# Graph Analytics

Hannes Voigt

# Bio

## CURRENT POSITION

- Postdoc at Database System Group, Technische Universität Dresden

## EDUCATION

- Ph.D. in 2014
- Master in 2008

## EXPERIENCE

- Visiting scholar at SAP Labs, Palo Alto for one year in 2010
- Visiting scholar at University Waterloo, for 4 months in 2007

## INTERESTS AND ACTIVITIES

- Graph Data Management and Data Science
- LDBC Graph Query Language Standardization Task Force
- Collaboration with SAP Hana Graph Team

# Dresden Database Systems Group

VISIT: HTTPS://WWWDB.INF.TU-DRESDEN.DE/

# Trends in Data Management

# Everything is Data



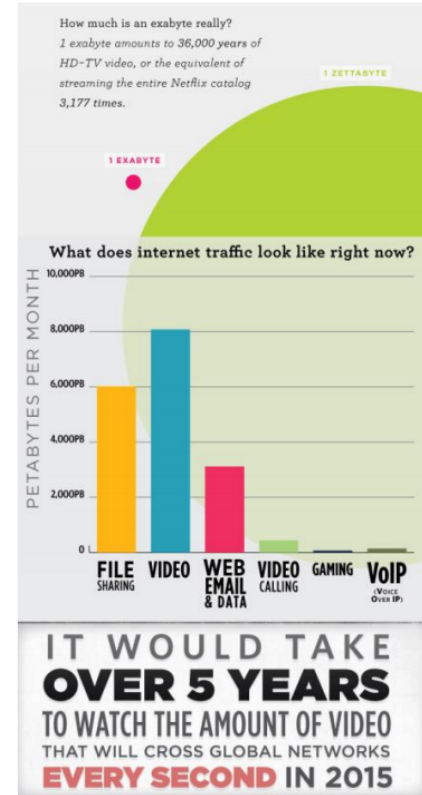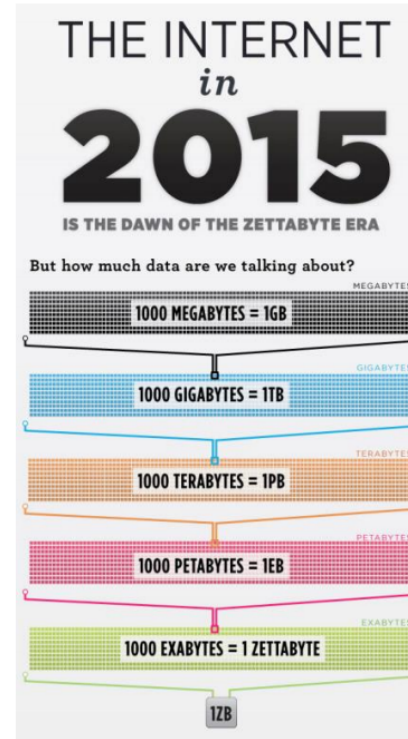[http://blogacronis.com/posts/data-everything-8-noble-truths]

# The Zettabyte Age

## Cisco Visual Networking Index

### The internet in 2020

- 26.3 billion networked devices
  - Up from 16.3 billion in 2015
  - 44% of all networked devices will be mobile-connected
- 25.1 GB average traffic per capita per month
  - Up from 9.9 GB in 2015
- 2.3 Zettabytes annual IP-Traffic
  - up from 870.3 Exabytes annual IP-Traffic in 2015
  - One zettabyte = stack of books from Earth to Pluto 20 times

# The End of Science

> The quest for knowledge used to begin with grand theories. Now it begins with massive amounts of data.
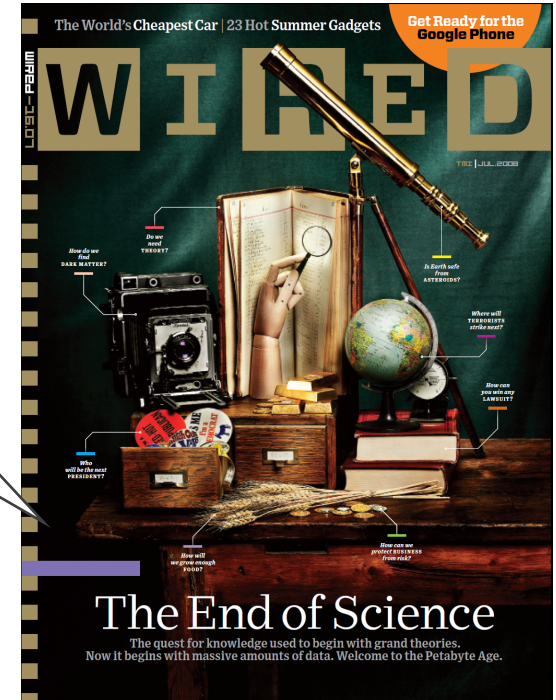> Zetta
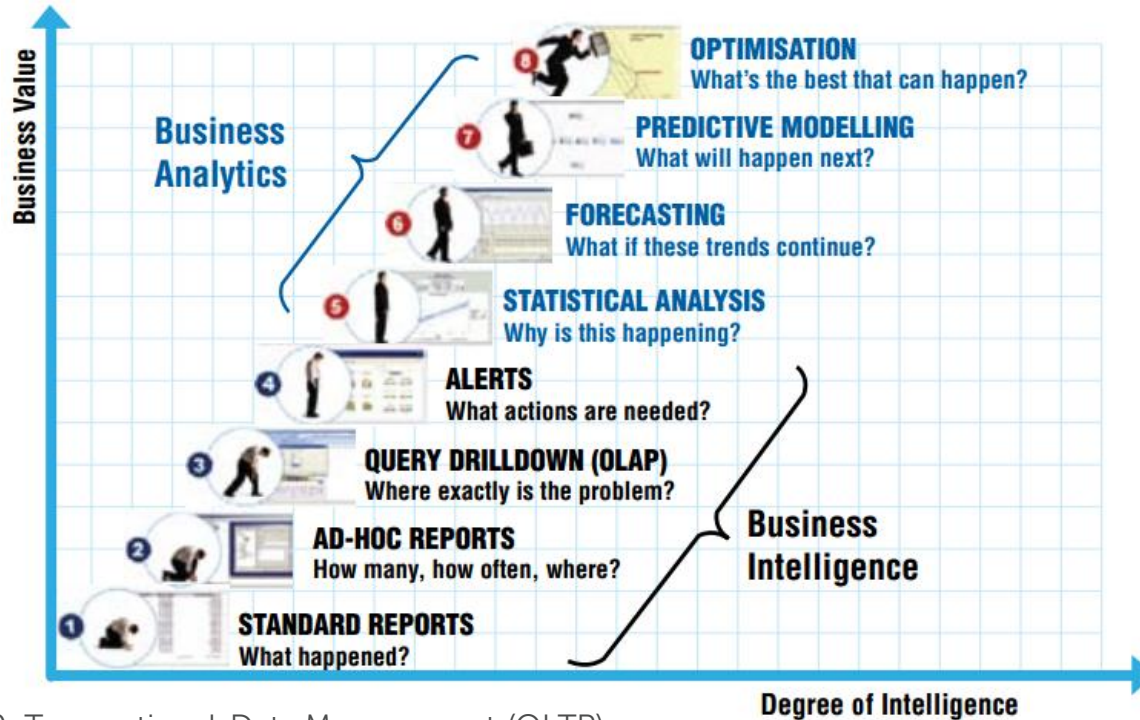> Welcome to the ~~Peta~~byte Age.

## New Realities

- Everything is digital data
- Rise of data-driven culture
- High-performant data analytics
- Exploit sophisticated statistical methods

## How do we structure/implement/live with this trend?

# Level of Analytics

0. Transactional Data Management (OLTP)

# Focus of Interest

## PROPERTIES OF ENTITIES

- Captured/measured values
- What are the sales figures/temperatures/etc.?
- Multidimensional data/time series/matrixes

## CONNECTIONS BETWEEN ENTITIES

- Network structure
- What do the friends of your customers buy?
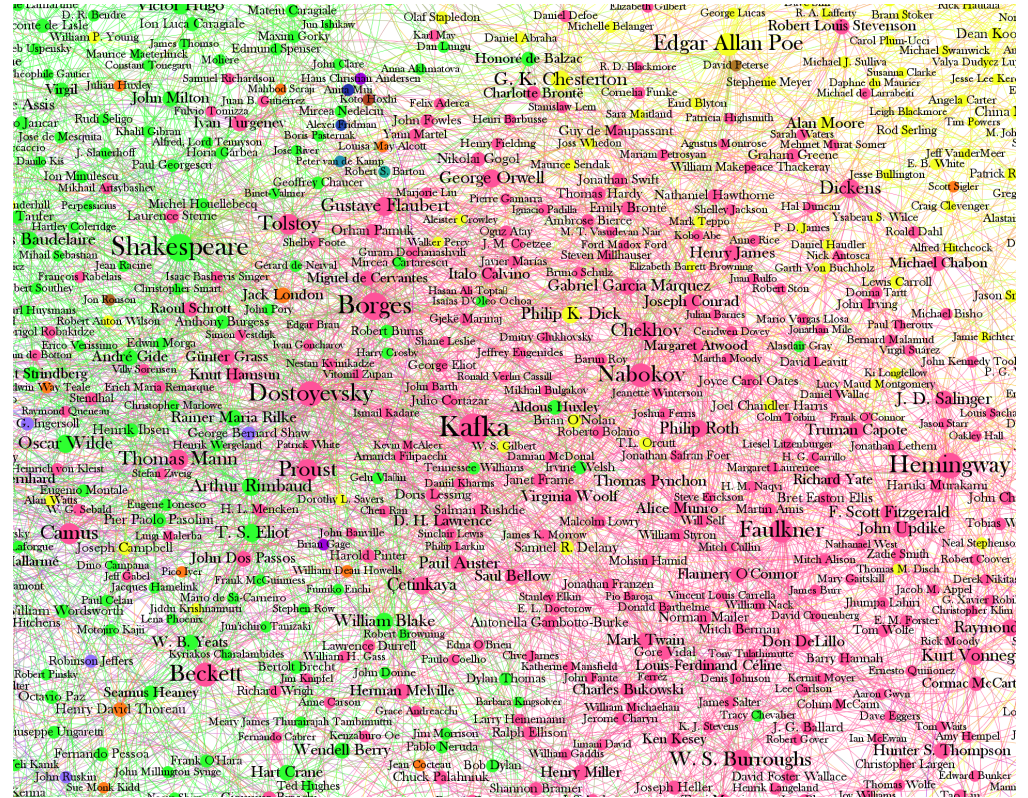- Graph data

# Example: Centrality Measures

## QUESTION: WHO ARE THE KEY PLAYERS IN A GRAPH

- Most social contacts (vaccination schedules)
- Most influential thinkers/papers (reading lists)
- Most important website (web search)
- Most important distributers (supply network)
- etc.

- Can we measure that?

## YES! WITH CENTRALITY MEASURES!

- Centrality measures identify the most important vertices within a graph



[http://brendangriffen.com/blog/gow-influential-thinkers]

# Example: Supply Chain Management

Backward Tracability (Level 3 Analytics)
• Find the parts causing the problem

# Example: Supply Chain Management



Forward Tracability (Level 4 Analytics)
• Alert other customer potentially affected by the problem

# Example: Supply Chain Management



Inventory

Response Times

When to send?
How much to send?
From where to where?

Suppliers | Manufacturers | Distribution Centers | Demand Markets

Domestic Manufacturer

Supply Chain Optimization (Level 8 Analytics)
- Customer A: 10-30% reduction in inventory
- Customer B: 8% reduction in transportation costs

[http://www.logicblox.com/solutions/supply-chain-optimization/]

13

# Business Processes

## BUSINESS PROCESSES ARE ESSENTIALLY GRAPHS

- Who does what in relationship with whom …



[http://integrella.com/services/solutions/business-optimization-bpm/]

## THE WHOLE ANALYTICS STACK DESIRED

- From tracking the state of processes (level 0)



- To optimizing the processes (level 8)



[https://xkcd.com/399/]

14

# Level of Analytics

[http://www.rosebt.com/blog/eight-levels-of-analytics-business-intelligence-to-business-analytics]



**Business Value**

**Business Analytics**

8 **OPTIMISATION** — What's the best that can happen?

7 **PREDICTIVE MODELLING** — What will happen next?

6 **FORECASTING** — What if these trends continue?

5 **STATISTICAL ANALYSIS** — Why is this happening?

4 **ALERTS** — What actions are needed?

3 **QUERY DRILLDOWN (OLAP)** — Where exactly is the problem?

2 **AD-HOC REPORTS** — How many, how often, where?

1 **STANDARD REPORTS** — What happened?

**Business Intelligence**

**Degree of Intelligence**

Graph Data Management Applications

0. Transactional Data Management (OLTP)

# Graph Data Model

# Graph Building Blocks

## NODES (DOTS)

- Like an entity in ER
- Exist on their own
- Have object identity

## EDGES (LINES)

- Like a relationship in ER
- Exist only between nodes
- Identity depends on the nodes they connect

# Graph Building Blocks

## VERTICES & EDGES

- $G = (V, E)$ with
  $E \subseteq \{e \mid e \in \mathcal{P}(V) \wedge |e| = 2\}$

- Vertices have identity
- Edge depend on vertices

## VERTEX LABELS

- $G = (V, E, L_V, f_V)$ with
  $f_V : V \rightarrow L_V$ (or $f_V : V \rightarrow \mathcal{P}(L_V)$)

Film
Actor
Actor
Film

- Label are not unique[1]

## VERTEX PROPERTIES

- Vertices have set of key-value pairs

name = Stolen 3
screen = 2015
name = Brain Pit
eye color = brown
name = Roy Winslet
born = 1975
name = Sequle 4
oscars = 2

## DIRECTIONALITY

- $E \subseteq V \times V$

## EDGE LABELS (OR WEIGHTS)

- $G = (V, E, L_E, f_E)$ with
  $f_E : V \rightarrow L_E$

act in
act in
friends
act in

## EDGE PROPERTIES

- Edges have set of key-value pairs

year = 2004
character = Oncle Bob
since = 1985
year = 2005
salary = 2.5 m

[1] Labels that are required to be unique are not labels but vertex identity

# Resource Description Framework (RDF)

- Data descripted in triples subject-predicate-object
- Subjects and objects are vertices (URIs $U$ or value literals $L$ (objects only))
- Predicates are edge labels (URIs $U$)
- RDF dataset $\subseteq U \times U \times \{U \cup L\}$
- Edges are directed
- No vertex labels
  (note, every literal is per se unique)
- No properties

```
@prefix eric: <http://www.w3.org/People/EM/contact#> .
@prefix contact: <http://www.w3.org/2000/10/swap/pim/contact#> .
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .

eric:me contact:fullName "Eric Miller" .
eric:me contact:mailbox <mailto:e.miller123(at)example> .
eric:me contact:personalTitle "Dr." .
eric:me rdf:type contact:Person .
```

[http://commons.wikimedia.org/wiki/File:Rdf_graph_for_Eric_Miller.png]

- RDF Schema (RDFS)
- Set of predefined predicates and classes to describe data schema in RDF

Linked Open Data

[http://lod-cloud.net/]

# Property Graph Model

- Directed Graph
- Vertices are proper entities with
  - Label (often as type Property)
  - Properties
- Edges are "rich" relationships with
  - Label
  - Properties

- No standard
- Various different implementations
  - TinkerPop/Gremlin
  - Neo4j (allows multiple vertex labels)
  - Green-Marl (no labels)
  - ...

