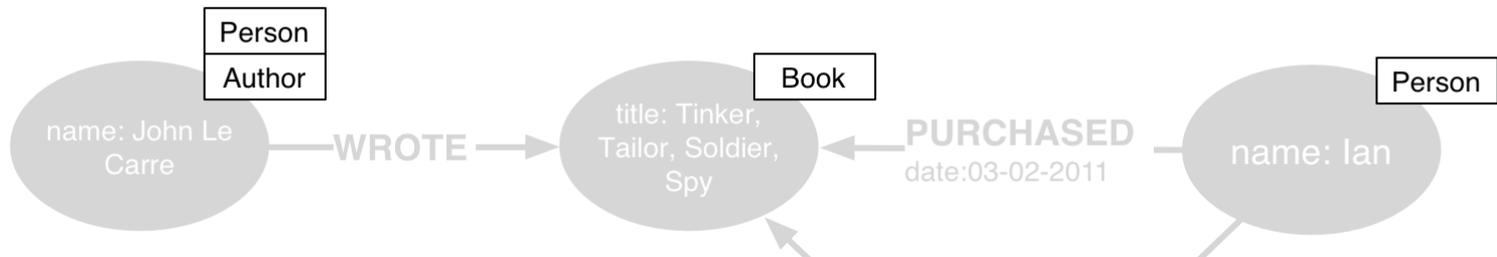
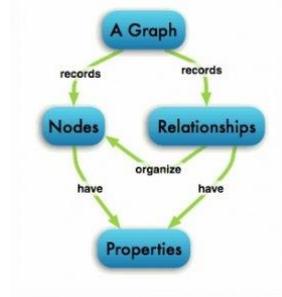
4) Graph Stores - Neo4j

Labels

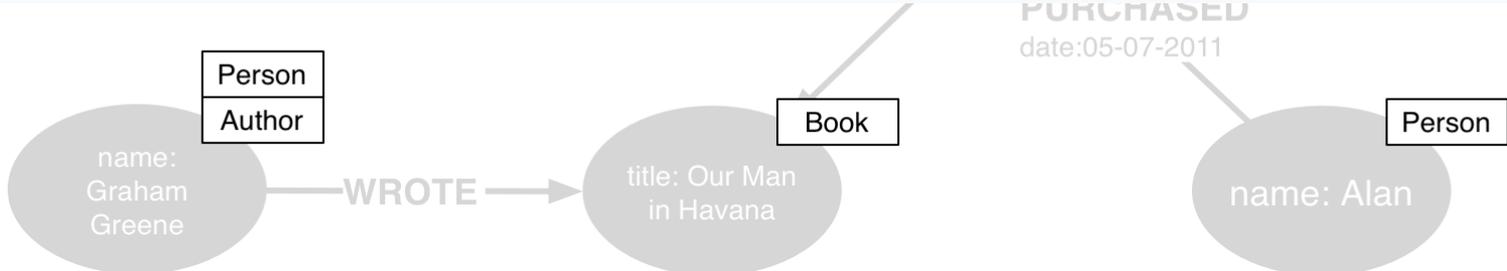


4) Graph Stores - Neo4j

Think Gmail labels

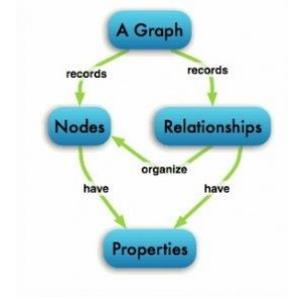


<input type="checkbox"/> ☆ ▢ john.tiger	nodejs [nodejs] node knockout source code - is the source code available for node knockout submission:
<input type="checkbox"/> ☆ ▢ Anders, Xavier, Peter (5)	ruby Leading zeros - Hey! How can I do this: number_pool = (00000..99999).to_a So the first number is 00
<input type="checkbox"/> ☆ ▢ Magnar .. Paul, Stefan (9)	clojure [ANN] Optimus - a Ring middleware for frontend performance optimization. - I just open source
<input type="checkbox"/> ☆ ▢ Alejandro .. Nickolay (3)	haskell [Haskell-cafe] free vs. operational vs. free-operational - Dear Café, I've been reading lately about



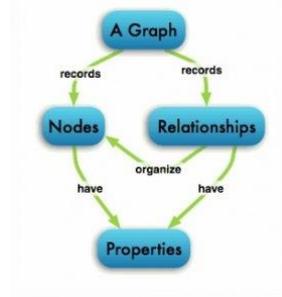
4) Graph Stores - Neo4j

Labels



- Every node can have *zero or more* labels
- Used to represent *roles* (e.g. user, product, company)
 - Group nodes
 - Allow us to associate *indexes* and *constraints* with groups of nodes

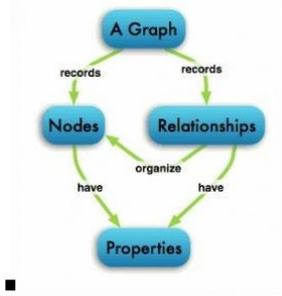
4) Graph Stores - Neo4j



Introducing Cypher

- Declarative graph pattern matching language
- SQL-like syntax
- ASCII art based
- Able to read and mutate the data, as well as perform various aggregate functions such as count and so on

4) Graph Stores - Neo4j

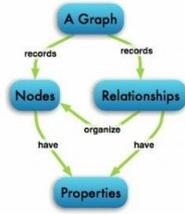


Two nodes, one relationship

```
MATCH (a)-[:ACTED_IN]->(m)
RETURN a.name, m.title;
```



4) Graph Stores – Neo4j



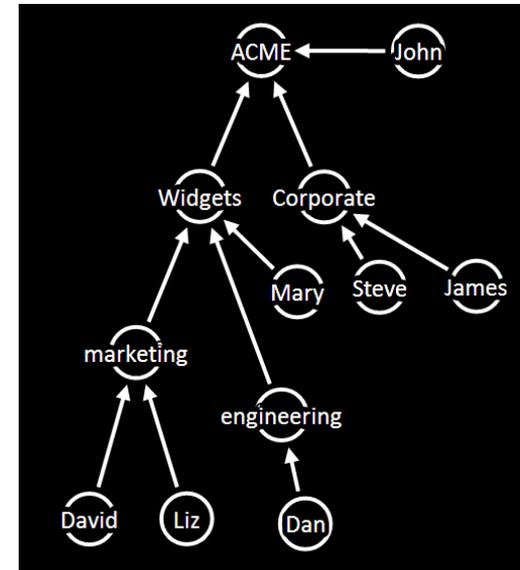
× Graph data (see diagram):

- Person (label) → “WORKS_IN”- Department (label)
- Department → “PARENT”- Department

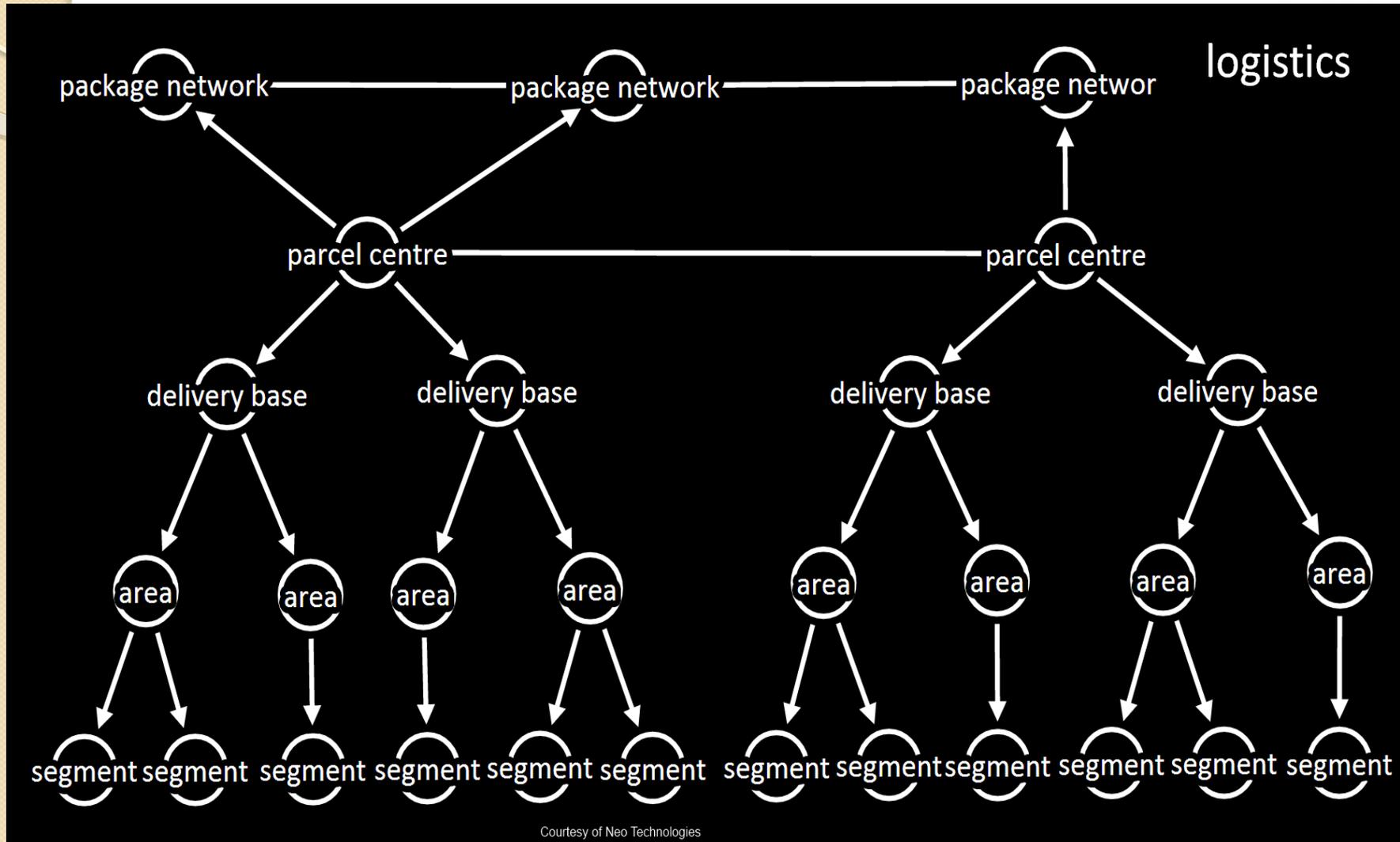
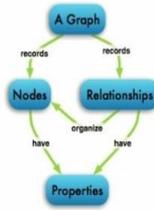
× Cypher query:

```
match (emp:Person) -[:WORKS_IN]->
  (sub_dept:Department)
  -[:PARENT*0..3]->
  (dept:Department {name:"Widgets"})
return emp.name
```

Results: David, Liz, Dan, Mary

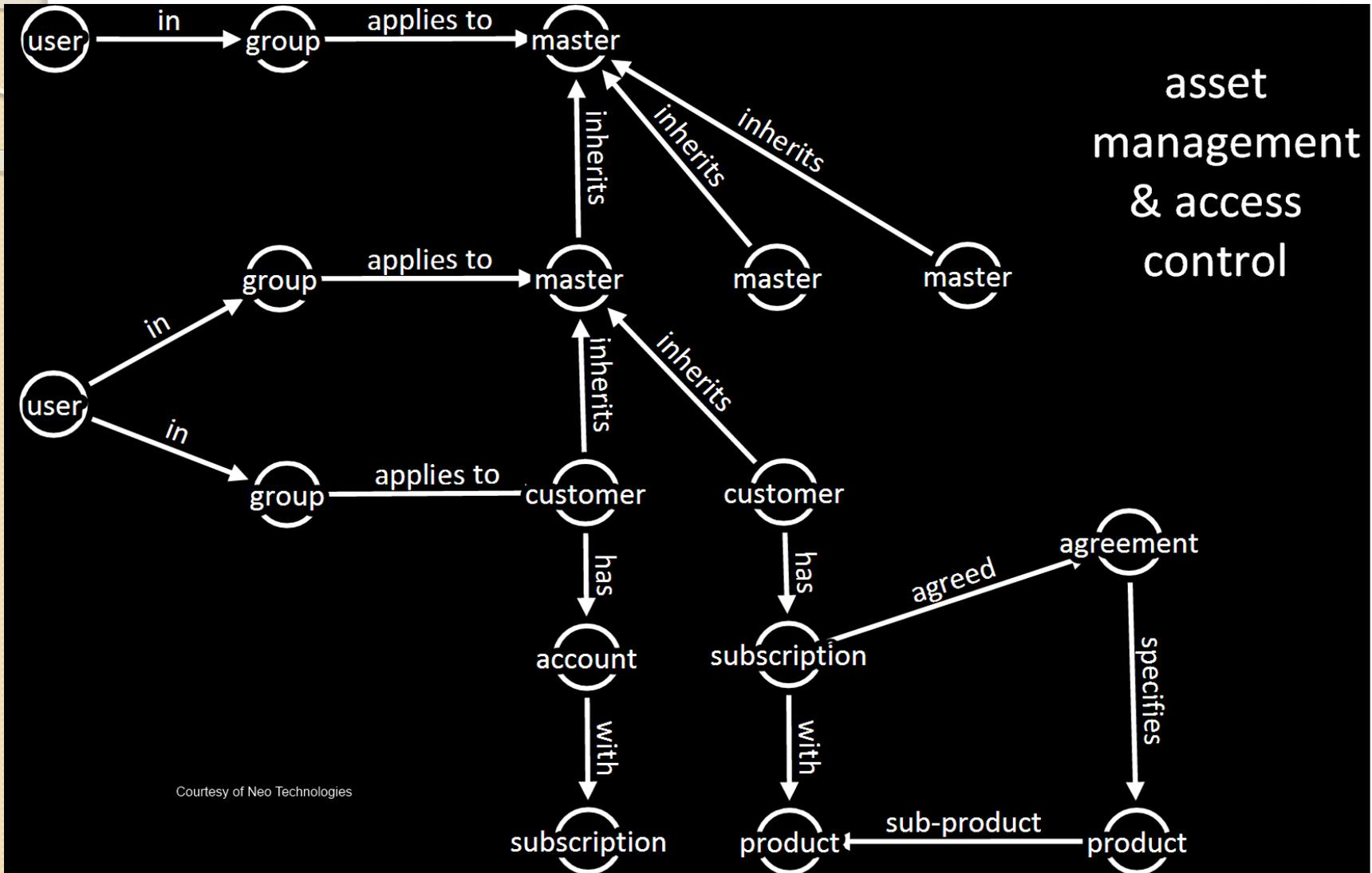
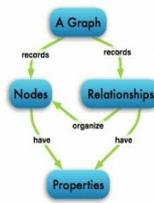