**TinkerPop**

# Possible Graph Data Models

| | | Structure | | | Plain Data | | Structured Data | |
|---|---|---|---|---|---|---|---|---|
| | | Directionality | Loops & Cycles | Multiple Edges | Vertex Labels | Edge Labels | Vertex Properties | Edge Properties |
| **Pure Structure Graph Models** | Basic | ○ | ○ | ○ | ○ | ○ | ○ | ○ |
| | DAG | ✓ | ○ | ○ | ○ | ○ | ○ | ○ |
| | ⋮ | | | | | | | |
| **Plain Data Graph Models** | RDF | ✓ | ✓ | ✓ | ○ | ✓ | ○ | ○ |
| | Pregel Graph Model, Graph-Oriented Object Data model (GOOD) | ✓ | ✓ | ✓ | ✓ | ✓ | ○ | ○ |
| **Structured Data Graph Models** | Green-Marl Graph Model | ✓ | ✓ | ✓ | ○ | ○ | ✓ | ✓ |
| | OrientDB | ✓ | ✓ | ✓ | ○ | ✓ | ✓ | ✓ |
| | Property Graph, Neo4j, TinkerPop | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

# Relational Representation

## RDF

- Triples fits in three-column relational table

| Subject | Predicate | Object |
|---|---|---|
| <http://www.w3.org/People/EM/contact#me> | <http://www.w3.org/2000/10/swap/pim/contact#fullName> | "Eric Miller" |
| <http://www.w3.org/People/EM/contact#me> | <http://www.w3.org/2000/10/swap/pim/contact#mailbox> | <mailto:e.miller123(at)example> |
| <http://www.w3.org/People/EM/contact#me> | <http://www.w3.org/2000/10/swap/pim/contact#personalTitle> | "Dr." |
| <http://www.w3.org/People/EM/contact#me> | <http://www.w3.org/1999/02/22-rdf-syntax-ns#type> | <http://www.w3.org/2000/10/swap/pim/contact#Person> |

## PROPERTY GRAPH

- Two universal tables: one for the vertices, one for the edges

| ID | Type | Color | Name | RAM | Nationality | Source | Target | Type | Rating |
|---|---|---|---|---|---|---|---|---|---|
| 1 | Product | black | "Apple iPad MC707LL/A" | 64 GB | | 1 | 7 | in | |
| … | … | … | … | … | … | … | … | … | … |
| 4 | Category | | "Cell Phones & Accessories" | | | 5 | 4 | part of | |
| 5 | Category | | "Phones" | | | 7 | 6 | part of | |
| … | … | … | … | … | … | … | … | … | … |

- Alternative: One universal tables per vertex and edge label (type)

# Relational Representation

## GRAPH DATA WITH FIXED SCHEMA

- One table for every vertex type and every edge type

- Think of: Two-universal-tables schema partitioned by type

- Edge types representing 1:N relationship can be presented with a simple foreign key



Social Network Benchmark

[http://ldbcouncil.org/]

# Graph-structure data

**OFTEN HIDDEN IN NON-GRAPH DATA**

- E.g. TPC-H scenario
- Customer that also is a supplier

# Graph Models vs. Relational Model

## IDENTITY OF ENTITIES

- Relational: Value-based identity
  - One or more attributes serves as identity and are declared as such per type by primary key constraint
  - Values of identity attributes are user-given
- Graph: Object identity
  - Fixed (visible or hidden) attribute serves as identity
  - Values of identity attribute are either system-generated (e.g. object id) or user-given (e.g. URI)
- Distinction is blurred by bag-semantics of SQL

## REFERENCING MECHANISMS

- Relational: Value-based reference
  - References are expressed by value equality
  - Necessity of referential integrity can be declared by a foreign key constraint
- Graph: Explicit association
  - References are expressed with a dedicated association element -> edges
  - Dedicated association element has referential integrity built in

## LIBERTY OF REFERENCING

- Relational

  Schema

  Data

- Graph

# Graph Querying

# Querying Graphs



## QUESTION

"How is friend of John Doe?"
"Couples where both like 'House of Cards'?"
"Shortest connection between
John Doe and Joe Dohn?"
...

## ANSWER

[http://www.bordalierinstitute.com/images/yeastProteinInteractionNetwork.jpg]

# Level of Analytics

[http://www.rosebt.com/blog/eight-levels-of-analytics-business-intelligence-to-business-analytics]



**Business Value**

**Business Analytics**

- **8 OPTIMISATION** — What's the best that can happen?
- **7 PREDICTIVE MODELLING** — What will happen next?
- **6 FORECASTING** — What if these trends continue?
- **5 STATISTICAL ANALYSIS** — Why is this happening?
- **4 ALERTS** — What actions are needed?
- **3 QUERY DRILLDOWN (OLAP)** — Where exactly is the problem?
- **2 AD-HOC REPORTS** — How many, how often, where?
- **1 STANDARD REPORTS** — What happened?

**Business Intelligence**

**Degree of Intelligence**

0. Transactional Data Management (OLTP)

batch
long running
imperative

**OFFLINE GRAPH ANALYTICS**

Graph Data Management Applications

ad-hoc
short running
declarative

**ONLINE GRAPH QUERYING**

28

# Graph Workloads

## ONLINE GRAPH QUERYING

- OLTP-Style
- Short, read and write, high selectivity
- Interest in proximity of one or more start nodes
- E.g., Loading and updating of you Facebook page, Facebook Graph Search



[http://blackfin360.com/2013/01/15/facebook-graph-search/]

## OFFLINE GRAPH ANALYTICS

- OLAP style
- Long, expensive, mainly read, low selectivity
- Topological analysis of whole graph
- E.g. Page rank (centrality), shortest path, connected components, ...



[http://commons.wikimedia.org/wiki/File:PageRank-hi-res.png]

# Graph Query Concepts

## ONLINE GRAPH QUERYING



## DECLARATIVE

- DataLog, SPARQL, RQL, Cypher, ...
- Focuses on the What
- Abstracts from the How
- Limited compared to normal programming languages
- Allows optimization: Optimizer takes care of the How
- Hides technical, low-level concerns, e.g. selectivities, parallelization, etc.

*I'm going to make him an offer he can't refuse.**

## OFFLINE GRAPH ANALYTICS



PageRank

## IMPERATIVE

- Gremlin, GreenMarl, Travel, Pregel, GraphLab, ...
- DSLs or APIs
- Focuses on the How
- Sets of commands
  - Graph traversal and access
  - General-propose programming language constructs
- (Almost) no restriction in expression power compared to normal programming languages
- Comfortable graph navigation and access

*Leave the gun. Take the cannoli.****

*[declarative sentence spoken by Don Corleone in The Godfather; http://grammar.about.com/od/d/g/declsenterm.htm]

**[Imperative sentences spoken by Clemenza in The Godfather; http://grammar.about.com/od/il/g/impersent09.html]

# Level of Analytics

[http://www.rosebt.com/blog/eight-levels-of-analytics-business-intelligence-to-business-analytics]



0. Transactional Data Management (OLTP)

batch
long running
imperative

IMPERATIVE
GRAPH
ANALYTICS

Graph Data
Management
Applications

ad-hoc
short running
declarative

DECLARATIVE
GRAPH
QUERYING

# Online Graph Querying – Pattern Matching

# Graph Pattern Matching



**PATTERN $p$**

Graph with place holders (A,B)

**MATCHING $p$ ON $G$**

Finds all subgraphs in $G$ that fit to $p$

# Similarity of two Graphs

ARE GRAPHS *G* AND *H* EQUAL/SIMILAR?



MANY SIMILARITY CRITERIA

- Isomorphism
- Homomorphism
- Simulation
- Bisimilarity
- ...

# Graphs Homomorphism



- Given two graphs $G(V_G, E_G)$ and $H(V_H, E_H)$
- $G$ and $H$ are homomorph, if there is a surjective function $\sigma : V_G \rightarrow V_H$ (left-total, right-total, right-unique) such that $(v_i, v_j) \in E_G \rightarrow \left(\sigma(v_i), \sigma(v_j)\right) \in E_H$ ($G$ preserves adjacency of $H$, i.e. an edge in $G$ has to exist in $H$ as well)

- Note, there may be multiple functions $\sigma$, i.e., multiple homomorphism between two graphs

# Graph Isomorphism

## GRAPH HOMOMORPHISM



- Given two graphs $G(V_G, E_G)$ and $H(V_H, E_H)$
- $G$ and $H$ are homomorph,
- if there is a surjective function $\sigma: V_G \rightarrow V_H$ (left-total, right-total, right-unique)
  such that $(v_i, v_j) \in E_G \rightarrow \big(\sigma(v_i), \sigma(v_j)\big) \in E_H$
  ($G$ preserves adjacency of $H$)

## GRAPH ISOMORPHISM



- Given two graphs $G(V_G, E_G)$ and $H(V_H, E_H)$
- $G$ and $H$ are isomorph,
- if there is a bijective function $\sigma: V_G \rightarrow V_H$ (left-total, left-unique, right-total, and right-unique)
  such that $(v_i, v_j) \in E_G \leftrightarrow \big(\sigma(v_i), \sigma(v_j)\big) \in E_H$
  ($G$ preserves adjacency and non-adjacency of $H$ and vice versa)

Cool!

Can we apply this to find subgraphs?

[see also video tutorial on graph gsomorphism by Sarada Herke:
https://www.youtube.com/watch?v=yFpRpxOry-A]

# Subgraph Homomorphism Query



- Given a query graphs $Q(V_G, E_G)$ and data graph $G(V_H, E_H)$
- Graph $R(V_R, E_R)$ is a result for $Q$ if
  - $V_R \subseteq V_G$ and $E_R \subseteq E_G$ $(R$ is a subgraph of $G)$ and
  - there is a surjective function $\sigma: V_G \rightarrow V_H$ $(Q$ and $R$ are homomorph$)$ such that
    $(v_i, v_j) \in E_Q \leftrightarrow \left(\sigma(v_i), \sigma(v_j)\right) \in E_R$ $(Q$ preserves adjacency of $R$ with no extra edges in $R)$ and
    $\forall (v_Q, v_R) \in \sigma, v_Q \sim v_R$ (vertex properties match) and $\forall (v_i, v_j) \in E_Q, (v_i, v_j) \sim \left(\sigma(v_i), \sigma(v_j)\right)$ (edge properties match)
- $Q$ can have more vertices and edges than $R$, i.e., $|V_Q| \geq |V_R|$ and $|E_Q| \geq |E_R|$ holds.

- Note: $R$ is not given, $R$ has to be determined by the query mechanism -> search problem

# Subgraph Isomorphism Query

- Given a query graphs $Q(V_G, E_G)$ and data graph $G(V_H, E_H)$
- Graph $R(V_R, E_R)$ is a result for $Q$ if
  - $V_R \subseteq V_G$ and $E_R \subseteq E_G$ ($R$ is a subgraph of $G$) and
  - there is a bijective function $\sigma: V_G \to V_H$ ($Q$ and $R$ are isomorph) such that
    $(v_i, v_j) \in E_Q \leftrightarrow \left(\sigma(v_i), \sigma(v_j)\right) \in E_R$ ($Q$ preserves adjacency of $R$ with no extra edges in $R$) and
    $\forall (v_Q, v_R) \in \sigma, v_Q \sim v_R$ (vertex properties match) and $\forall (v_i, v_j) \in E_Q, (v_i, v_j) \sim \left(\sigma(v_i), \sigma(v_j)\right)$ (edge properties match)
- $Q$ and $R$ will have the same number of vertices and edges, i.e., $|V_Q| = |V_R|$ and $|E_Q| = |E_R|$ holds.

> **Single vertex in $V_G$ can be matched
> multiple times in a homomorphic subgraph but only once in an isomorphic subgraph.**

# Why Homomorphism is useful…

EXAMPLE: LOOK FOR ALL PAIRS OF FRIENDS AND THE CITY EACH FRIEND LIVES IN

QUERY



DATA



ISOMORPHISM FINDS ONLY FRIENDS LIVING IN DIFFERENT CITIES
- (Leipzig, Chris, Anne, Berlin), (Leipzig, Chris, Mary, Berlin) … and permutation of these

HOMOMORPHISM ADDITIONALLY FINDS FRIENDS LIVING IN THE SAME CITY
- (Berlin, Mary, Anne, Berlin) … and permutation of these

# Induced Subgraph Isomorphism

### EXAMPLE

- Data graph:            Query graph:



- Does it have a match? How many?

- One solution:



- What about the other edges? Could we forbid them?
- With induced subgraph isomorphism semantics, example query has no match!

### INDUCED SUBGRAPH

- Vertex-induced
- Is a subset of the vertices of a graph together with any edges whose endpoints are both in this subset.

### INDUCED SUBGRAPH ISOMORPHISM

- Stricter isomorphism
- Given query graph $Q(V_Q, E_Q)$ and data graph $G(V_G, E_G)$
- Graph $R(V_R, E_R)$ is a result for $Q$ if
  - $V_R \subseteq V_G$ and $\boldsymbol{E_R = \{(v_i, v_j) | (v_i, v_j) \in E_G \wedge v_i, v_j \in V_R\}}$ ($R$ is a vertex-induced subgraph of $G$) and
  - there is a bijective function $\sigma: V_Q \to V_R$ ($Q$ and $R$ are isomorph) such that …

- In graph query languages, typically explicit negation use instead
- Induced subgraph homomorphism also possible

Dresden Database
Systems Group

## EXAMPLE: FIND TRIANGLES

- Data graph:



- Query: All triangles



```
q(X,Y,Z) <- e(X,Y),e(Y,Z),e(Z,X).
```

- Result:

Duplicate results, same subgraph but different isomorphism $\sigma$

| x | y | z |
|---|---|---|
| 1 | 2 | 3 |
| 1 | 3 | 2 |
| 2 | 1 | 3 |
| 2 | 3 | 1 |
| 3 | 1 | 2 |
| 3 | 2 | 1 |
| 2 | 3 | 4 |
| 2 | 4 | 3 |
| 3 | 2 | 4 |
| 3 | 4 | 2 |
| 4 | 2 | 3 |
| 4 | 3 | 2 |

## DUPLICATE RESULTS

- Given two results $R(V_R, E_R)$ and $S(V_S, E_S)$
- $R$ and $S$ are equivalent iff $V_R = V_S$ and $E_R = E_S$ (both denote the same subgraph)

## SUBGRAPH ISOMORPHISM W/O DUPLICATES

- Given query graph $Q(V_Q, E_Q)$ and data graph $G(V_G, E_G)$
- Graph $R(V_R, E_R)$ is a result for $Q$ if
  - $V_R \subseteq V_G$ and $E_R \subseteq E_G$ ($R$ is a subgraph of $G$) and
  - there is a bijective function $\sigma: V_Q \to V_R$ ($Q$ and $R$ are isomorph) such that $(v_i, v_j) \in E_Q \leftrightarrow \left(\sigma(v_i), \sigma(v_j)\right) \in E_R$ ($R$ preserves adjacency of $Q$ with no extra data) and properties match 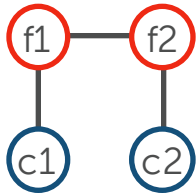and $\forall v_i, v_j \in V_Q: \left(v_i <_{V_Q} v_j\right) \leftrightarrow \left(\sigma(v_i) <_{V_R} \sigma(v_j)\right)$ assuming a total order $<_V$ on a vertex set $V$ (allows only one $\sigma$ per subgraph)

TECHNISCHE
UNIVERSITÄT
DRESDEN

# Comparison of Semantics



```
q(X,Y,Z) <-
```

| | INDUCED-SUBGRAPH ISOMORPHISM | ⊆ | SUBGRAPH ISOMORPHISM | ⊆ | SUBGRAPH HOMOMORPHISM |
|---|---|---|---|---|---|
| w/ duplicates | `e(X,Y),e(X,Z),` `!e(Y,X),!e(Y,Z),` `!e(Z,Y),!e(Z,X),` `X!=Y,Y!=Z,Z!=X.` | | `e(X,Y),e(X,Z),` `X!=Y,Y!=Z,Z!=X.` | | `e(X,Y),e(X,Z).` |
| w/o duplicates | `e(X,Y),e(X,Z),` `!e(Y,X),!e(Y,Z),` `!e(Z,Y),!e(Z,X),` `X<Y,Y<Z.` | | `e(X,Y),e(X,Z),` `X<Y,Y<Z.` | | `e(X,Y),e(X,Z),` `X<=Y,Y<=Z.` |

# SPARQL

## QUERY LANGUAGE FOR RDF DATA

- Selection of subgraph with a triple patterns
- Triple pattern is a set of triples containing variables
- One variable binding of a pattern forms a tuple
- All unique variable binding form a table
- Projection to variable of interest yields query result

SELECT ?p, ?s
WHERE
?p type Person
?p likes ?f
?f type Person
?f sex ?s

| ?p | ?s |
|----|----|

| S | P | O |
|---|---|---|
| Lucy | born | 1982 |
| Peter | born | 1985 |
| Jen | born | 1987 |
| Lucy | sex | F |
| Peter | sex | M |
| Jen | sex | F |
| Lucy | likes | Peter |
| Peter | likes | Lucy |
| Jen | likes | Lucy |
| Peter | likes | Jen |

# SPARQL

## Query language for RDF data

- Selection of subgraph with a triple patterns
- Triple pattern is a set of triples containing variables
- One variable binding of a pattern forms a tuple
- All unique variable binding form a table
- Projection to variable of interest yields query result

| S | P | O |
|---|---|---|
| Lucy | born | 1982 |
| Peter | born | 1985 |
| Jen | born | 1987 |
| Lucy | sex | F |
| Peter | sex | M |
| Jen | sex | F |
| Lucy | likes | Peter |
| Peter | likes | Lucy |
| Jen | likes | Lucy |
| Peter | likes | Jen |