4) Graph Stores – use case



Courtesy of Neo Technologies

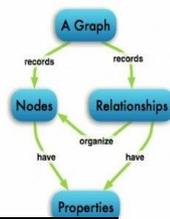
4) Graph Stores – use case



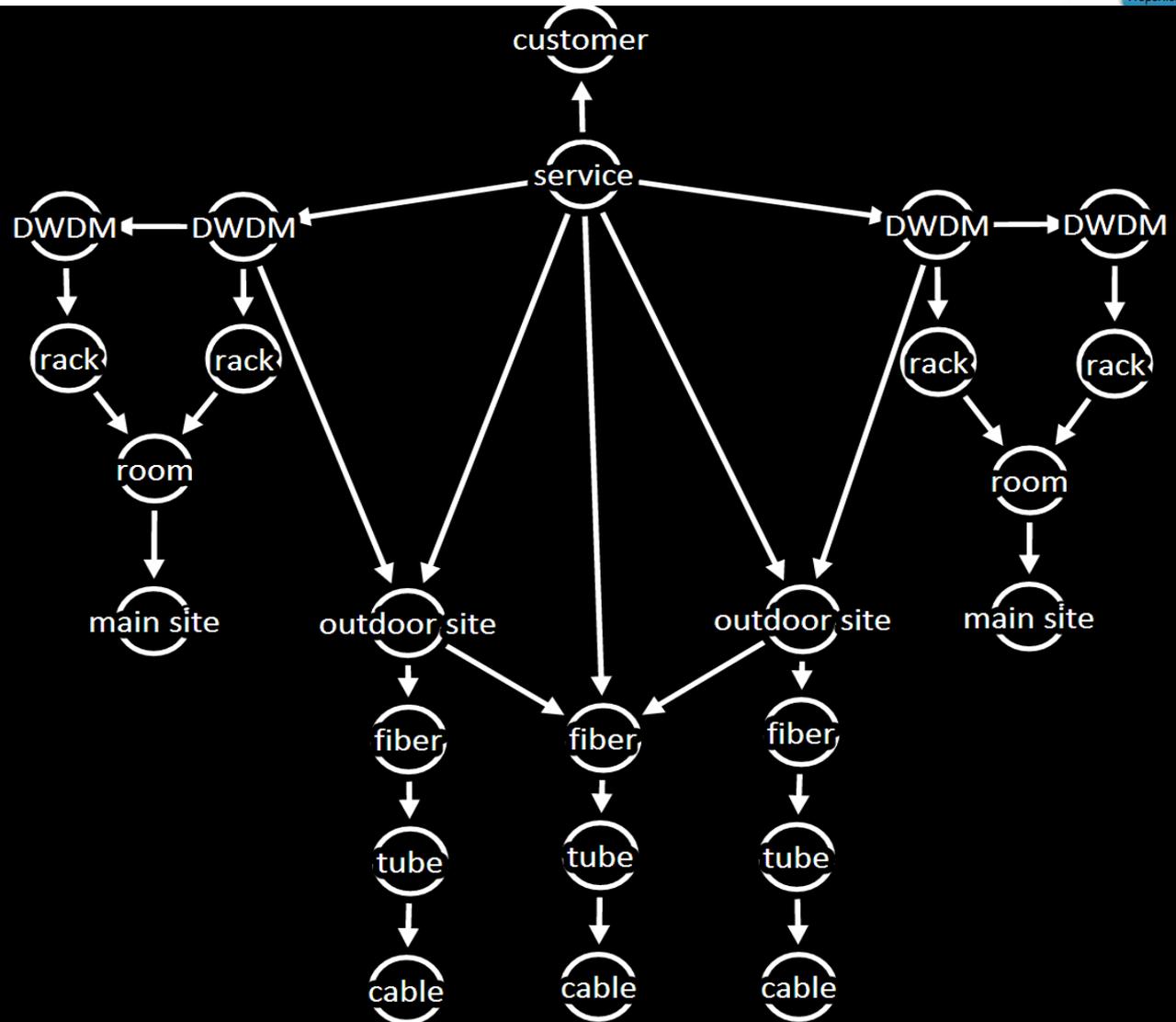
asset management & access control

Courtesy of Neo Technologies

4) Graph Stores – use case

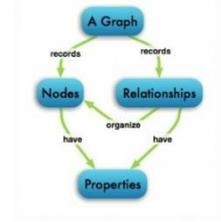


network
impact
analysis



Courtesy of Neo Technologies

4) Graph Stores



- × Less about the volume of data or availability
- × More about how your data is related
- × Densely-connected, variably structured domains**
 - × Lots of join tables? Connectedness**
 - × Lots of sparse tables? Variable structure**
- × Path finding**
- × Deep joins**
- × Use in any case where the **relationship** between the data is just as important as the data itself.
- × Don't use if your data is simple or tabular.
- × More use cases for graphs at <http://neo4j.com/customers/>

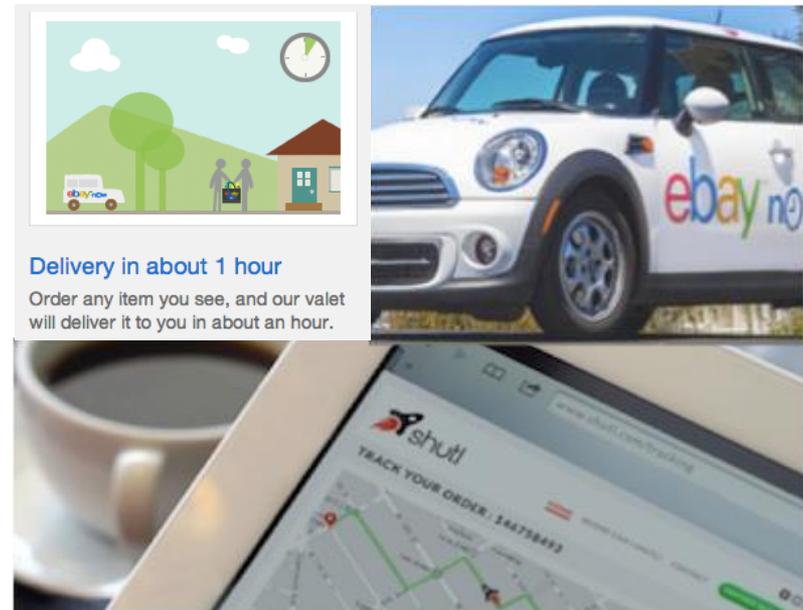
****where a real time response is needed**

Background

- As eBay seeks to expand its global retail presence. Quick & predictable delivery is an important competitive cornerstone
- To counter & upstage Amazon Prime, eBay acquired U.K.-based Shutl to form the core of a new delivery service, launching eBay Now (www.ebay.com/now) prior to Christmas 2013
- Founded in 2009, Shutl was the U.K. Leader in same-day delivery, with 70% of the market

Business problem

- Enable customer-selected delivery inside 90min
- Maintain a large network routes covering many carriers and couriers. Calculate multiple routing operations simultaneously, in real time, across all possible routes
- Scale to enable a variety of services, including same-day delivery, consumer-to-consumer shipping (www.shutl.it) and more predictable delivery times



Delivery in about 1 hour

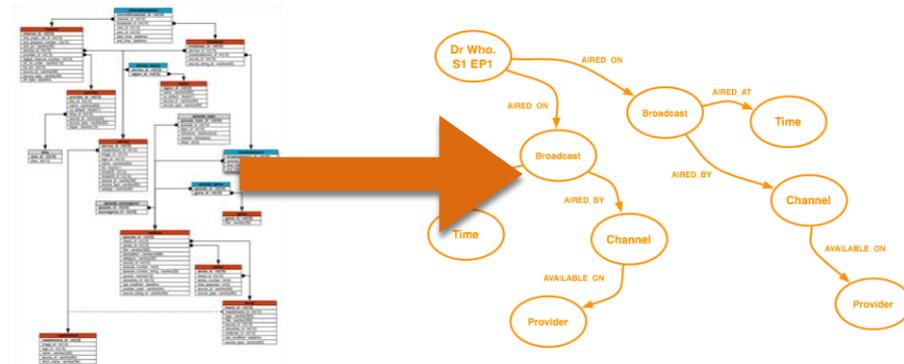
Order any item you see, and our valet will deliver it to you in about an hour.

Solution & Benefits

- Neo4j runs at the heart of the system, calculating all possible routes in real time for every order
- The Neo4j-based solution is **thousands of times faster than the prior MySQL solution**
- Queries require 10-100 times less code, **improving time-to-market & code quality**
- Neo4j makes it possible to add functionality that was previously not possible, and to easily extend the platform over time

Background

- Zeebox is a well-established UK startup that offers second screen applications to end-users, advertisers and broadcasters
- Founded by true media experts, Zeebox aims to reinvent TV since the advent of ... TV.



Business problem

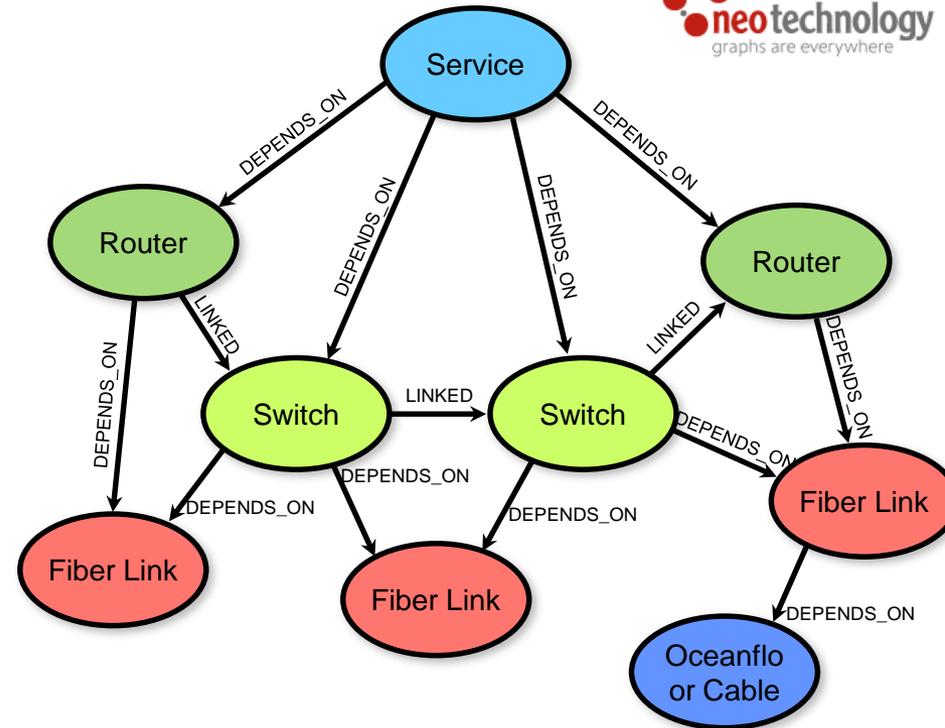
- Data complexity was growing exponentially as more broadcasters and more shows were being added
 - leading to development time increases for applications - a key strategic disadvantage in a fast-moving industry
- Query times on the MySQL based model were starting to explode
 - risk of having worse end-user experience. This was “make or break” with respect to Zeebox’ offering and market position

Solution & Benefits

- Neo4j 2.0 offered a much simpler, natural way to model, implement and query their electronic program guide data
 - leading to faster development cycles
 - no “wedging” of the model into an artificial relational representation
- Future-safe solution: adding more channels/broadcasters/programs does not complicate the model unnecessarily
- Query times went from 80 seconds (MySQL) to 42 milliseconds (neo4j 2.0 traversal)

Background

- Second largest communications company in France
- Part of Vivendi Group, partnering with Vodafone



Business problem

- Infrastructure maintenance took one full week to plan, because of the need to model network impacts
- Needed rapid, automated “what if” analysis to ensure resilience during unplanned network outages
- Identify weaknesses in the network to uncover the need for additional redundancy
- Network information spread across > 30 systems, with daily changes to network infrastructure
- Business needs sometimes changed very rapidly

Solution & Benefits

- Flexible network inventory management system, to support modeling, aggregation & troubleshooting
- Single source of truth (Neo4j) representing the entire network
- Dynamic system loads data from 30+ systems, and allows new applications to access network data
- Modeling efforts greatly reduced because of the near 1:1 mapping between the real world and the graph
- Flexible schema highly adaptable to changing business requirements



Industry: Online Dating
Use case: Social, Recommendations
New York, NY



Background

- AYI, one of the largest dating applications with over 68 million installs, offers an integrated Facebook, iPhone, Android, and Web experience.
- Recently ranked the 36th Fastest Growing Company in North America on Deloitte 's 2012 Technology Fast 500™.



Business problem

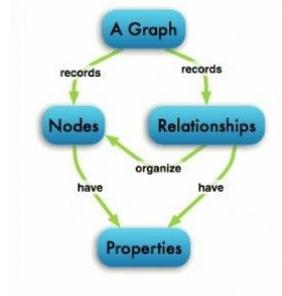
- Performant friends of friends etc. queries (2 and 3 hops out) which become very difficult considering sheer volume of data.
- From their 28 million active users and their friends lists, they had 1.04 billion unique users in the database. Each of these 28 million users has a relationship going from them to their friends which equates to 9.7 billion relationships.