SELECT  ?p, ?s
WHERE   ?p type Person
        ?p likes ?f
        ?f type Person
        ?f sex ?s

| ?p | ?s |
|---|---|
| Lucy | M |
| Peter | F |
| Jen | F |
| Peter | F |

SELECT  ?p, ?fof
WHERE   ?p type Person
        ?p likes ?f
        ?f type Person
        ?f like ?fof

| ?p | ?fof |
|---|---|

# SPARQL

## Query language for RDF data

- Selection of subgraph with a triple patterns
- Triple pattern is a set of triples containing variables
- One variable binding of a pattern forms a tuple
- All unique variable binding form a table
- Projection to variable of interest yields query result

| S | P | O |
|---|---|---|
| Lucy | born | 1982 |
| Peter | born | 1985 |
| Jen | born | 1987 |
| Lucy | sex | F |
| Peter | sex | M |
| Jen | sex | F |
| Lucy | likes | Peter |
| Peter | likes | Lucy |
| Jen | likes | Lucy |
| Peter | likes | Jen |

SELECT    ?p, ?s
WHERE    ?p type Person
       ?p likes ?f
       ?f type Person
       ?f sex ?s

| ?p | ?s |
|---|---|
| Lucy | M |
| Peter | F |
| Jen | F |
| Peter | F |

SELECT    ?p, ?fof
WHERE    ?p type Person
       ?p likes ?f
       ?f type Person
       ?f like ?fof

| ?p | ?fof |
|---|---|
| Lucy | Jen |
| Peter | Lucy |
| Peter | Peter |
| Jen | Peter |
| Lucy | Lucy |

# Cypher neo4j

## MATCH-Clause

- Primary way of getting data from a Neo4j database
- Allows you to specify the patterns
- Named pattern element, e.g. (p:Person), will be bound to the match instance
- Query can have multiple MATCH-clauses

## WHERE-Clause (Optional)

- Allows additional complex predicates in the pattern
- Allows joining two matches

## RETURN-Clause

- Projects to the result set
- Allows projection to nodes, edges, and properties

## ORDER BY-Clause (like in SQL)



| p.name | f.sex |
|--------|-------|

```
MATCH (p:Person)-[:Likes]->(f:Person)
RETURN p.name, f.sex
```

# Cypher neo4j

## MATCH-Clause

- Primary way of getting data from a Neo4j database
- Allows you to specify the patterns
- Named pattern element, e.g. (p:Person), will be bound to the match instance
- Query can have multiple MATCH-clauses

## WHERE-Clause (Optional)

- Allows additional complex predicates in the pattern
- Allows joining two matches

## RETURN-Clause

- Projects to the result set
- Allows projection to nodes, edges, and properties

## ORDER BY-Clause (like in SQL)



: Person
name=Lucy
born=1982
sex=F

Likes
stars=5

Likes
stars=5

Likes
stars=4

: Person
name=Jen
born=1987
sex=F

: Person
name=Peter
born=1985
sex=M

Likes
stars=3

```
MATCH (p:Person)-[:Likes]->(f:Person)
RETURN p.name, f.sex
```

| p.name | f.sex |
|--------|-------|
| Lucy   | M     |
| Peter  | F     |
| Jen    | F     |
| Peter  | F     |

| p.name | fof.name |
|--------|----------|

```
MATCH (p:Person)-[:Likes]->(:Person) -[:Likes]->(fof:Person)
RETURN p.name, fof.name
```

# Cypher neo4j

## MATCH-CLAUSE

- Primary way of getting data from a Neo4j database
- Allows you to specify the patterns
- Named pattern element, e.g. (p:Person), will be bound to the match instance
- Query can have multiple MATCH-clauses

## WHERE-CLAUSE (OPTIONAL)

- Allows additional complex predicates in the pattern
- Allows joining two matches

## RETURN-CLAUSE

- Projects to the result set
- Allows projection to nodes, edges, and properties

## ORDER BY-CLAUSE (LIKE IN SQL)



```
MATCH (p:Person)-[:Likes]->(f:Person)
RETURN p.name, f.sex
```

| p.name | f.sex |
|--------|-------|
| Lucy   | M     |
| Peter  | F     |
| Jen    | F     |
| Peter  | F     |

```
MATCH (p:Person)-[:Likes]->(:Person) -[:Likes]->(fof:Person)
RETURN p.name, fof.name
```

| p.name | fof.name |
|--------|----------|
| Lucy   | Jen      |
| Peter  | Lucy     |
| Peter  | Peter    |
| Jen    | Peter    |
| Lucy   | Lucy     |

# Online Graph Querying – Query Types

# Types of Graph Pattern Queries

[Peter T. Wood. Query Languages for Graph Database. SIGMOD Record 41(1), 50–60, March 2012]

## OVERVIEW

- Conjunctive queries (standard subgraph matching)
- Regular path queries (reachability)
- Conjunctive regular path query
- ...

## GRAPH DATA MODEL FOR FOLLOWING

- RDF-like data
- Graph $G(V, E, \Sigma)$ with
  - V being the set of vertices,
  - $E \subseteq V \times \Sigma \times V$ being the set of labeled edges,
  - and $\Sigma$ being the set (or alphabet) of labels



Edge labels
**lo**catedIn
**li**vesIn
**b**ornIn
**w**on

# Conjunctive Queries (Std. Matching)

## IDEA

- Query is given as a set of edge predicates
- Each edge predicate consists of a pair vertex variables and an edge label
- A set of variable bindings is a valid answer iff all predicates hold on the data graph

## DEFINITION

- Query $Q$ is an expression

$$ans(z_1, \ldots, z_n) \leftarrow \bigwedge_{1 \leq i \leq m} (x_i, a_i, y_i)$$

- Each $x_i \in X$ and $y_i \in Y$ is a vertex variable or a constant from $V$
- Each $a_i \in \Sigma$ is an edge label
- Each $z_i$ is some $x_i$ or $y_i$

*Correspond to homomorphism*

## SEMANTICS

- Let $\sigma: X \cup Y \to V$ be a specific selection of variable bindings, i.e., a mapping to vertices of $G$
- Say relation $(G, \sigma) \models Q$ holds iff $(\sigma(x_i), a_i, \sigma(y_i)) \in E$ for $1 \leq i \leq m$, i.e., $\sigma$ maps the query pattern to valid subgraphs of $G$
- Then the query result $Q(G)$ is the set of tuples $(\sigma(z_1), \ldots, \sigma(z_n))$ such that $(G, \sigma) \models Q$

# Conjunctive Queries (Std. Matching)

## EXAMPLE

- All authors born in South Africa who have won both the Nobel and Booker prizes

$$ans(x) \leftarrow \begin{array}{l} (x, \text{hasWon}, \text{Nobel}), \\ (x, \text{hasWon}, \text{Booker}), \\ (x, \text{bornIn}, \text{South Africa}) \end{array}$$

- Visually:



- Result?



Edge labels
**lo**catedIn
**li**vesIn
**b**ornIn
**w**on



52

# Conjunctive Queries (Std. Matching)

## EXAMPLE

- All authors born in South Africa who have won both the Nobel and Booker prizes

$$ans(x) \leftarrow \begin{array}{l} (x, \text{hasWon}, \text{Nobel}), \\ (x, \text{hasWon}, \text{Booker}), \\ (x, \text{bornIn}, \text{South Africa}) \end{array}$$

- Visually:

- Result?



Edge labels
**lo**catedIn
**li**vesIn
**b**ornIn
**w**on

# Conjunctive Queries (Std. Matching)

**EXTENSION TO PROPERTY GRAPHS**

- Add predicates on properties
- For vertex properties
  - All authors born in South Africa before 1930 who have won Nobel and Booker prizes

$$ans(x) \leftarrow \begin{array}{l} (x, \text{hasWon}, \text{Nobel}), \\ (x, \text{hasWon}, \text{Booker}), \\ (x, \text{bornIn}, \text{South Africa}), \\ \textcolor{red}{(x, \text{year}, < 1930)} \end{array}$$

  - Extra syntax required for non-equi predicates
- For edge properties
  - Extra syntax required on existing edge predicates
  - All authors born in South Africa who have won Nobel once and Booker twice

$$ans(x) \leftarrow \begin{array}{l} (x, \text{hasWon}: \textcolor{red}{(x = 1)}, \text{Nobel}), \\ (x, \text{hasWon}: \textcolor{red}{(x = 2)}, \text{Booker}), \\ (x, \text{bornIn}, \text{South Africa}), \end{array}$$



Edge labels
**lo**catedIn
**li**vesIn
**b**ornIn
**w**on

54

# Regular Path Query (Reachability)

[Peter T. Wood. Query Languages for Graph Database. SIGMOD Record 41(1), 50−60, March 2012]

## IDEA

- Query is given as a path predicate consisting of a pair vertex variables and a path expression
- Path expression is a regular expression of edge labels
- A pair of variable bindings is a valid answer iff the respective vertices are connect in hold on the data graph by a path conforming to the path expression

## DEFINITION

- Query $Q$ is an expression $\quad\quad\quad\quad ans(x,y) \leftarrow (x,r,y)$
- $x \in V$ and $y \in V$ are vertex variables
- $r \in \Sigma^*$ is a regular expression over alphabet of edge labels $\Sigma$

## SEMANTICS

- A path $p$ between $v_0$ and $v_m$ in $G$ is a sequence $v_0 a_0 v_1 a_1 v_2 \ldots v_{m-1} a_{m-1} v_m$, with $v_i \in V$, $a_i \in \Sigma$, and $(v_i, a_i, v_{i+1}) \in E$
- Let $\lambda(p) \in \Sigma^*$ be the label of the path $p$, with $\lambda(p) = a_0 a_1 \ldots a_{m-1}$
- Let $L(r)$ be the language denoted by the regular expression $r$, i.e. set all of all possible path labels denoted by $r$
- Path $p$ satisfies $r$ if $\lambda(p) \in L(r)$, i.e. the path's label satisfies the regular expression
- Then the query result $Q(G)$ is the set of all pairs of nodes $(x,y)$ in $G$ such there is a path from $x$ to $y$ which satisfies $r$

# Regular Path Query (Reachability)

## EXAMPLE

- All authors and where they live in or are born

$$ans(x, y) \leftarrow (x, (b|li) \cdot lo^*, y)$$

- Visually:



- Result?

| Neruda | South America |
|---|---|
| Coetzee | Australia |
| Coetzee | South Africa |
| Gordimer | South Africa |
| Carey | Australia |



Edge labels
**lo**catedIn
**li**vesIn
**b**ornIn
**w**on

# Regular Path Query (Reachability)

[Peter T. Wood. Query Languages for Graph Database. SIGMOD Record 41(1), 50−60, March 2012]

## EXAMPLE

- All authors and where they live in or are born

$$ans(x, y) \leftarrow (x, (\text{b}|\text{li}) \cdot \text{lo}^*, y)$$

- Visually:



- Result?

| Neruda | South America |
|--------|---------------|
| Coetzee | Australia |
| Coetzee | South Africa |
| Gordimer | South Africa |
| **Carey** | **Australia** |



Edge labels
**lo**catedIn
**li**vesIn
**b**ornIn
**w**on

# Regular Path Query (Reachability)

## Extension to vertex label

- Path expressions include vertex labels, e.g., $(\text{bornIn}|\text{livesIn}) \cdot \text{City} \cdot \text{locatedIn}^* \cdot \text{Continent}$
- Let $L_V$ and $L_E$ be the set of vertex and edge labels and $\lambda$ the labeling function
- A regular expression $r$ is over alphabet of edge and vertex labels pairs $(L_E \times L_V)^*$
- Let the label of a path be $\lambda(p) = a_0\lambda(v_1)a_1 \dots a_{m-1}\lambda(v_m)$ with $a_i \in L_E$ and $\lambda(v_i) \in L_V$ and $\lambda(p) \in (L_E \times L_V)^*$
- As before: Path $p$ satisfies $r$ if $\lambda(p) \in L(r)$, i.e. the path's label satisfies the regular expression

## Extension to properties

- Path expression include property predicates,
  e.g., $\text{livesIn:}[\text{since} < 1990] \cdot \text{City:}[\text{population} > 10\text{Mio}] \cdot \text{locatedIn}^+ \cdot \text{Continent}$
- Path expressions quickly become hard to read, cf. XPath and XQuery
- In Datalog rules:

$$ans(x,y) \leftarrow livesIn(x,z), eSince(x,z,s), s < 1990,$$
$$city(z), vPopulation(z,p), p > 10\text{Mio}, loStar(z,y), continent(y)$$
$$loStar(x,y) \leftarrow locatedIn(x,y)$$
$$loStar(x,y) \leftarrow locatedIn(x,z), loStar(z,y)$$

# Conjunctive Regular Path Queries (CRPQs)

## IDEA

- Query is given as a set of path predicates
- Each path predicate consists of a pair vertex variables and a regular expression of edge labels
- A set of variable bindings is a valid answer iff all path predicates hold on the data graph

## DEFINITION

- Query $Q$ is an expression

$$ans(z_1, \ldots, z_n) \leftarrow \bigwedge_{1 \leq i \leq m} (x_i, r_i, y_i)$$

- Each $x_i \in X$ and $y_i \in Y$ is a vertex variable or a constant from $V$
- Each $r_i \in \Sigma^*$ is a regular expression over alphabet of edge labels $\Sigma$
- Each $z_i$ is some $x_i$ or $y_i$

## SEMANTICS

- Let $\sigma: X \cup Y \rightarrow V$ be a specific selection of variable bindings, i.e., a mapping to vertices of $G$
- Say relation $(G, \sigma) \vDash Q$ holds iff, for $1 \leq i \leq m$ there exists a path $p_i$ in $G$ from $\sigma(x_i)$ to $\sigma(y_i)$ such that $\lambda(p) \in L(r)$
- Then the query result $Q(G)$ is the set of tuples $\big(\sigma(z_1), \ldots, \sigma(z_n)\big)$ such that $(G, \sigma) \vDash Q$

# Conjunctive Regular Path Queries (CRPQs)

## EXAMPLE

- All winners of Nobel and Booker and where they live

$$ans(x, y) \leftarrow \begin{matrix} (x, \text{hasWon}, \text{Nobel}), \\ (x, \text{hasWon}, \text{Booker}), \\ (x, \text{li} \cdot \text{lo}^*, y) \end{matrix}$$

- Visually:



- Result?

| Coetzee | Australia |
|---------|-----------|



Edge labels
**lo**catedIn
**li**vesIn
**b**ornIn
**w**on

# Conjunctive Regular Path Queries (CRPQs)

## EXAMPLE

- All winners of Nobel and Booker and where they live

$$ans(x, y) \leftarrow \begin{matrix} (x, \text{hasWon}, \text{Nobel}), \\ (x, \text{hasWon}, \text{Booker}), \\ (x, \text{li} \cdot \text{lo}^*, y) \end{matrix}$$

- Visually:



- Result?

| Coetzee | Australia |
|---------|-----------|



Edge labels
**lo**catedIn
**li**vesIn
**b**ornIn
**w**on

# Further Query Types

## Union conjunctive Queries (UQs)

- Adds disjunction
- Example: $ans(x) \leftarrow (x, \text{hasWon}, \text{Booker}) \lor (x, \text{hasWon}, \text{Nobel})$ which give all price winners
- Multiple conjunctive queries with intersecting variable sets
- Result is the union of the result each conjunctive query projected to the intersection of all variable sets

## Two-way regular path queries (2RPQs)

- Allows to express backward traversal of edge types
- Example: $ans(x, y) \leftarrow (x, \text{hasWon} \cdot -\text{hasWon}, y)$ which gives all author pairs where both have won the same price

## Combination of all: Union conjunctive two-way regular path queries (UC2RPQs)

- Class of Queries that can be expressed with SPARQL 1.1 and Neo4j Cypher (differences in the exact semantics)
- Can also be expressed in the relational world with Datalog or SQL incl. recursive common table expressions

# Level of Analytics

[http://www.rosebt.com/blog/eight-levels-of-analytics-business-intelligence-to-business-analytics]



Business Value

**Business Analytics**

⑧ **OPTIMISATION**
What's the best that can happen?

⑦ **PREDICTIVE MODELLING**
What will happen next?

⑥ **FORECASTING**
What if these trends continue?

⑤ **STATISTICAL ANALYSIS**
Why is this happening?

④ **ALERTS**
What actions are needed?

③ **QUERY DRILLDOWN (OLAP)**
Where exactly is the problem?

② **AD-HOC REPORTS**
How many, how often, where?

① **STANDARD REPORTS**
What happened?

**Business Intelligence**

Degree of Intelligence

0. Transactional Data Management (OLTP)

batch
long running
imperative

Graph Data
Management
Applications

ad-hoc
short running
declarative

Aggregations &
composability
needed

Graph
Matching

# Break

# Level of Analytics

[http://www.rosebt.com/blog/eight-levels-of-analytics-business-intelligence-to-business-analytics]



**Graph Data Management Applications**

batch
long running
imperative

ad-hoc
short running
declarative

Aggregations &
composability
needed

**Graph Matching**

0. Transactional Data Management (OLTP)

# Example: Air Traffic Surveillance

**CAPTURED OF SURVEILLANCE DATA** ➡ **INFORMATION OF INTEREST**

- Fine granularity data
- Low abstraction
- Geo position, timestamp, plane id

- Stepwise abstraction from base data to aggregated information



[http://science.howstuffworks.com/transport/flight/modern/air-traffic-control2.htm]

Composability: Operations closed on graphs

Dangerous routes

Flight routes    Critical landing

En route   Dangerous approach   Landing    Delays

Surveillance data    Airports    Flight plan

Base data

Derived data

TECHNISCHE UNIVERSITÄT DRESDEN

# Composability

# Graph Transformation Rules

MATCH -> VARIABLE BINDINGS -> PRODUCTION

| Match Pattern | Variable Bindings | Production Pattern |
|---|---|---|

| $x_1$ | ... | $x_n$ |
|---|---|---|
| | | |
| | | |

## DIFFERENT TYPES OF VARIABLES

- Value variable
- Vertex variable
- (Edge variables)
- (Path variables)
- ((Sub)Graph variables)

## PRODUCTION

- Existing vertices from bound vertex variables
- New vertices with new unbound vertex variables
- Edges either implicit (via vertex variable) or explicit (with edge variables)
- Existing values from bound value variable

# Graph Transformation Rule

EXAMPLE

Match Pattern              Production Pattern

$$\begin{pmatrix} h\colon \text{Person}(\text{sex} = \text{m}), \\ w\colon \text{Person}(\text{sex} = \text{f}), \\ h - w\colon \text{married}(\text{since} = d) \end{pmatrix} \rightarrow \begin{pmatrix} h, w, \\ m\colon \text{Marriage}(\text{date} = d), \\ m \rightarrow h\colon \text{husband}, m \rightarrow w\colon \text{wife} \end{pmatrix}$$

Vertex variables

Person          Person

$h$ ——— $w$

sex=m    since=$d$    sex=f

Value variables

# Graph Transformation Rule

EXAMPLE

Match Pattern                          Production Pattern

$$\begin{pmatrix} h: \mathrm{Person}(\mathrm{sex} = \mathrm{m}), \\ w: \mathrm{Person}(\mathrm{sex} = \mathrm{f}), \\ h - w: \mathrm{married}(\mathrm{since} = d) \end{pmatrix} \rightarrow \begin{pmatrix} h, w, \\ m: \mathrm{Marriage}(\mathrm{date} = d), \\ m \rightarrow h: \mathrm{husband}, m \rightarrow w: \mathrm{wife} \end{pmatrix}$$

Unbound vertex variable $m$ produces a new vertex for every match

Unbound vertex pairs $m \rightarrow h$ and $m \rightarrow w$ produce new edges for every match

Vertex variables

Person          Person

$h$  since=$d$  $w$

sex=m           sex=f

Value variables

$h$        $w$

husband   Marriage   wife

$m$

date=$d$

# Result Presentation

## RULE (WITHOUT TRANSFORMATION)

- e.g., Pairs of friends:     ○—○

## ISOLATED MATCHES

- Each match separately independently of vertex identity

○—○  ○—○  ○—○  ○—○  ○—○  ○—○

- Vertices taking part in multiple matches have to duplicated
- Good for querying paths, further combining individual matches and result iteration

## MERGED MATCHES

- All matches form a (partitioned) graph based

- No vertex duplication necessary
- Keeps topology of source graph
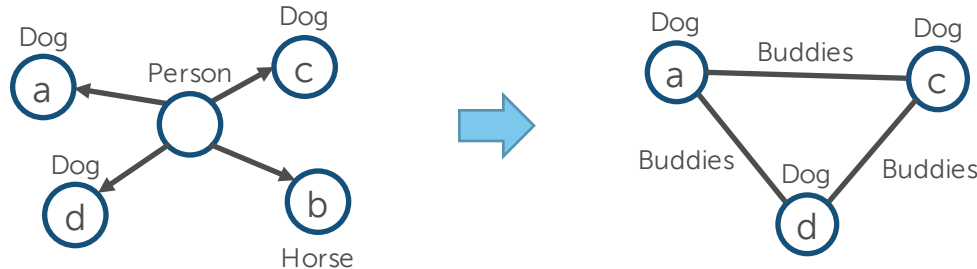
# Merged Transformation Results

## EXAMPLE

- Rule

$$\big((p\colon \text{Person}, d_1\colon \text{Dog}, d_2\colon \text{Dog}, p \to d_1, p \to d_2) \to (d_1, d_2, d_1 - d_2\colon \text{Buddies})\big)$$



- Data

# Side Note on Syntax

**You do not like the syntax?**

- How about

  ```
  SELECT NODE d1, NODE d2, UNDIRECTED EDGE d1 TO d2(Buddies)
  FROM NODE p(Person), NODE d1(Dog), NODE d2(Dog), EDGE p TO d1, EDGE p TO d2
  ```

- Or

  ```
  SELECT d1, d2, d1--d2:Buddies FROM p:Person, d1:Dog, d2:Dog, p->d1, p->d2
  ```

- Or
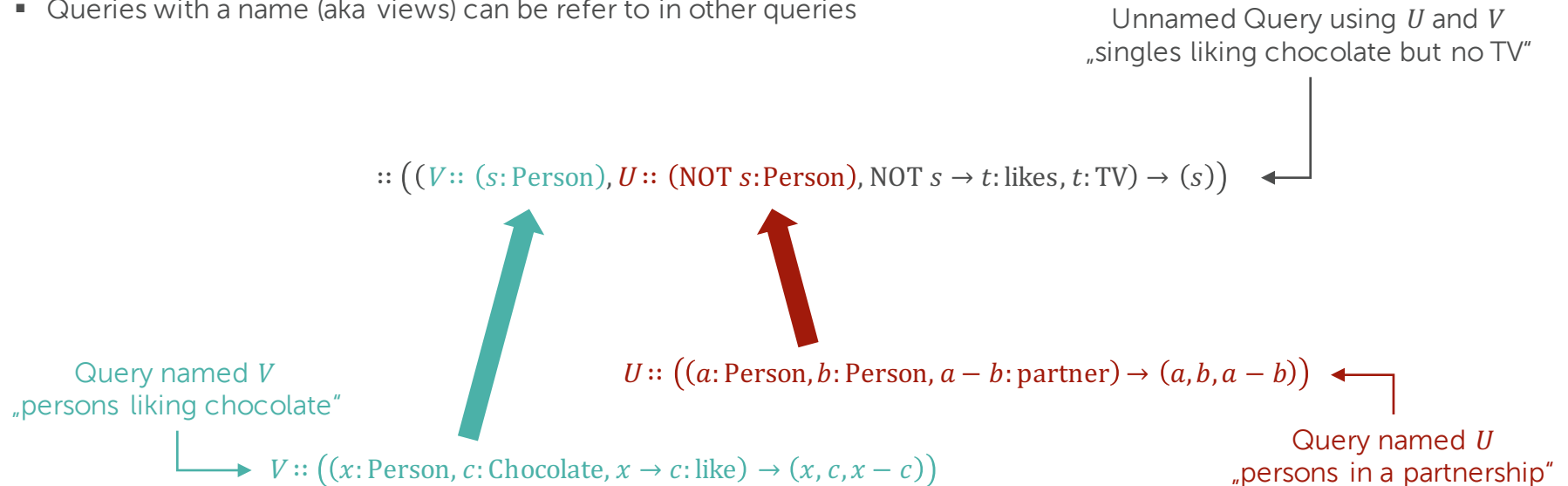
  ```
  SELECT (d1)-[:Buddies]-(d2) FROM (d1:Dog)<--(:Person)-->(d2:Dog)
  ```

- Or

  ```
  (V(p,Person),V(d1,Dog),V(d2,Dog),E(p,>,d1),E(p,>,d2)) -> (V(d1),V(d2),E(d1,-,d2)).
  ```
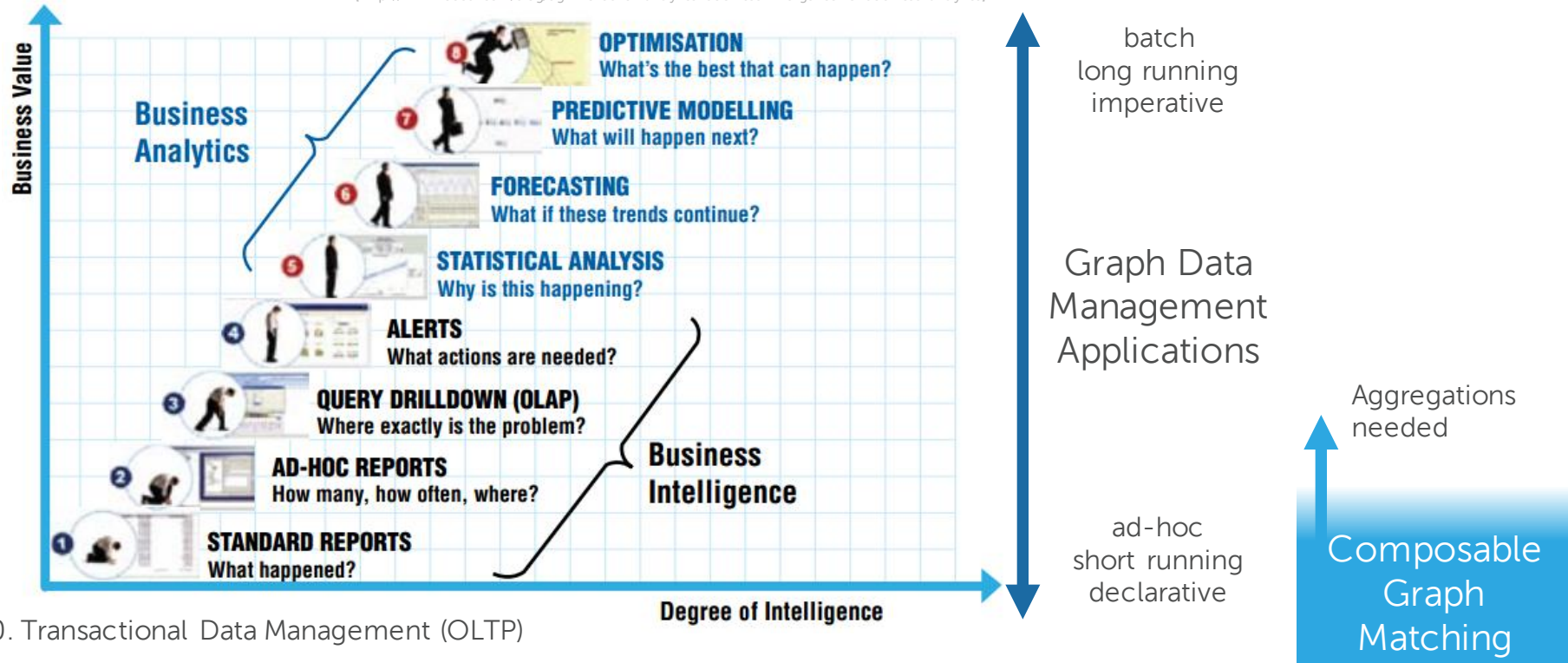
- There is some syntactical freedom as along as you stick with the language principles

73

# Composability

- Queries with a name (aka views) can be refer to in other queries

Unnamed Query using $U$ and $V$
„singles liking chocolate but no TV"

$$:: \left( \left( V :: (s: \text{Person}), U :: (\text{NOT } s: \text{Person}), \text{NOT } s \rightarrow t: \text{likes}, t: \text{TV}) \rightarrow (s) \right) \right)$$

Query named $V$
„persons liking chocolate"

$$V :: \left( (x: \text{Person}, c: \text{Chocolate}, x \rightarrow c: \text{like}) \rightarrow (x, c, x - c) \right)$$

$$U :: \left( (a: \text{Person}, b: \text{Person}, a - b: \text{partner}) \rightarrow (a, b, a - b) \right)$$

Query named $U$
„persons in a partnership"

# Level of Analytics

[http://www.rosebt.com/blog/eight-levels-of-analytics-business-intelligence-to-business-analytics]



batch
long running
imperative

Graph Data
Management
Applications

Aggregations
needed

ad-hoc
short running
declarative

Composable
Graph
Matching

0. Transactional Data Management (OLTP)
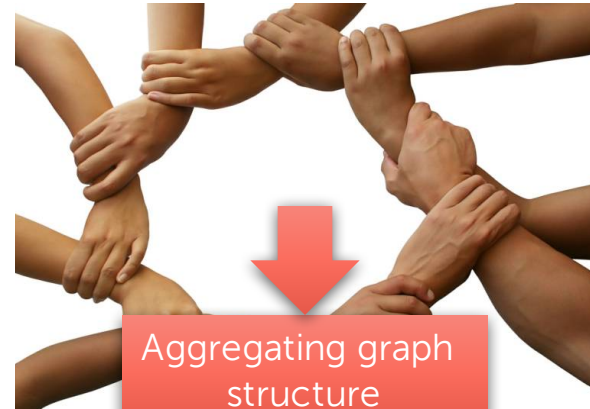
# Aggregating Graphs

## PROPERTIES OF ENTITIES

- Captured/measured values
- What are the sales figures/temperatures/etc.?
- Multidimensional data/time series/matrixes

## CONNECTIONS BETWEEN ENTITIES

- Network structure
- What do the friends of your customers buy?
- Graph data



Aggregating graph data
(vertex and edge properties)



Aggregating graph structure

# Aggregating Graph Data

# Aggregation in Graph Transformation

## GERNALE

- Per (vertex) production predicate
- Edges inherit grouping if they connect one or two grouping vertices
- Edges can have own grouping attributes, that are added to the inherit grouping attributes
- Based on match variables
- All match variables not used in grouping can be used in the same predicate only in an aggregation function
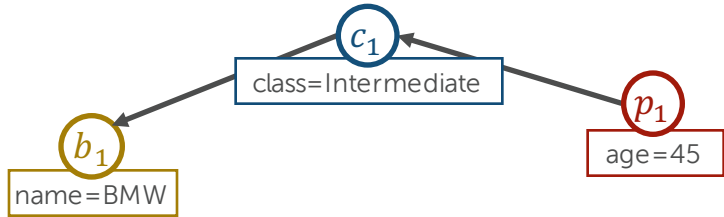
## SYNTAX

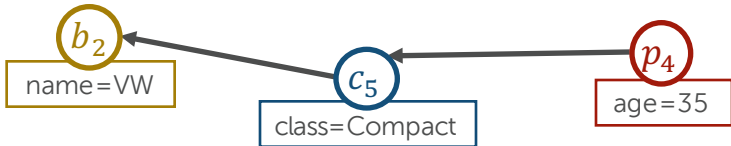- <productionGroupPredicate> ::= <productionGroupVertexVariable> "@" {<variable> ","} <labels>? <attributes>?
- <variable> ::= <matchVertexVariable> | <valueVariable>

## EXAMPLE

- $$\begin{pmatrix} p\colon \text{Person}(\text{age} = a\ ), \\ b\colon \text{Profession}, p \to b\colon \text{works}{-}\text{in} \end{pmatrix} \to \begin{pmatrix} g_a\,@a(\text{name} = a), \\ g_b\,@b(\text{name} = b.\,\text{name}), \\ g_a - g_b\big(\text{number} = \text{CNT}(p)\big) \end{pmatrix}$$