Solution & Benefits

- Real time recommendations serving production apps and website
- Outperformend Solr and MySql on all performance benchmarks for friends of friends
- Enabled the possibility to go even more hops out
- Ability to continously scale and add functionality for a constantly expanding user base

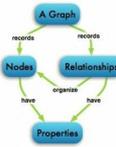
4) Graph Stores - Neo4j

Resources



- Neo4j Manual: <http://neo4j.com/docs/milestone/>
- Example models to look at:
 - Chapter 3 of Graph Databases (book available online at www.graphdatabases.com)
 - Neo4j Manual: <http://neo4j.com/docs/milestone/data-modeling-examples.html>
- Getting started: <http://neo4j.com/developer/get-started/>
- Online training: <http://neo4j.com/graphacademy/>
- Meetups / User group (last Wed of the month) at <http://www.meetup.com/graphdb-london> (free talks and training sessions)

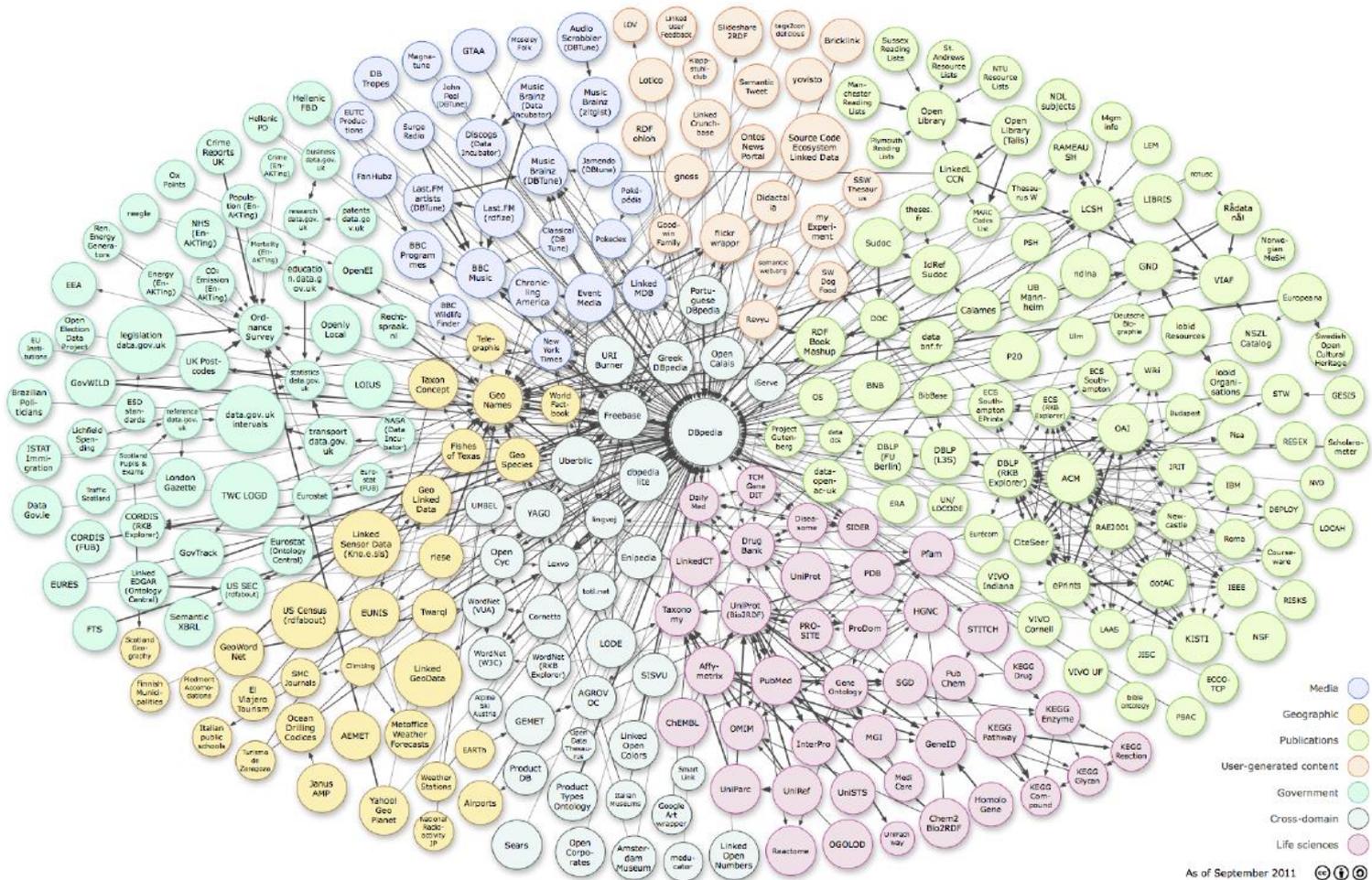
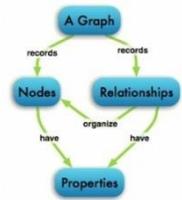
4) Graph Stores – Triple Stores



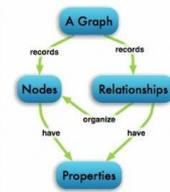
- × Comments from the web: “Semantic Web and RDF is ‘legacy’ / not interesting etc” - **WRONG**
- × Triple stores:
 - The foundation of many Semantic Web systems
 - Encoded in RDF format
 - Each row is a ‘node – link – node’ structure (subject – predicate - object)
 - They also focus on the ability to join graphs together automatically by matching the identifiers of nodes.
 - By merging two graphs from unrelated systems joins can be performed automatically. For example the first graph stores node A links to B and a second graph links B to C, and the union of these graphs shows a relationship of A to C.
- × Examples: Stardog, Virtuoso, Sesame, Jena
- × RDF data is queried via the protocol and query language SPARQL, which incorporates the use of ontologies for inferencing (designed by the W3C RDF Data Access Working Group)