# Aggregation

$$\left( \begin{array}{c} p: \text{Person}(\text{age} = a), \\ c: \text{Car}(\text{class} = k), \\ p \rightarrow c: \text{drives}, \\ b: \text{Brand}(), \\ c \rightarrow b: \text{madeBy} \end{array} \right) \rightarrow \left( \begin{array}{c} \phantom{xxxxxxxxxxx} \end{array} \right)$$
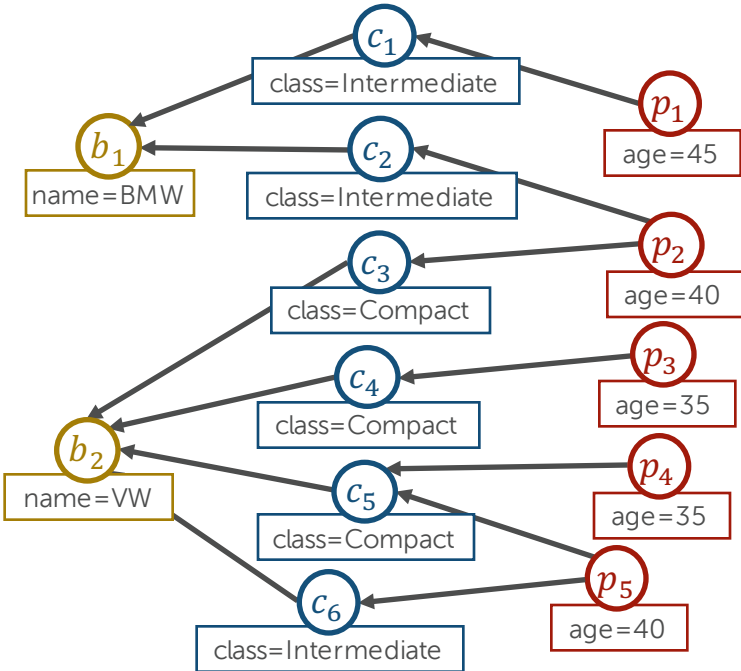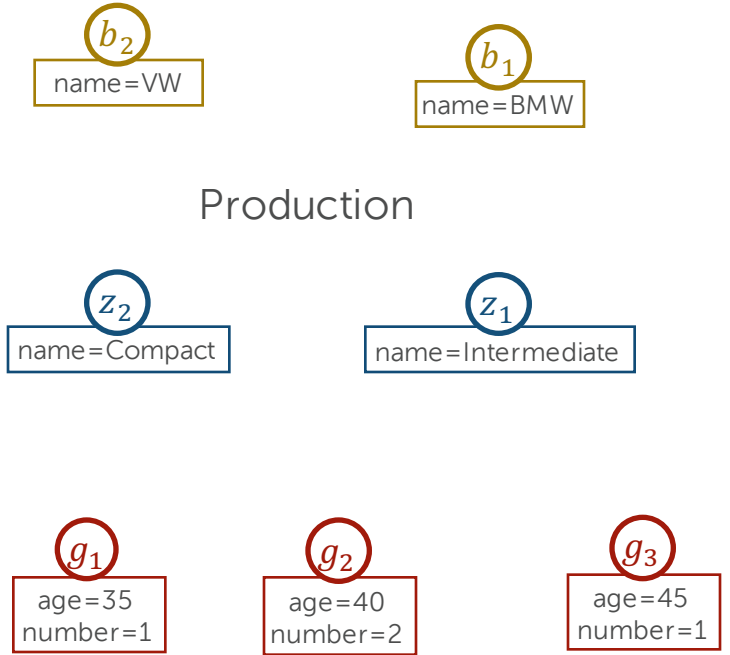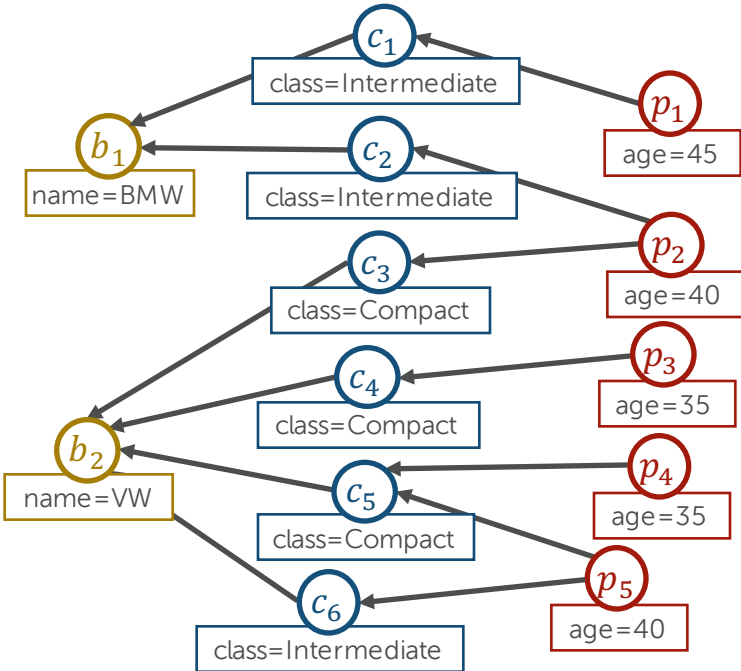
Matches

# Aggregation

$$\left( \begin{matrix} p\colon \text{Person}(\text{age} = a), \\ c\colon \text{Car}(\text{class} = k), \\ p \to c\colon \text{drives}, \\ b\colon \text{Brand}(), \\ c \to b\colon \text{madeBy} \end{matrix} \right) \to \left( \phantom{xxxx} \right)$$

# Aggregation

$$\begin{pmatrix} p: \text{Person}(\text{age} = a), \\ c: \text{Car}(\text{class} = k), \\ p \rightarrow c: \text{drives}, \\ b: \text{Brand}(), \\ c \rightarrow b: \text{madeBy} \end{pmat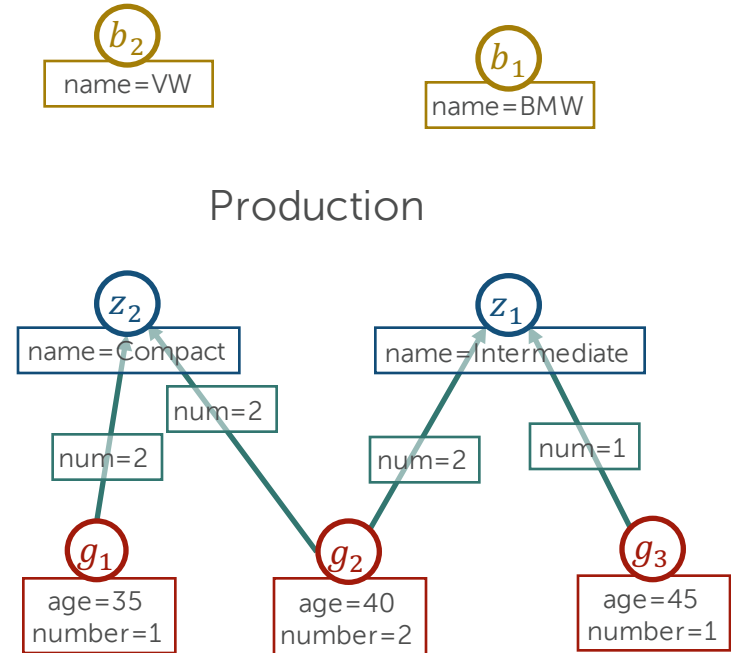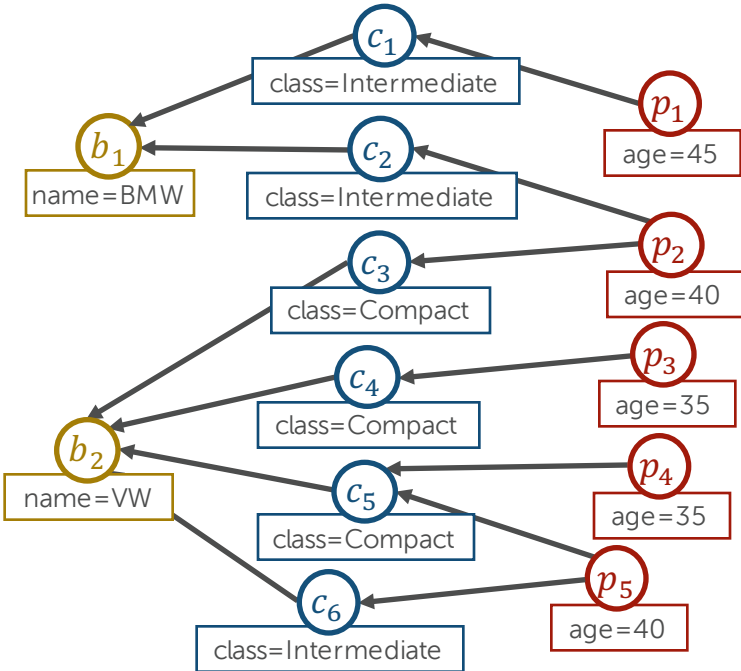rix} \rightarrow \begin{pmatrix} g@a: \text{AgeGroup}(\text{age} = a, \text{number} = \text{CNT}(p)), \\ z@k: \text{Class}(\text{name} = k), b \end{pmatrix}$$

Production

TECHNISCHE
UNIVERSITÄT
DRESDEN

# Aggregation

$$\begin{pmatrix} p: \text{Person}(age = a), \\ c: \text{Car}(class = k), \\ p \rightarrow c: \text{drives}, \\ b: \text{Brand}(), \\ c \rightarrow b: \text{madeBy} \end{pmatrix} \rightarrow \begin{pmatrix} g@a: \text{AgeGroup}(age = a, \text{number} = \text{CNT}(p)), \\ z@k: \text{Class}(\text{name} = k), b \\ g \rightarrow z(\text{numdrivers} = \text{CNT}(p)), \end{pmatrix}$$



$c_1$ class=Intermediate

$b_1$ name=BMW

$c_2$ class=Intermediate

$p_1$ age=45

$c_3$ class=Compact

$p_2$ age=40

$c_4$ class=Compact

$p_3$ age=35

$b_2$ name=VW

$c_5$ class=Compact

$p_4$ age=35

$c_6$ class=Intermediate

$p_5$ age=40

Production

$b_2$ name=VW

$b_1$ name=BMW

$z_2$ name=Compact

$z_1$ name=Intermediate

num=2

num=2

num=2

num=1

$g_1$ age=35 number=1

$g_2$ age=40 number=2

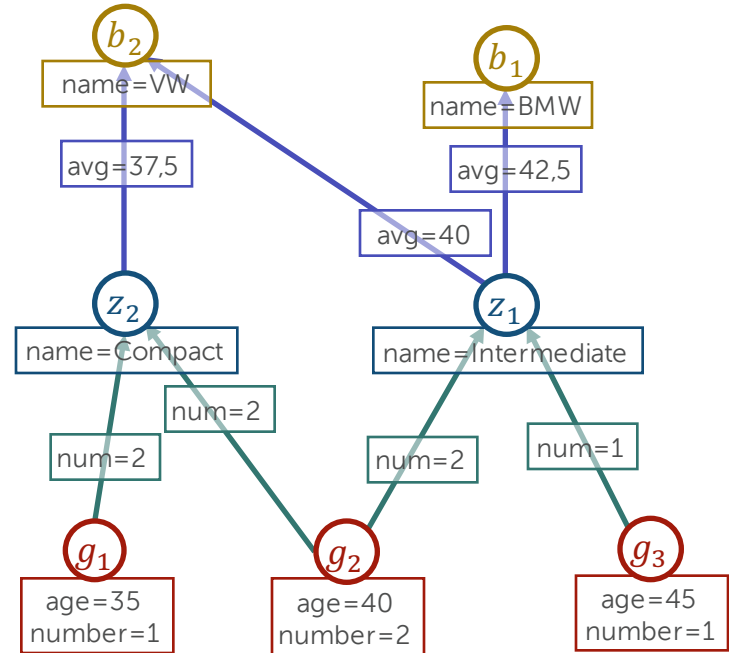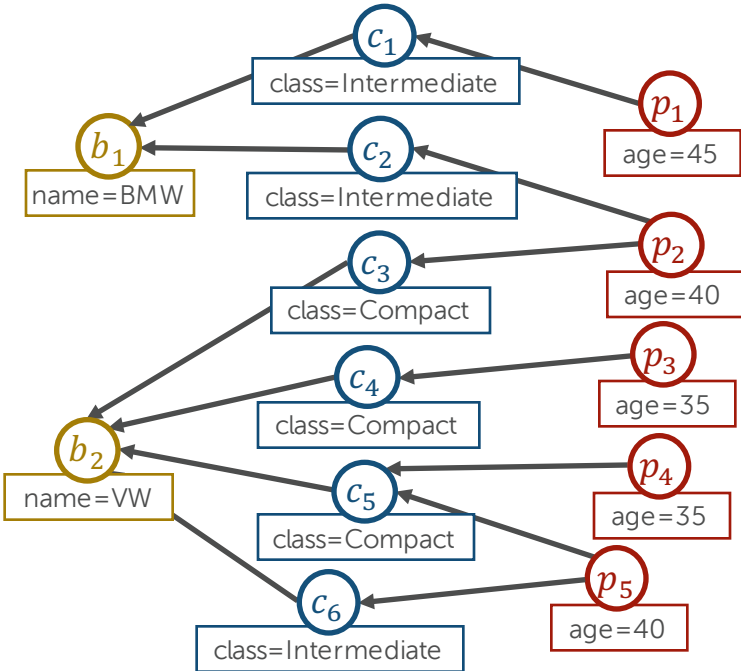$g_3$ age=45 number=1

TECHNISCHE UNIVERSITÄT DRESDEN
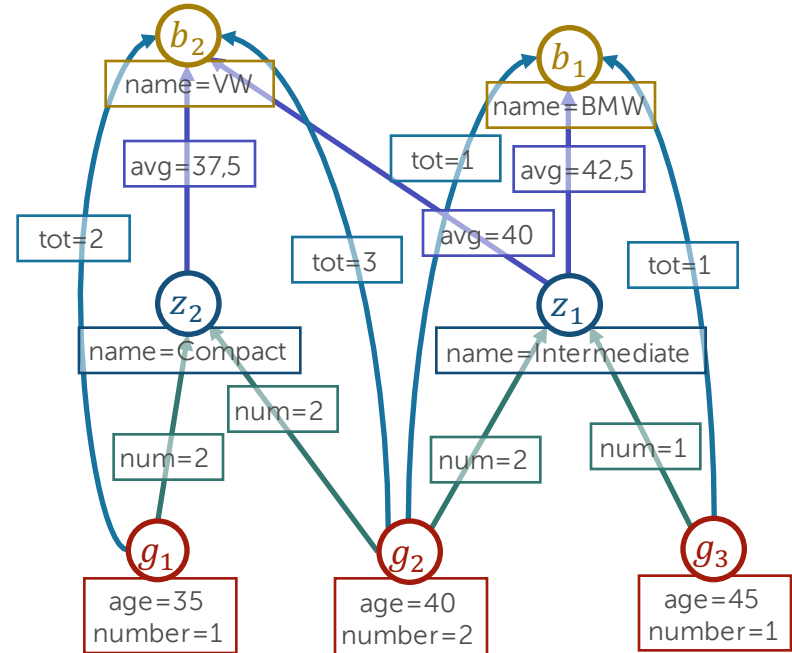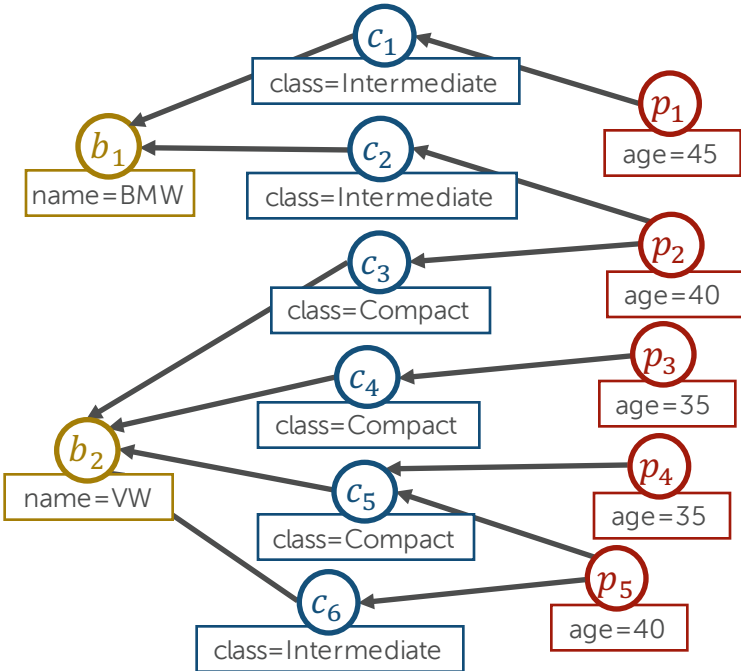
# Aggregation



$$\begin{pmatrix} p: \text{Person}(\text{age} = a), \\ c: \text{Car}(\text{class} = k), \\ p \rightarrow c: \text{drives}, \\ b: \text{Brand}(), \\ c \rightarrow b: \text{madeBy} \end{pmatrix} \rightarrow \begin{pmatrix} g@a: \text{AgeGroup}(\text{age} = a, \text{number} = \text{CNT}(p)), \\ z@k: \text{Class}(\text{name} = k), b \\ g \rightarrow z(\text{numdrivers} = \text{CNT}(p)), \\ z \rightarrow b(\text{avgage} = \text{AVG}(p.\text{age})), \end{pmatrix}$$