4) Graph Stores – LOD

Linked Data cloud



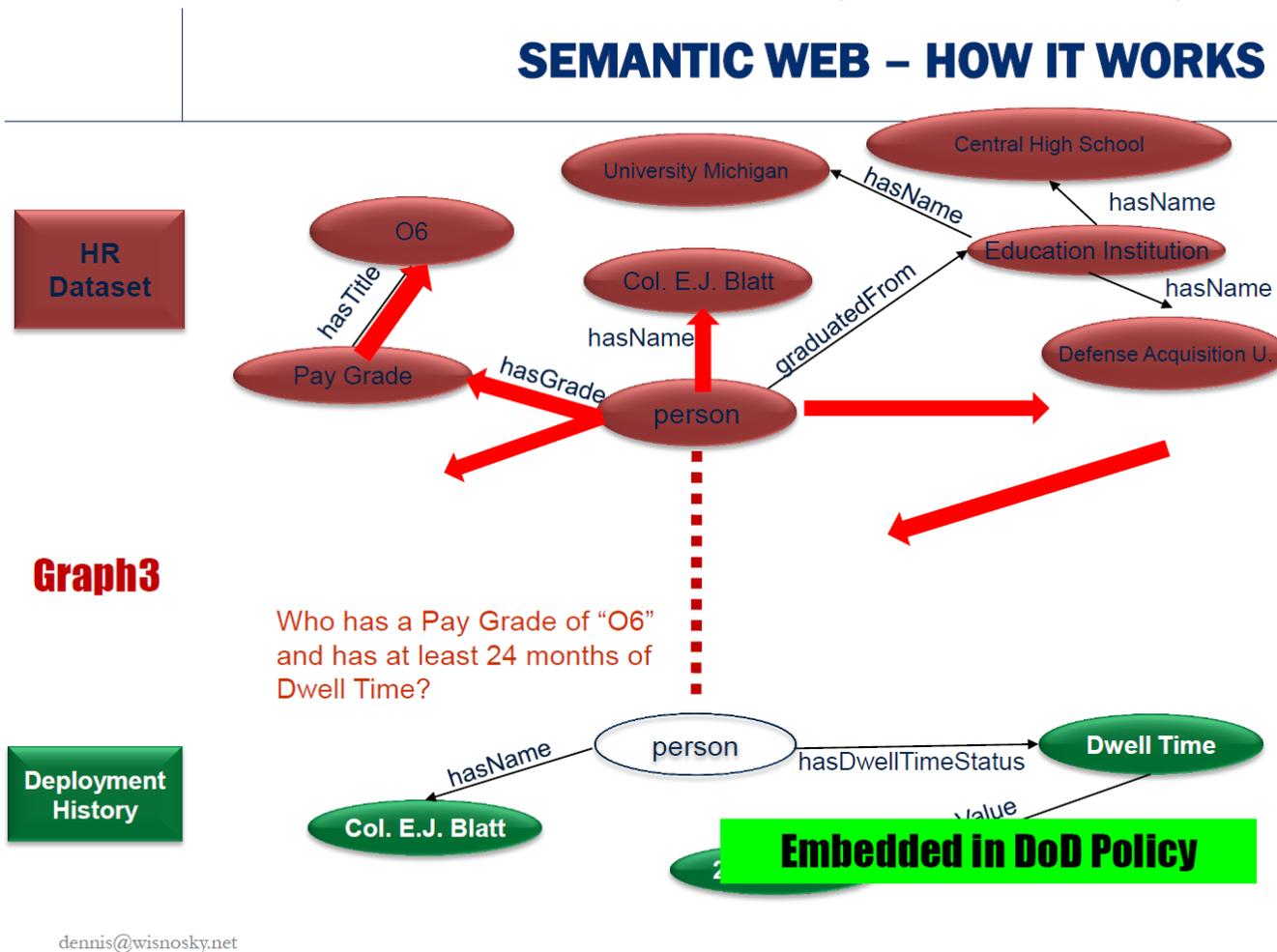
Linking Open Data cloud diagram, by Richard Cyganiak and Anja Jentzsch. <http://lod-cloud.net/>



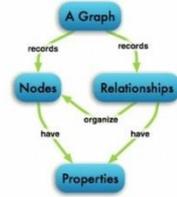
4) Graph Stores – Semantic Web (SW)

US DoD (Department of Defense) is using SW technologies...

SEMANTIC WEB – HOW IT WORKS



4) Graph Stores – SW: why should you care?



× In Life & Health Science –

- Genome projects
- Generating disease models
- New patient treatments
- New drugs

× Financial Services

× New trading strategies

× Intelligence & security industries

- New fraud patterns

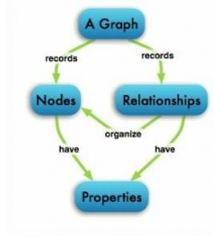
× Space sciences (NASA's Jet Propulsion Laboratory) and clinical trials: provenance of data and experiments

× Netherlands 'actionable intelligence' example:

- Using an ontology to model a wedding in Afghanistan
- Can be used to detect whether an event is really a wedding or a cover for a bombing: red-flag it if a component is missing
- Similar for detecting anomalies in airline passengers to counter terrorism

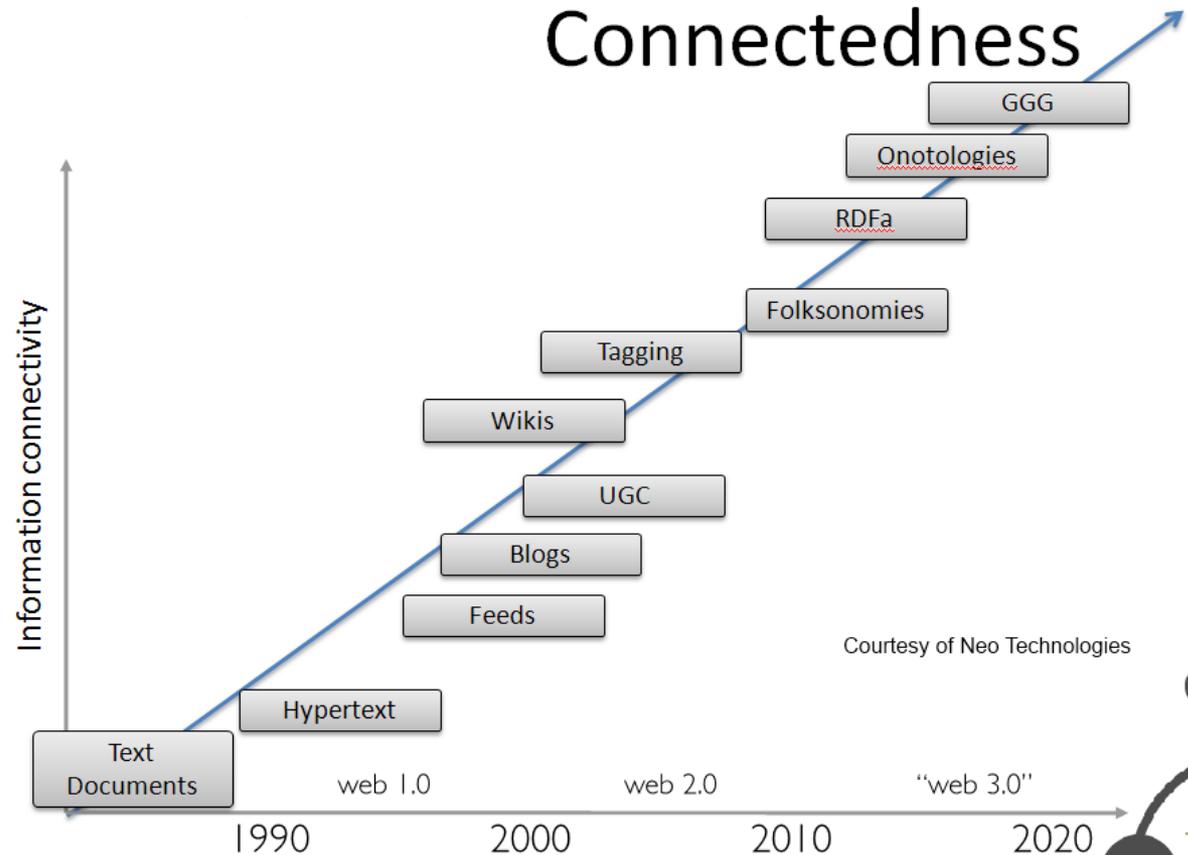
× Google Knowledge Graph

4) Graph Stores – future

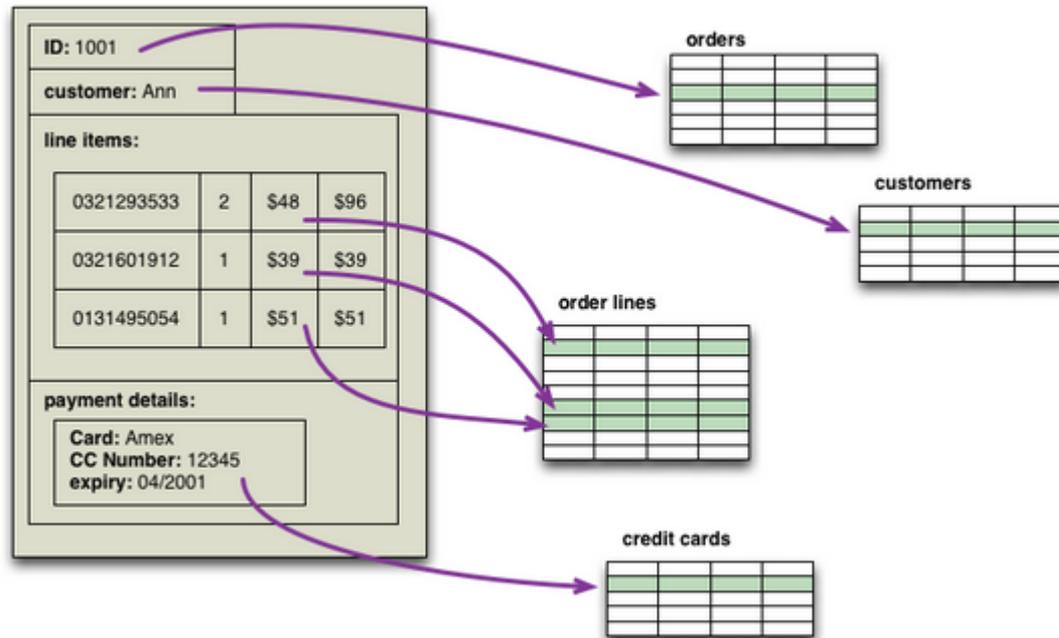


- × Internet → net of computers
- × Word Wide Web → web of documents
- × (GGG) Giant Global Graph → graph of metadata

“I called this graph the Semantic Web, but maybe it should have been Giant Global Graph.” - Tim Berners-Lee - 2007



Aggregate-oriented data stores revisited



“

- × “There is a significant downside - the whole approach works really well when data access is aligned with the aggregates, but what if you want to look at the data in a different way? Order entry naturally stores orders as aggregates, but analyzing product sales cuts across the aggregate structure. The advantage of not using an aggregate structure in the database is that it allows you to slice and dice your data different ways for different audiences.” – Martin Fowler

Aggregates and Polyglot Persistence

- × “This is why aggregate-oriented stores talk so much about MapReduce - which is a programming pattern that's well suited to running on clusters. MapReduce jobs can reorganize the data into different groups for different readers - what many people refer to as materialized views. But it's more work to do this than using the relational model.”
- × “This is part of the argument for [Polyglot Persistence](#) - use aggregate-oriented databases when you are manipulating clear aggregates (especially if you are running on a cluster) and use relational databases (or a graph database) when you want to manipulate that data in different ways.”
-- (Martin Fowler)

Polyglot Persistence

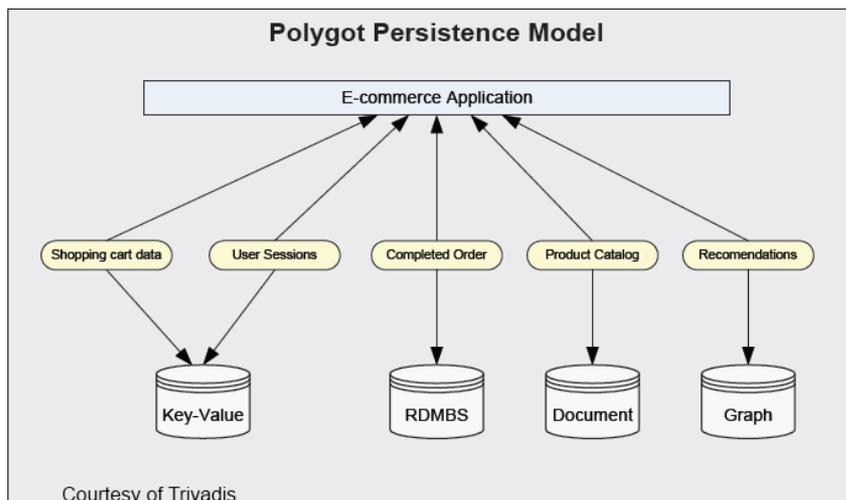
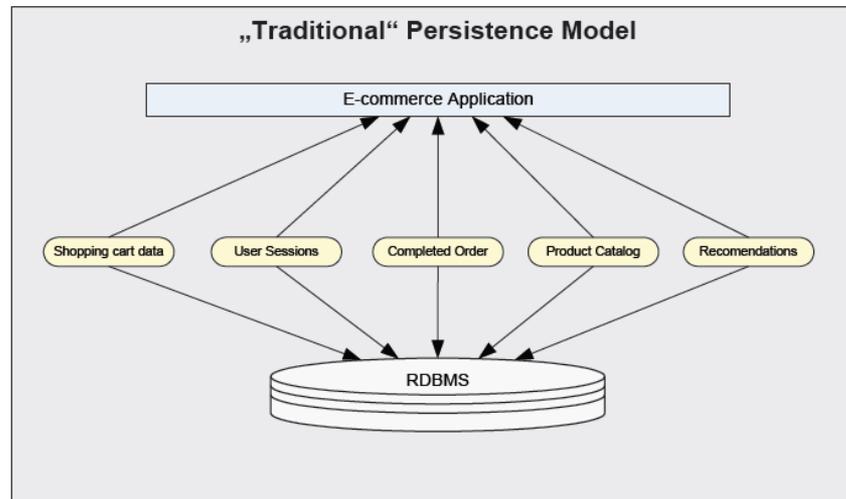
× Polyglot Programming and Persistence:

- Coined by Neal Ford in 2006
- Expresses the idea that applications should be written in a mix of languages to take advantage of the fact that different languages are suitable for tackling different problems: i.e. defines a hybrid approach to data storage
- Complex applications combine different types of problems, so picking the right language for the job may be more productive than trying to fit all aspects into a single language --- (recall the “hammer and nails” idea!)
- So, **we want to avoid pushing a square peg into a round hole....**
- Use multiple data storage technologies; particularly when there is no need for same properties of availability, consistency or backup requirements
- Selection based on the way data is being used by individual applications
- Can occur both over the enterprise as well as within a single application

× <http://martinfowler.com/bliki/PolyglotPersistence.html>

Polyglot Persistence

Using a variety of data technologies for different kinds of data, applications and use-cases



Polyglot Persistence

- × Polyglot persistence usually comes about in one of the following ways:
 - We have an existing relational system, and find that we want to query the data along 'connectedness' lines (as an example). We find that this gets harder to code over time, and that it becomes less and less performant in real time. This is the time to look at replacing those parts of the system with a graph store (for this example). It is more cost-effective to just move parts of the system over, as opposed to everything – and far less risky for the business.
 - We are building a new system, and realise that the data fits into distinct categories, for example: huge volumes of simple time-series data that don't need to be inter-related, giant multimedia files, and then something closely knit and highly interconnected. As the data volumes grow and SLAs become more rigorous, it starts to make sense to store data in a place that's optimised for that type of data. In this example, you might use something like Cassandra and Neo4j.
- × An example is Telenor, one of the world's 10 largest telecommunications firms, replaced part of a Sybase application with Neo4j for hierarchical queries that needed to run very fast, but kept much of their existing database around.