# Aggregation

$$\begin{pmatrix} p: \text{Person}(\text{age} = a), \\ c: \text{Car}(\text{class} = k), \\ p \to c: \text{drives}, \\ b: \text{Brand}(), \\ c \to b: \text{madeBy} \end{pmatrix} \to \begin{pmatrix} g@a: \text{AgeGroup}(\text{age} = a, \text{number} = \text{CNT}(p)), \\ z@k: \text{Class}(\text{name} = k), b \\ g \to z(\text{numdrivers} = \text{CNT}(p)), \\ z \to b(\text{avgage} = \text{AVG}(p.\text{age})), \\ g \to b(\text{totalcars} = \text{CNT}(c)) \end{pmatrix}$$

# On Tables

$$\begin{pmatrix} p: \text{Person}(\text{age} = a), \\ c: \text{Car}(\text{class} = k), \\ p \to c: \text{drives}, \\ b: \text{Brand}(), \\ c \to b: \text{madeBy} \end{pmatrix} \to \begin{pmatrix} \\ \\ \\ \\ \end{pmatrix}$$

# Grouping Lattice

# Multidimensional Graph Aggregation

$c_1$

$c_2$

$c_3$

$c_4$

$c_5$

$c_6$

$$\left( \begin{array}{c} \text{Fact} \\ c: \mathrm{Car}( \qquad ), \\ \text{-} \qquad \text{-} \qquad \text{-} \end{array} \right) \rightarrow \left( \phantom{xxxxxxxxxx} \right)$$

# Multidimensional Graph Aggregation

[https://commons.wikimedia.org/wiki/File:Rubiks_revenge_solved.jpg]

$c_1$

class=Intermediate

$c_2$

class=Intermediate

$c_3$

class=Compact

$c_4$

class=Compact

$c_5$

class=Compact

$c_6$

class=Intermediate

Fact

Dimension *Class*

$$\left( \phantom{xxxxxxxxxx} c : \mathrm{Car}(\mathrm{class} = k), \phantom{xxxxxxxxxx} \right) \rightarrow \left( \phantom{xxxxxxxxxxxxxx} \right)$$

# Multidimensional Graph Aggregation

$$\left( \begin{array}{l} \text{Fact} \\ c : \mathrm{Car}(\mathrm{class} = k), \\ p : \mathrm{Person}(\mathrm{age} = a), p \to c, \end{array} \right) \to \left( \phantom{xxx} \right)$$

Dimension *Class*

Dim. *Driver Age*

$c_1$ class=Intermediate — $p_1$ age=45

$c_2$ class=Intermediate — $p_2$ age=40

$c_3$ class=Compact

$c_4$ class=Compact — $p_3$ age=35

$c_5$ class=Compact — $p_4$ age=35

$c_6$ class=Intermediate — $p_5$ age=40

# Multidimensional Graph Aggregation

$$
\begin{array}{c}
\text{Fact} \\
\text{Dimension } \textit{Class} \\
\text{Dim. } \textit{Driver Age} \\
\text{Dim. } \textit{Brand}
\end{array}
\left(
\begin{array}{c}
c: \mathrm{Car}(\mathrm{class} = k), \\
p: \mathrm{Person}(\mathrm{age} = a), p \rightarrow c, \\
b: \mathrm{Brand}, c \rightarrow b,
\end{array}
\right) \rightarrow \left( \phantom{xxxxxxxx} \right)
$$

$p_1$  age=45

$c_1$

$b_1$  class=Intermediate

BMW

$c_2$

class=Intermediate  $p_2$  age=40

$c_3$

class=Compact  $p_3$  age=35

$c_4$

class=Compact  $p_4$  age=35

$b_2$

VW  $c_5$

class=Compact  $p_5$  age=40

$c_6$

class=Intermediate

$$\left( \begin{array}{l} c: \text{Car}(\text{class} = k), \\ p: \text{Person}(\text{age} = a), p \to c, \\ b: \text{Brand}, c \to b, \\ s: \text{City}, c \to s, l: \text{Country}, s \to l \end{array} \right) \to \left( \phantom{xx} \right)$$

Fact

Dimension *Class*

Dim. *Driver Age*

Dim. *Brand*

Dimension *Location*

[https://commons.wikimedia.org/wiki/File:Rubiks_revenge_solved.jpg]

# Multidimensional Graph Aggregation

[https://commons.wikimedia.org/wiki/File:Rubiks_revenge_solved.jpg]

$$\begin{pmatrix} \text{Fact} \\ c\colon \mathrm{Car}(\text{class} = k), \\ p\colon \mathrm{Person}(\text{age} = a), p \to c, \\ b\colon \mathrm{Brand}, c \to b, \\ s\colon \mathrm{City}, c \to s, l\colon \mathrm{Country}, s \to l \end{pmatrix} \to \begin{pmatrix} \\ (\mathrm{noCars} = \mathrm{COUNT}(c)), \\ \\ \\ \end{pmatrix}$$

Dimension *Class*

Dim. *Driver Age*

Dim. *Brand*

Dimension *Location*

Measure *noCars*

$c_1$ — class=Intermediate

$p_1$ — age=45

$b_1$ — BMW

$c_2$ — class=Intermediate

$p_2$ — age=40

$s_3$ — Paris

$l_1$ — France

$c_3$ — class=Compact

$p_3$ — age=35

$s_1$ — Berlin

$c_4$ — class=Compact

$p_4$ — age=35

$b_2$ — VW

$c_5$ — class=Compact

$l_2$ — Germany

$p_5$ — age=40

$s_2$ — Dresden

$c_6$ — class=Intermediate

TECHNISCHE UNIVERSITÄT DRESDEN

# Multidimensional Graph Aggregation

$$\begin{pmatrix} c: \mathrm{Car}(\mathrm{class} = k), \\ p: \mathrm{Person}(\mathrm{age} = a), p \to c, \\ b: \mathrm{Brand}, c \to b, \\ s: \mathrm{City}, c \to s, l: \mathrm{Country}, s \to l \end{pmatrix} \to \begin{pmatrix} \mathrm{noCars} = \mathrm{COUNT}(c), \\ g@a: \mathrm{AgeGroup}\,(\mathrm{age} = a), x \to g, \\ z@k: \mathrm{Class}(\mathrm{name} = k), x \to z, \\ b, x \to b, l, x \to l \end{pmatrix}$$

Fact

Dimension *Class*

Dim. *Driver Age*

Dim. *Brand*

Dimension *Location*

Measure *noCars*

# Multidimensional Graph Aggregation

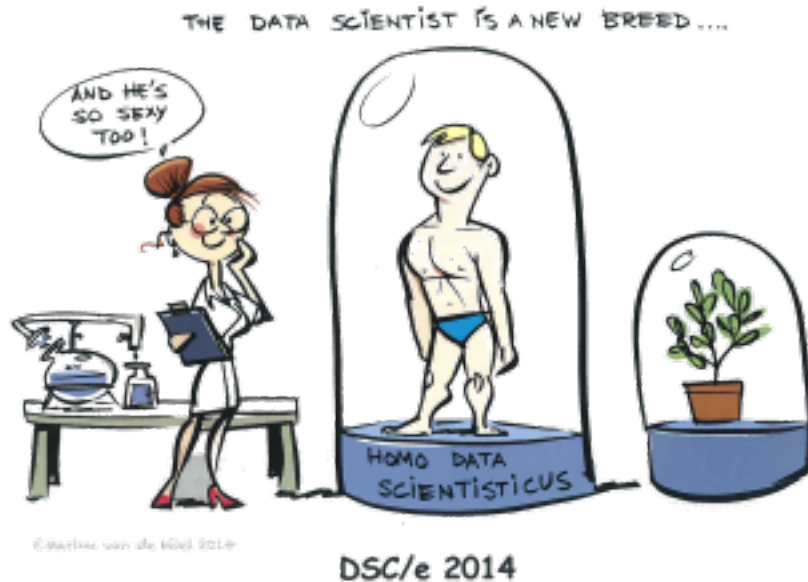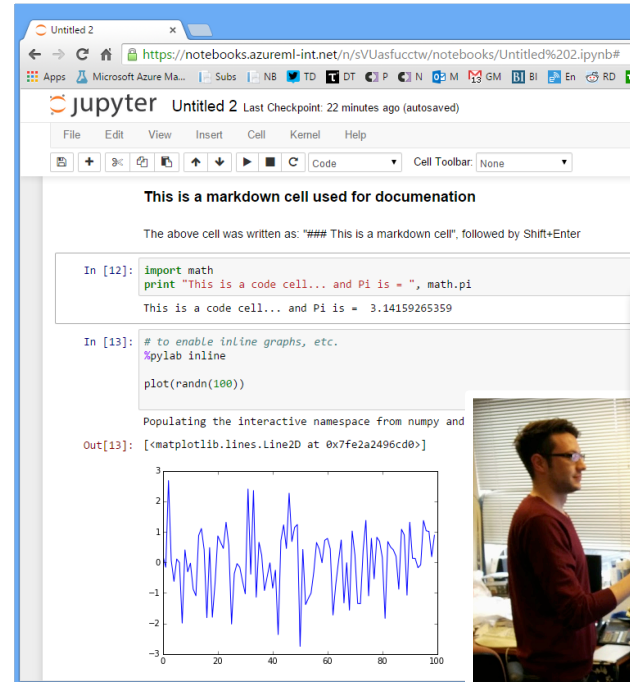# Multidimensional Graph Aggregation

# Interactive Multidimensional Graph Exploration

# Interactive Exploration



**CONSIDER A DATA SCIENTIST DOING MULTIDIMENSIONAL DATA EXPLORATION**

[http://wwwis.win.tue.nl/~wvdaalst/data_science/data_science.html]

[http://jupyter.org/]

[Crotty et al. Vizdom VLDB 2015]

# Interactive Exploration

[http://jupyter.org/]

[Crotty et al, Vizdom VLDB 2015]

# Interactive Exploration



[http://jupyter.org/]

Drilldown

Slice

Adhoc

Redefine

Dice

Incremental
State Change

[http://jupyter.org/]

[Crotty et al. Vizdom VLDB 2015]

**Drilldown**

**Slice**

**Redefine**

**Dice**

[http://jupyter.org/]

[Crotty et al. Vizdom VLDB 2015]

# SynopSys/SPARQLytics

[Michael Rudolf et al. SynopSys: Foundations for Multidimensional Graph Analytics. BIRTE 2014]

Create artifacts

## Artifacts Repository

Fact *Car:* $\langle c: \text{Car}() \rangle$

Dimension *Class:* $\langle c: \text{Car}(\text{class} = k) \rangle$

Dimension *Driver Age:* $\langle p: \text{Person}(\text{age} = a), p \rightarrow c \rangle$

Dimension *Brand:* $\langle b: \text{Brand}, c \rightarrow b \rangle$

Dimension *Location:* $\langle s: \text{City}, c \rightarrow s, l: \text{Country}, s \rightarrow l \rangle$

Measure *noCars:* $\langle \text{noCars} = \text{COUNT}(c) \rangle$

# SynopSys/SPARQLytics

Create artifacts

## Artifacts Repository

Fact *Car:* $\langle c: \mathrm{Car}()\rangle$

Dimension *Class:* $\langle c: \mathrm{Car}(\mathrm{class} = k)\rangle$

Dimension *Driver Age:* $\langle p: \mathrm{Person}(\mathrm{age} = a), p \to c\rangle$

Dimension *Brand:* $\langle b: \mathrm{Brand}, c \to b\rangle$

Dimension *Location:* $\langle s: \mathrm{City}, c \to s, l: \mathrm{Country}, s \to l\rangle$

Measure *noCars:* $\langle \mathrm{noCars} = \mathrm{COUNT}(c)\rangle$

Define cube by selecting predefined artifact

Rollup/ Drilldown/ etc.

Compute (generates the actual query)

$$\begin{pmatrix} c: \mathrm{Car}(\mathrm{class} = k), \\ p: \mathrm{Person}(\mathrm{age} = a), p \to c, \\ b: \mathrm{Brand}, c \to b, \\ s: \mathrm{City}, c \to s, l: \mathrm{Country}, s \to l \end{pmatrix} \to \begin{pmatrix} x@b, k, a, l(\mathrm{noCars} = \mathrm{COUNT}(c)), \\ g@a: \mathrm{AgeGroup}(\mathrm{age} = a), x \to g, \\ z@k: \mathrm{Class}(\mathrm{name} = k), x \to z, \\ b, x \to b, l, x \to l \end{pmatrix}$$

# SynopSys/SPARQLytics

[Michael Rudolf et al SynopSys: Foundations for Multidimensional Graph Analytics. BIRTE 2014]

Dresden Database
Systems Group

Create artifacts

Redefine cube/Define another cube

## Artifacts Repository

Fact *Car:* $\langle c: \mathrm{Car}() \rangle$

Dimension *Class:* $\langle c: \mathrm{Car}(\mathrm{class} = k) \rangle$

Dimension *Driver Age:* $\langle p: \mathrm{Person}(\mathrm{age} = a), p \rightarrow c \rangle$

Dimension *Brand:* $\langle b: \mathrm{Brand}, c \rightarrow b \rangle$

Dimension *Location:* $\langle s: \mathrm{City}, c \rightarrow s, l: \mathrm{Country}, s \rightarrow l \rangle$

Measure *noCars:* $\langle \mathrm{noCars} = \mathrm{COUNT}(c) \rangle$

Define cube
by selecting
predefined
artifact

Rollup/
Drilldown/
etc.

Compute
(generates the actual query)

$$\begin{pmatrix} c: \mathrm{Car}(\mathrm{class} = k), \\ p: \mathrm{Person}(\mathrm{age} = a), p \rightarrow c, \\ b: \mathrm{Brand}, c \rightarrow b, \\ s: \mathrm{City}, c \rightarrow s, l: \mathrm{Country}, s \rightarrow l \end{pmatrix} \rightarrow \begin{pmatrix} x@b, k, a, l(\mathrm{noCars} = \mathrm{COUNT}(c)), \\ g@a: \mathrm{AgeGroup}(\mathrm{age} = a), x \rightarrow g, \\ z@k: \mathrm{Class}(\mathrm{name} = k), x \rightarrow z, \\ b, x \rightarrow b, l, x \rightarrow l \end{pmatrix}$$

TECHNISCHE
UNIVERSITÄT
DRESDEN

# SynopSys/SPARQLytics – Experiments

## LDBC Social Network Benchmark – BI Workload

Q1 - Posting summary
Q2 - Top tags for country, age, gender, time
Q3 - Tag evolution
Q4 - Popular topics in a country
Q5 - Top posters in a country
Q6 - Most active Posters of a given Topic
Q7 - Most authoritative users on a given topic
Q8 - Related Topics
Q9 - Forum with related Tags
Q10 - Central Person for a Tag
Q11 - Unrelated Replies
Q12 - Trending Posts
Q13 - Popular Tags per month in a country
Q14 - Top thread initiators
Q15 - Social Normals
Q16 - Experts in Social Circle
Q17 - Friend Triangles
Q18 - How many persons have a given number of p
Q19 - Stranger's Interaction
Q20 - High level topics
Q21 - Zombies in a country
Q22 - International Dialog
Q23 - Holiday Destinations
Q24 - Messages By Topic And Continent

Interactive Exploration

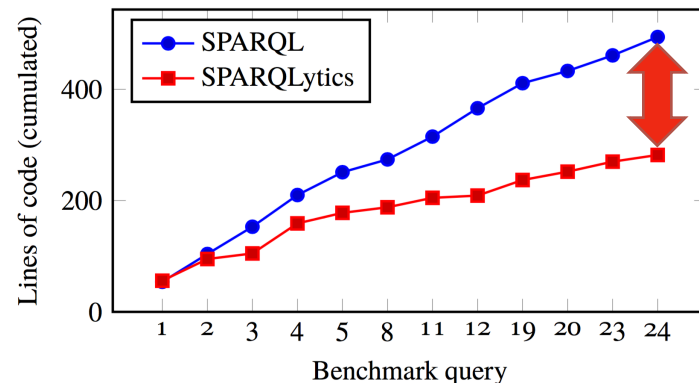Very low reuse, in practice likely to be much higher!

Unfavorable setting for our approach!

Already considerably fewer lines of code needed (factor 2)

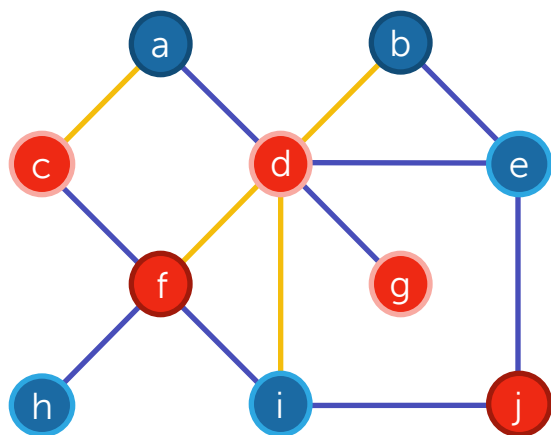| Metric | Maximum | Average |
|---|---|---|
| Re-use per dimension | 3 | 1.6 |
| Re-use per measure | 5 | 1.8 |
| Dimensions per cube | 8 | 3.5 |
| Levels per dimension | 3 | 1.5 |
| Measures per cube | 4 | 1.6 |

# Aggregating Graph Structure

# Graph Aggregation

[Peixiang Zhao et al: Graph Cube: On Warehousing and OLAP Multidimensional Networks, SIGMOD 2011]

## SUMMARIZE THE STRUCTURE OF A GRAPH IN A SMALLER GRAPH

- Group all vertices and all edge
- Represent the relationship of the groups in a graph

$$\left( \gamma_{\text{gender},\text{COUNT}(*)} V, \gamma_{\emptyset,\text{COUNT}(*)} E \right)$$



Schema: Male/Teacher   Female/Teacher   Male/Lawyer   Female/Lawyer   Co-workers   Friends
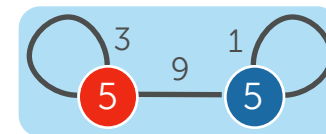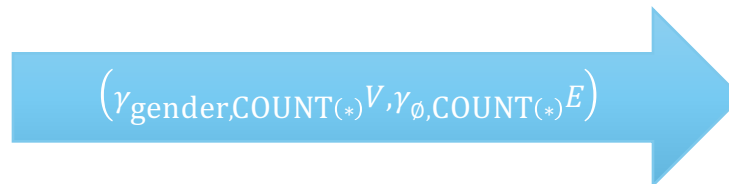
# Graph Aggregation

[Peixiang Zhao et al: Graph Cube: On Warehousing and OLAP Multidimensional Networks, SIGMOD 2011]

## SUMMARIZE THE STRUCTURE OF A GRAPH IN A SMALLER GRAPH

- Group all vertices and all edge
- Represent the relationship of the groups in a graph

$$\left(\gamma_{\text{gender},\text{COUNT}(*)}V, \gamma_{\emptyset,\text{COUNT}(*)}E\right)$$

$$\left(\gamma_{\text{gender},\text{job},\text{COUNT}(*)}V, \gamma_{\emptyset,\text{COUNT}(*)}E\right)$$



Schema:   Male/Teacher   Female/Teacher   Male/Lawyer   Female/Lawyer   Co-workers   Friends

108

# Graph Aggregation

[Peixiang Zhao et al: Graph Cube: On Warehousing and OLAP Multidimensional Networks, SIGMOD 2011]

## SUMMARIZE THE STRUCTURE OF A GRAPH IN A SMALLER GRAPH

- Group all vertices and all edge
- Represent the relationship of the groups in a graph



$$\left(\gamma_{\text{gender,COUNT}(*)}V, \gamma_{\emptyset,\text{COUNT}(*)}E\right)$$

$$\left(\gamma_{\text{gender,job,COUNT}(*)}V, \gamma_{\emptyset,\text{COUNT}(*)}E\right)$$

$$\left(\gamma_{\text{gender,COUNT}(*)}V, \gamma_{\text{status,COUNT}(*)}E\right)$$

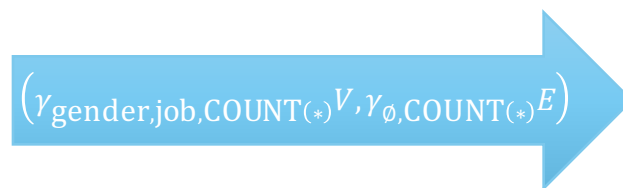Schema:   Male/Teacher   Female/Teacher   Male/Lawyer   Female/Lawyer   Co-workers   Friends
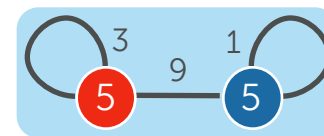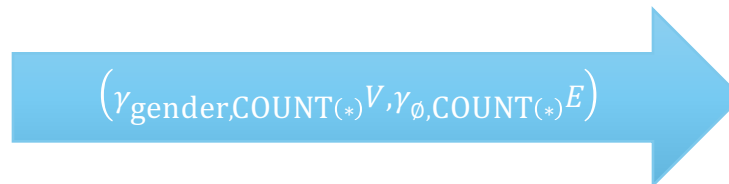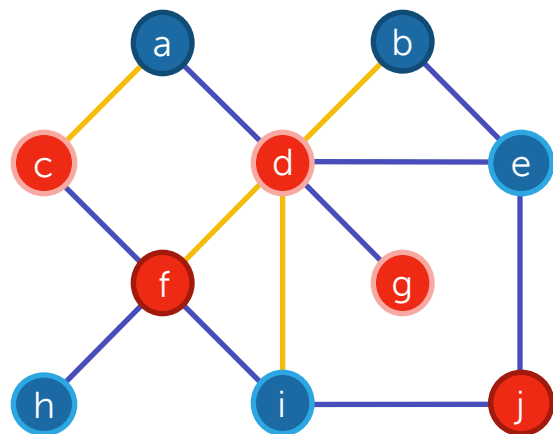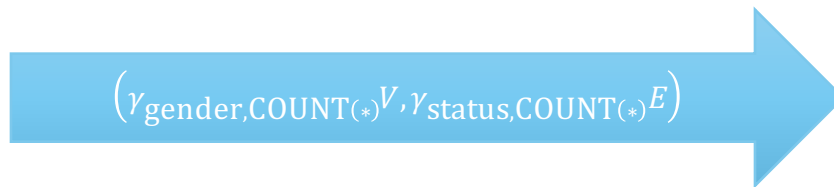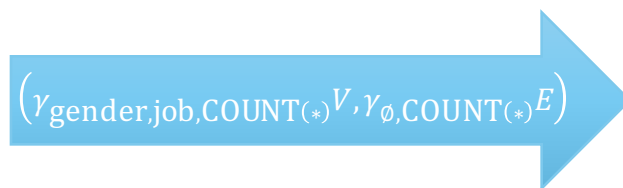
# Graph Aggregation

[Peixiang Zhao et al: Graph Cube: On Warehousing and OLAP Multidimensional Networks, SIGMOD 2011]

## SUMMARIZE THE STRUCTURE OF A GRAPH IN A SMALLER GRAPH

- Group all vertices and all edge
- Represent the relationship of the groups in a graph

$$\left(\gamma_{\text{gender,COUNT}(*)}V, \gamma_{\emptyset,\text{COUNT}(*)}E\right)$$

Generalization

$$\left(\gamma_{\text{gender,job,COUNT}(*)}V, \gamma_{\emptyset,\text{COUNT}(*)}E\right)$$

$$\left(\gamma_{\text{gender,COUNT}(*)}V, \gamma_{\text{status,COUNT}(*)}E\right)$$



Schema: Male/Teacher  Female/Teacher  Male/Lawyer  Female/Lawyer  Co-workers  Friends

110

# Graph Aggregation

## SUMMARIZE THE STRUCTURE OF A GRAPH IN A SMALLER GRAPH

- Group all vertices and all edge
- Represent the relationship of the groups in a graph

$$\left(\gamma_{\text{gender,COUNT}(*)}V, \gamma_{\emptyset,\text{COUNT}(*)}E\right)$$

$$\left(\gamma_{\text{gender,SUM}(*)}V, \gamma_{\emptyset,\text{SUM}(cnt)}E\right)$$

Generalization

$$\left(\gamma_{\text{gender,job,COUNT}(*)}V, \gamma_{\emptyset,\text{COUNT}(*)}E\right)$$

$$\left(V, \gamma_{\emptyset,\text{SUM}(cnt)}E\right)$$

$$\left(\gamma_{\text{gender,COUNT}(*)}V, \gamma_{\text{status,COUNT}(*)}E\right)$$

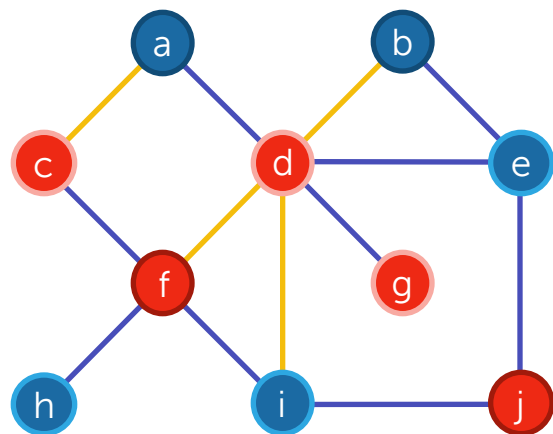Schema:   Male/Teacher   Female/Teacher   Male/Lawyer   Female/Lawyer   Co-workers   Friends

111

# Graph Aggregation

## GRAPH CUBE (CUBOID)

- Cube of all possible aggregation of a graph

- Example grouping attributes: $\{id, gender, job, status\}$

Grouping Lattice

$\emptyset$

$\{gender\}$     $\{job\}$     $\{status\}$

$\{gender, job\}$     $\{gender, status\}$     $\{job, status\}$

$\{gender, job, status\}$

$Base = \{id, gender, job, status\}$

# Graph Aggregation

## GRAPH CUBE (CUBOID)

- Cube of all possible aggregation of a graph

- Example grouping attributes: $\{id, gender, job, status\}$

Grouping Lattice

$\emptyset$

$\{gender\}$      $\{job\}$      $\{status\}$

$\{gender, job\}$    $\{gender, status\}$    $\{job, status\}$

$\{gender, job, status\}$

$Base = \{id, gender, job, status\}$

# Graph Aggregation

## GRAPH CUBE (CUBOID)

- Cube of all possible aggregation of a graph

- Example grouping attributes: $\{id, gender, job, status\}$

Grouping Lattice

# Graph Aggregation

## GRAPH CUBE (CUBOID)

- Cube of all possible aggregation of a graph

- Example grouping attributes: $\{id, gender, job, status\}$



Grouping Lattice

Graph Cube

$\emptyset$

$\{gender\}$  $\{job\}$  $\{status\}$

$\{gender, job\}$  $\{gender, status\}$  $\{job, status\}$
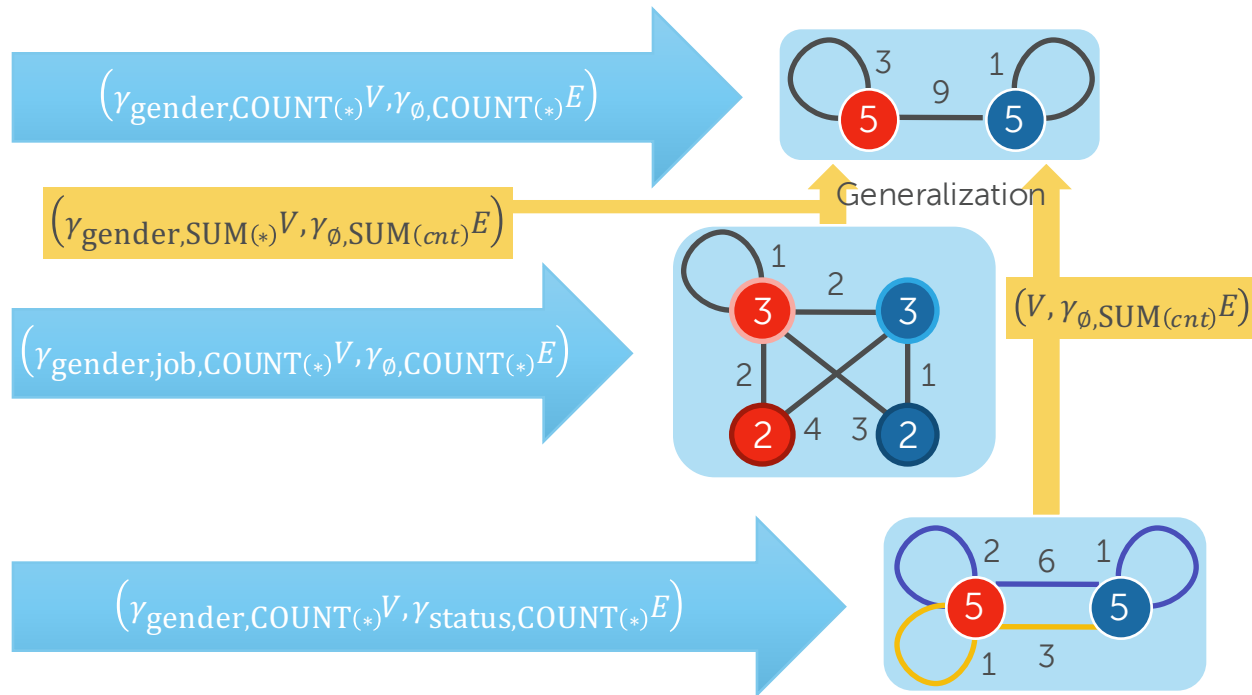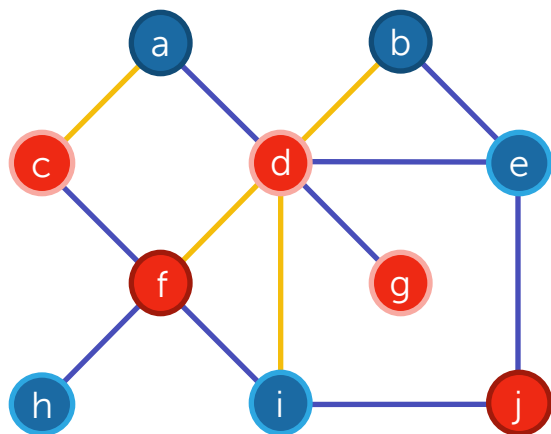
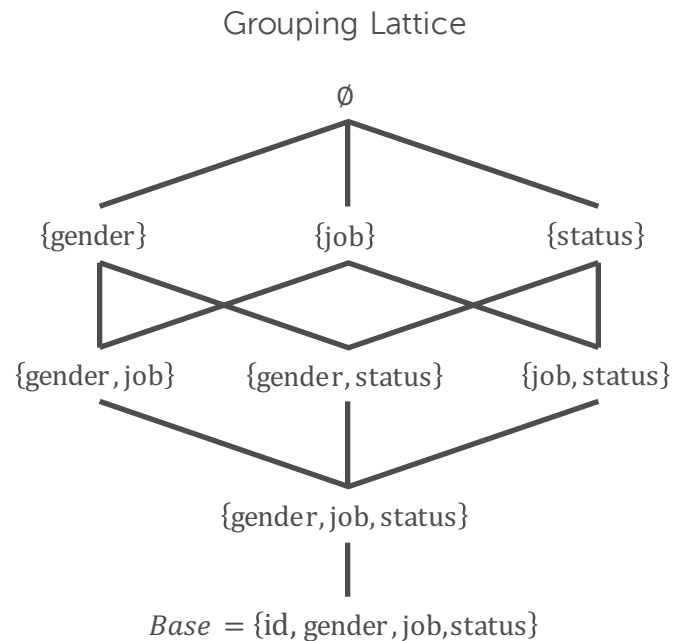$\{gender, job, status\}$

$Base = \{id, gender, job, status\}$

# Graph Aggregation

[Peixiang Zhao et al: Graph Cube: On Warehousing and OLAP Multidimensional Networks, SIGMOD 2011]

## GRAPH CUBE (CUBOID)

- Cube of all possible aggregation of a graph

- Example grouping attributes: $\{id, gender, job, status\}$



Grouping Lattice

Graph Cube

$\emptyset$

$\{gender\}$   $\{job\}$   $\{status\}$

$\{gender, job\}$   $\{gender, status\}$   $\{job, status\}$

$\{gender, job, status\}$

Cube definition only requires list of grouping attributes and measures (e.g. count) for vertices and edges

# Level of Analytics

[http://www.rosebt.com/blog/eight-levels-of-analytics-business-intelligence-to-business-analytics]



**Business Value**

**Business Analytics**

**8 OPTIMISATION** — What's the best that can happen?

**7 PREDICTIVE MODELLING** — What will happen next?

**6 FORECASTING** — What if these trends continue?

**5 STATISTICAL ANALYSIS** — Why is this happening?

**4 ALERTS** — What actions are needed?

**3 QUERY DRILLDOWN (OLAP)** — Where exactly is the problem?

**2 AD-HOC REPORTS** — How many, how often, where?

**1 STANDARD REPORTS** — What happened?

**Business Intelligence**

**Degree of Intelligence**

0. Transactional Data Management (OLTP)

batch
long running
imperative

Graph Data
Management
Applications

ad-hoc
short running
declarative

Integration with
Machine
Learning and
Optimization
Methods

Composable
Graph
Matching w/
Aggregation

# Level of Analytics

[http://www.rosebt.com/blog/eight-levels-of-analytics-business-intelligence-to-business-analytics]



Vertex-
centric
Programming

Graph Data
Management
Applications

Composable
Graph
Matching w/
Aggregation

batch
long running
imperative

ad-hoc
short running
declarative

0. Transactional Data Management (OLTP)

# Vertex-centric Programming

# Vertex-centric Programming Model

## COMPUTE - COMMUNICATE

- Based on stateless user-defined function(s)
- Collection data from adjacent vertices
- Compute new state of vertex (update)
- Send data to adjacent vertices
- A vertex can be set to inactive to vote for termination of the whole computation

Think like a vertex!