NOSQL Challenges

× Immaturity

- NOSQL tools are still young, full of the rough edges that new tools have, although this is changing
- Not much experience: we don't know how to use them well
- No patterns and best practices exist yet

× Organisational Change

- How will the different data groups in an enterprise react to this new technology
- IT ecosystem
- Not as many standards as relational (yet)
- Security not as comprehensive as relational (yet)

Choosing – Design Mindset

- × Don't think of RDBMS vs. NOSQL, think of the “right tool for the right job”
 - Flexibility: Do you benefit from having schema-less design? Do you need the flexibility of exploratory queries?
 - How connected is your data
 - How complex is your data
 - Performance: How much data? How many writes? How many reads? Connectedness of data important?
 - Transactions? Can you live without them?
 - Do I need complex queries and sorts?
 - Am I mapping a complex relationship tree?
 - What is more important: **C**onsistency or **A**vailability?
 - Am I updating records, or just inserting new records?
 - Could I wait a few seconds if a primary node failed?

“The whole point of seeking alternatives [to RDBMS systems] is that you need to solve a problem that relational databases are a bad fit for.” – Eric Evans, Rackspace

Using NOSQL - Strategies

- × Good to use the "Goldilocks Pilot Project Strategy"
 - Not too big, not too small, just the right size
 - Duration
 - Sponsorship
 - Importance
 - Skills
 - Mentorship
- × Remember that a flexible schema does **not** mean there is an excuse for a badly-designed application!
- × When using NOSQL, and when one doesn't know about the schema beforehand, use a "**schema strategy**" to ensure you adhere to a good design throughout

Data Management vs Data Analytics



Data Management – the actual
'body' of the cake

Data Analytics – the icing
on the cake; cannot exist
in isolation

Data Analytics – a few uses...

- × eCommerce optimisation (which basically means "how do I get more ££ from the customers")
- × Targeted advertisements
- × Financial services, such as risk-modelling
- × Detecting anomalies
- × Automated metadata generation – generating graphs. Graph distance algorithms to determine how related two entities are, and updated in real-time
- × Recommendation systems:
 - Recently heard a talk about Amazon's recommendation system and how 'flat' it was - the fact that, within one household, there may be multiple user types. Thus, it would be cleverer if Amazon could distinguish between these, so the parents, browsing at night, do not get recommendations about colouring books because the kids looked at those during the day...
- × System monitoring. An example would be detecting whether there are emerging problems with an oil rig

Data Analytics

× Modes:

- **Real-time** – transactional processing; showing a ‘fairly suitable’ advert to a customer browsing your web site.
- **Offline** – batch processing; improvements and refinements to the model over time, so that in the future, you can show the ‘best’ advert to a customer browsing your web site

Data Analytics – Hadoop

- × Synonymous with NOSQL
- × Apache – open source
- × A combination of:
 - A distributed file manager – HDFS (Hadoop Distributed File System)
 - A parallel programming framework known as MapReduce (original paper from Google)
- × Usually runs on a cluster of commodity nodes
- × **Batch-oriented**: not suitable for ad-hoc querying!
- × Integration with many data stores – and not only the NOSQL ones...
- × Two thoughts regarding Hadoop:
 - “I will spread your data over many servers and keep it safe”
 - “I will facilitate a new idea that you should send the work to the data and not the other way around”



Data Analytics - Hadoop

- × Great for “boil the ocean”-type processing *and* analytics:
 - **Big Data Analytics:**
 - × Running queries over large semi-structured datasets
 - × Makes filtering and aggregation-type jobs very easy
 - **Big Data Processing:**
 - × Efficient and effective way to write data pipelines
 - × Easy way to parallelise computationally complex queries
 - × Scales nicely with the amount of data and cluster size
- × Not great for the following:
 - Real-time analytics or processing
 - × Even small queries take time
 - × Can't build into real-time data flows
 - Algorithms which are difficult to parallelise
 - × Almost anything can be expressed in a number of MR steps
 - × Almost always, MR is SUB-OPTIMAL! This is the case even if it is easier to abstract



Data Analytics – Hadoop tooling

- × Writing MapReduce jobs in Hadoop natively is very tricky
- × Instead, it is recommended to use tools to do this for you:

- **Hive**

- × Provides a SQL interface to Hadoop's MapReduce
- × Thus, MapReduce jobs are written in a SQL-like language



- **Pig**

- × Language is also very simple and high-level
- × UDFs can be written in Java, Python or Javascript – useful

- Other tools: Cascading, Crunch etc..



Data Analytics – Machine Learning

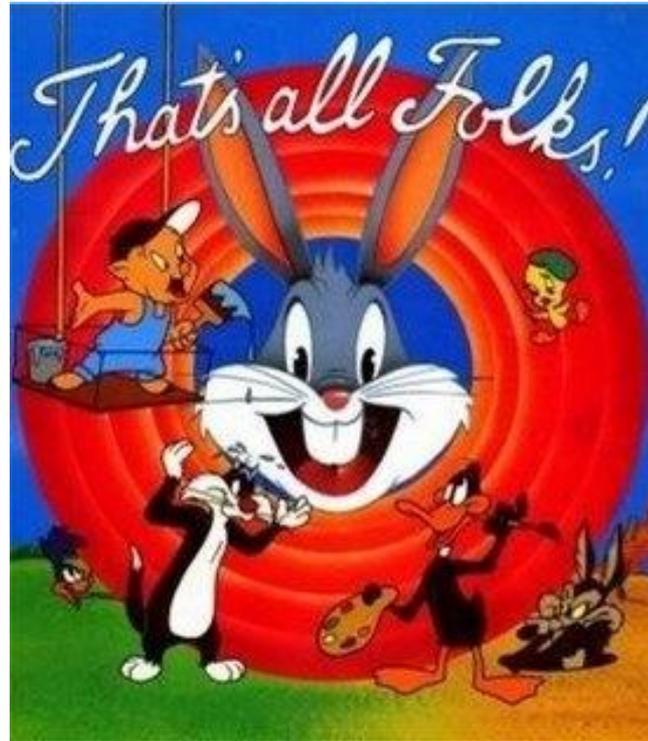
× **Apache Mahout:**

- A library of ML (Machine Learning) algorithms

× **Uses:**

- Log analysis
- Data warehouses
- NLP (Natural Language Processing)
- Search indexing (creation of indices offline)
- Machine learning





Thank you for your attention!