Vote to halt

Active → Inactive

Message received

- Processing terminates when all vertices are simultaneously inactive and there is no data in transit

- User has to provide Compute function
- Depending framework further function are necessary, controlling data collection and sending



Collect information → Compute → Distribute information

# Pregel/Giraph

## PREGEL

- Developed by Google
- "a scalable and fault-tolerant platform with an API that is sufficiently flexible to express arbitrary graph algorithms"
- For directed graphs with vertex and edge labels (byte strings)
- First framework using vertex-centric API
- Vertices exchange instruction messages along edges
- Bulk-Synchronous-Parallel (BSP) processing in super steps

## APACHE GIRAPH

- Open source implementation of Pregel

Euler's Seven Bridges of Königsberg (from 1735)



Pregel River

# Pregel/Giraph

[Grzegorz Malewicz et al: Pregel: a system for large-scale graph processing. SPAA 2009/PODC 2009/SIGMOD 2010]

## MESSAGE ABSTRACTION

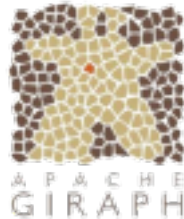- Data send between vertices, typically neighbors
- Think of message as data store on edges
- Exception
  - Combined message (cf. Combiner)
  - message not send directly, not along edge
- Receiving message = collecting data
- Sending message = sending data
- Message delivery done by framework

## SINGLE COMPUTE-FUNCTION

- Gets incoming messages as parameter
  - Think of reading data on incoming edges of current vertex
- Computes new vertex state
- Sends you new messages
  - Think of write data on outgoing edges of current vertex

Compute(incoming Message) {
        …
        Send Message
        …
}

# Pregel/Giraph

## WRITING A VERTEX-CENTRIC PROGRAM

- Subclassing the predefined `Vertex` class
- Virtual `compute()` method, which will be executed at each active vertex in every super step
- `Vertex` class provides `compute()` helper methods
- Get vertex id with `vertex_id()` and super step with `superstep()`
- Inspect the value associated with its vertex via `GetValue()` or modify it via `MutableValue()`
- Get outgoing edges with `getOutEdgeIterator()`
- Send messages to other vertices with `sendMessageTo(…)`
- Change vertex state from active to halt with `voteToHalt()`

```cpp
template <typename VertexValue,
          typename EdgeValue,
          typename MessageValue>
class Vertex {
 public:
  virtual void compute(MessageIterator* msgs) = 0;
  const string& vertex_id() const;
  int64 superstep() const;
  const VertexValue& getValue();
  VertexValue* mutableValue();
  OutEdgeIterator getOutEdgeIterator();
  void sendMessageTo(const string& dest_vertex,
                     const MessageValue& message);
  void voteToHalt();
};
```



*Vote to halt*

**Active**    **Inactive**

*Message received*

# Pregel/Giraph

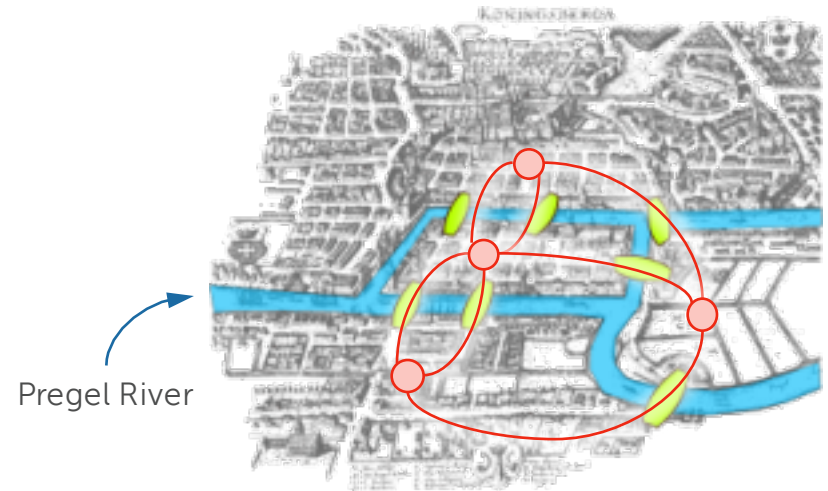[Grzegorz Malewicz et al.: Pregel: a system for large-scale graph processing. SPAA 2009/PODC 2009/SIGMOD 2010]

## COMBINERS

- User-written code
- Combines all message for a vertex $V$ into a single message
- Reduces overhead of message passing
- Enabled by subclassing the `Combiner` class and overriding a virtual `combine()` method
- No guarantees about
  - which (if any) messages are combined,
  - the groupings presented to the combiner, or
  - the order of combining
- Combiners operations should be commutative and associative operations.

## AGGREGATORS

- User-written code
- Mechanism for global communication, monitoring, and data
- Each vertex can provide a value to an aggregator in super step $S$
- Aggregator is used to combine these values to a single value
- resulting value is made available to all vertices in super step $S + 1$
- Predefined aggregators for min, max, sum , etc.
- Enabled by subclassing the `Aggregator` class
- Implementation specifies how
  - aggregated value is initialized from the first input value
  - multiple partially aggregated values are reduced to one
- Aggregator operation should be commutative and associative

# Pregel/Giraph

[Grzegorz Malewicz et al: Pregel: a system for large-scale graph processing. SPAA 2009/PODC 2009/SIGMOD 2010]

## EXAMPLE: SINGLE-SOURCE SHORTEST PATHS PROBLEM

- Finding a shortest path between a single source vertex and every other vertex in the graph
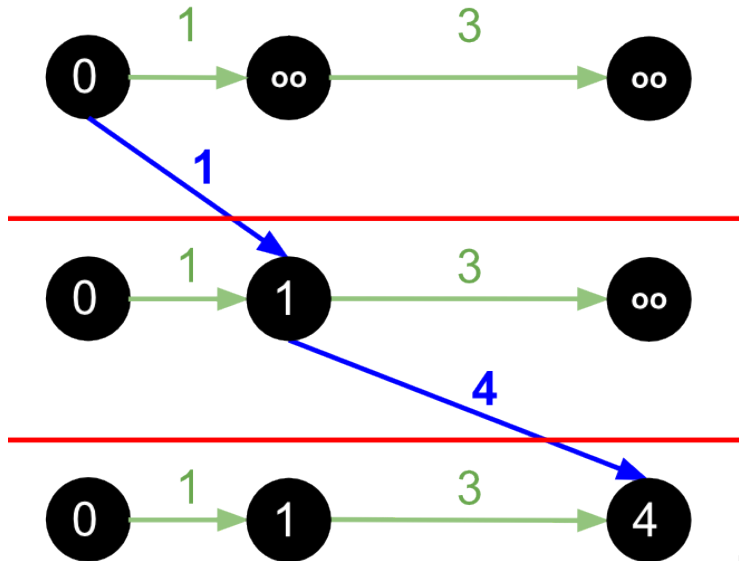
```cpp
class ShortestPathVertex : public Vertex<int, int, int> {
  void compute(MessageIterator* msgs) {
    int mindist = IsSource(vertex_id()) ? 0 : INF;    // If vertex is source shortest distance is 0
    for (; !msgs->Done(); msgs->Next())
      mindist = min(mindist, msgs->Value());           // Find shortest distance send with messages

    if (mindist < getValue()) {
      *MutableValue() = mindist;                       // If send distance is shorter that shortest distance already known, remember it
      OutEdgeIterator iter = GetOutEdgeIterator();
      for (; !iter.Done(); iter.Next())                // and for each out edge: multiply
        SendMessageTo(iter.Target(), mindist + iter.GetValue());   // own distance with edge length and send result to target vertex
    }
    VoteToHalt();
  }
};
```

Combiner:

```cpp
class MinIntCombiner : public Combiner<int> {
  virtual void combine(MessageIterator* msgs) {
    int mindist = INF;
    for (; !msgs->Done(); msgs->Next()) mindist = min(mindist, msgs->Value());
    Output("combined_source", mindist);
  }
};
```

# Pregel/Giraph

## EXAMPLE: SINGLE-SOURCE SHORTEST PATHS PROBLEM

- Finding a shortest path between a single source vertex and every other vertex in the graph



vertices with values

edges with values

messages

superstep barriers

[http://giraph.apache.org/intro.html]

126

# Pregel/Giraph

## EXAMPLE: SINGLE-SOURCE SHORTEST PATHS PROBLEM

- Finding a shortest path between a single source vertex and every other vertex in the graph

```java
public void compute(Iterable<DoubleWritable> messages) {
  double minDist = Double.MAX_VALUE;
  for (DoubleWritable message : messages) {
    minDist = Math.min(minDist, message.get());
  }
  if (minDist < getValue().get()) {
    setValue(new DoubleWritable(minDist));
    for (Edge<LongWritable, FloatWritable> edge : getEdges()) {
      double distance = minDist + edge.getValue().get();
      sendMessage(edge.getTargetVertexId(), new DoubleWritable(distance));
    }
  }
  voteToHalt();
}
```

[http://giraph.apache.org/intro.html]

127

# Pregel/Giraph
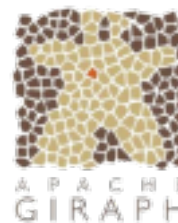
[Grzegorz Malewicz et al.: Pregel: a system for large-scale graph processing. SPAA 2009/PODC 2009/SIGMOD 2010]

## EXAMPLE: PAGE RANK

```
void compute(MessageIterator* msgs) {
  int sum = 0;
  for (; !msgs->Done(); msgs->Next())
    sum = sum + msgs->Value();
  rank = ALPHA + ((1-ALPHA)/N) * sum;
  *MutableValue() = rank;

  if (supersep() < MAX_STEPS) {
    nedges = <count number of out edges with iterator>;
    OutEdgeIterator iter = GetOutEdgeIterator();
    for (; !iter.Done(); iter.Next())
      SendMessageTo(iter.Target(), rank / nedges);
  } else {
    VoteToHalt();
  }
}
```

Sum page rank over incoming messages

Send new message over outgoing message or terminate

$$R[i] = \alpha + \frac{(1-\alpha)}{N} \sum_{(i,j) \in E} \frac{1}{L[j]} R[j]$$

# Vertex-centric Frameworks

### SCALE-OUT DISK-BASED

- Pregel
  [Grzegorz Malewicz et al: Pregel: a system for large-scale graph processing. SPAA 2009/PODC 2009/SIGMOD 2010]

- Giraph
  [http://giraph.apache.org/]

### SINGLE-BOX DISK-BASED

- GraphChi
  [Kyrola. Ligra: GraphChi: Large-Scale Graph Computation on Just a PC. OSDI 2012]

- TurboGraph
  [Han et al: TurboGraph: A Fast Parallel Graph Engine Handling Billion-scale Graphs in a Single PC. SIGKDD 2013]

### SCALE-OUT IN-MEMORY BASED

- GraphLab (asynchronous)
  [Low et al: Distributed GraphLab: A Framework for Machine Learning in the Cloud. VLDB 2012]

- PowerGraph (Gather, Sum, Apply, Scatter)
  [Gonzalez et al: PowerGraph: Distributed Graph-Parallel Computation on Natural Graphs. OSDI 2012]

- GraphX (on Apache Spark)
  [http://spark.apache.org/graphx/]

- Gelly (on Apache Flink)
  [https://ci.apache.org/projects/flink/flink-docs-master/apis/batch/libs/gelly.html]

### SHARED-MEMORY BOX

- Ligra
  [J. Shun and G. E. Blelloch. Ligra: A Lightweight Graph Processing Framework for Shared Memory. PPoPP 2013]
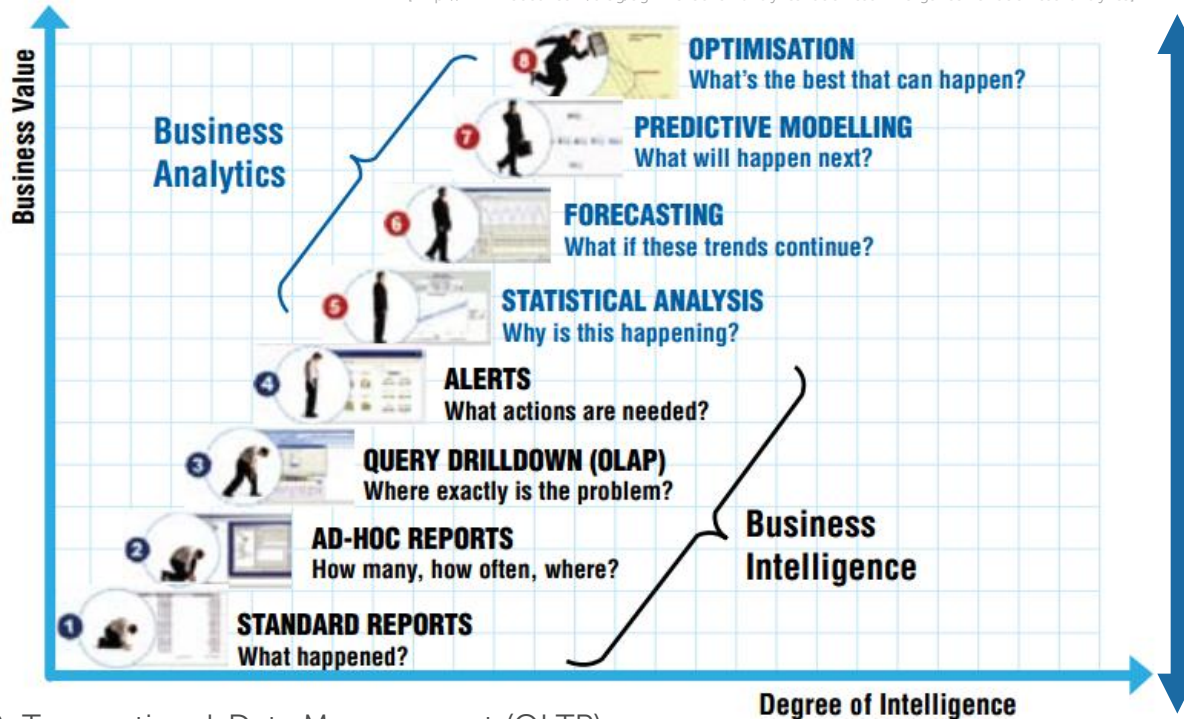
- X-Stream
  [Roy et al: Ligra: X-Stream: Edge-centric Graph Processing using Streaming Partitions. SIGOPS 2013]

- Polymer (NUMA optimized)
  [Zhang et al: NUMA-Aware Graph-Structured Analytics. PPoPP 2015]

# Level of Analytics

[http://www.rosebt.com/blog/eight-levels-of-analytics-business-intelligence-to-business-analytics]



**Business Value**

**Business Analytics**

8 **OPTIMISATION** What's the best that can happen?

7 **PREDICTIVE MODELLING** What will happen next?

6 **FORECASTING** What if these trends continue?

5 **STATISTICAL ANALYSIS** Why is this happening?

4 **ALERTS** What actions are needed?

3 **QUERY DRILLDOWN (OLAP)** Where exactly is the problem?

2 **AD-HOC REPORTS** How many, how often, where?

1 **STANDARD REPORTS** What happened?

**Business Intelligence**

**Degree of Intelligence**

0. Transactional Data Management (OLTP)

batch
long running
imperative

Graph Data
Management
Applications

ad-hoc
short running
declarative

Vertex-
centric
Programming

?

Composable
Graph
Matching w/
Aggregation

… and many more interesting topics around graphs and data analytics!!!

[http://www.bordalierinstitute.com/images/yeastProteinInteractionNetwork.jpg]

# Graph Analytics

Hannes Voigt