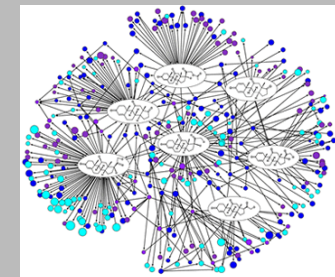


SCALABLE AND PRIVACY-PRESERVING DATA INTEGRATION - PART 3 -

ERHARD RAHM
MARTIN JUNGHANNS, ANDRÉ PETERMANN
UNIVERSITY OF LEIPZIG

www.scads.de

- ScaDS Dresden/Leipzig
- Big Data Integration
 - Scalable entity resolution / link discovery
 - Large-scale schema/ontology matching
 - Holistic data integration
- Privacy-preserving record linkage
 - Encryption of sensitive information
 - PPRL with linkage unit
 - Secure multi-party approaches
- Graph-based data integration and analytics
 - Introduction
 - Graph-based data integration / business intelligence (BIIG)
 - Hadoop-based graph analytics (GRADOOP)



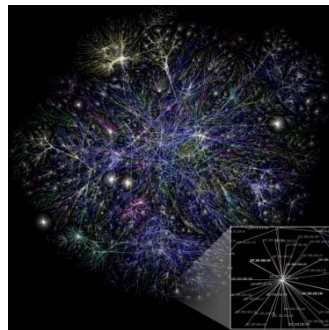


„GRAPHS ARE EVERYWHERE“

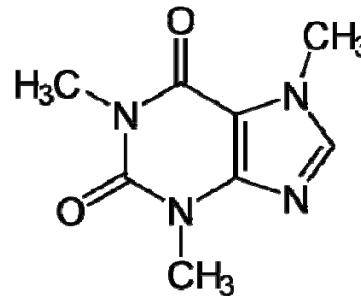
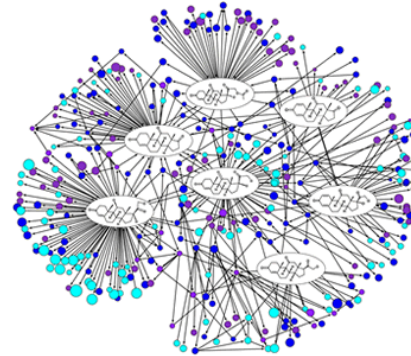
Social science



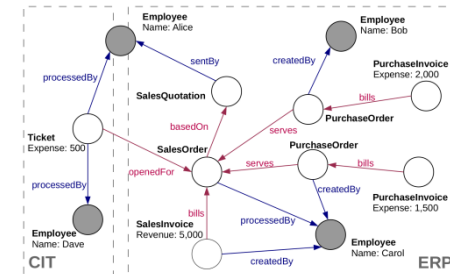
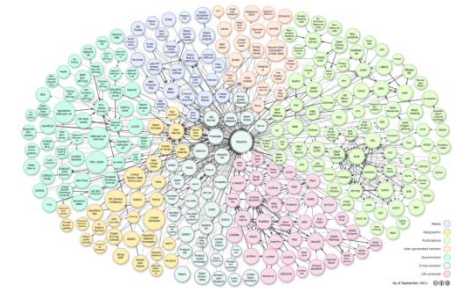
Engineering



Life science



Information science



Facebook

ca. 1.3 billion users
ca. 340 friends per user

Twitter

ca. 300 million users
ca. 500 million tweets per day

Internet

ca. 2.9 billion users

Gene (human)

20,000-25,000
ca. 4 million individuals

Patients

> 18 millions (Germany)

Illnesses

> 30.000

World Wide Web

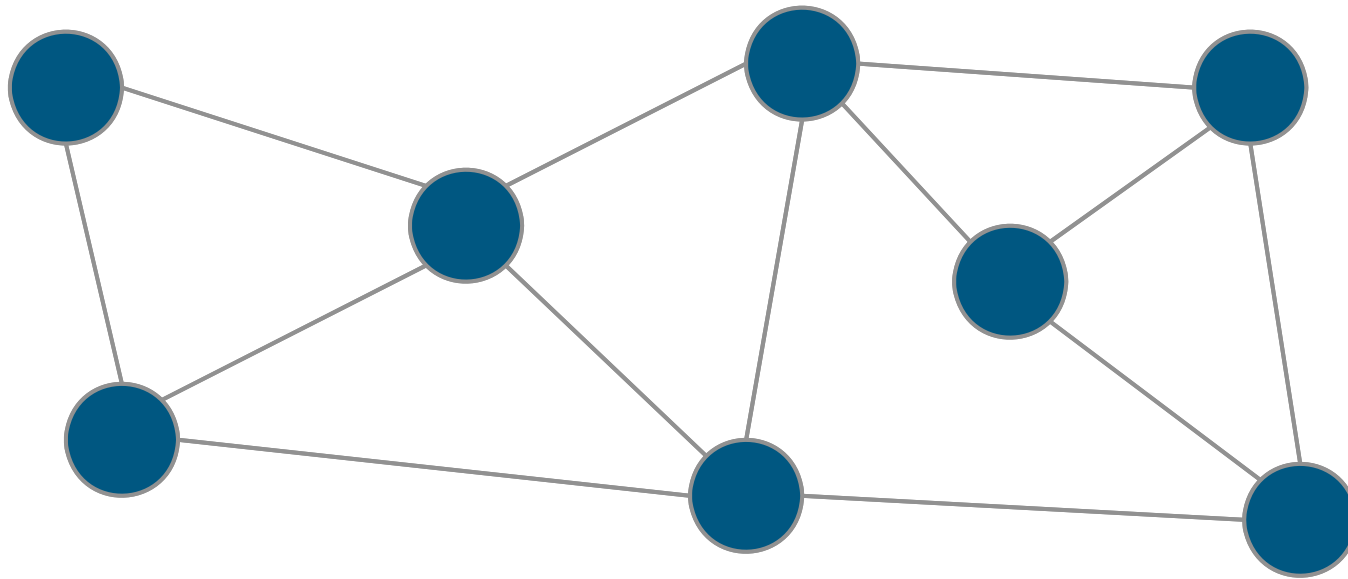
ca. 1 billion Websites

LOD-Cloud

ca. 90 billion triples



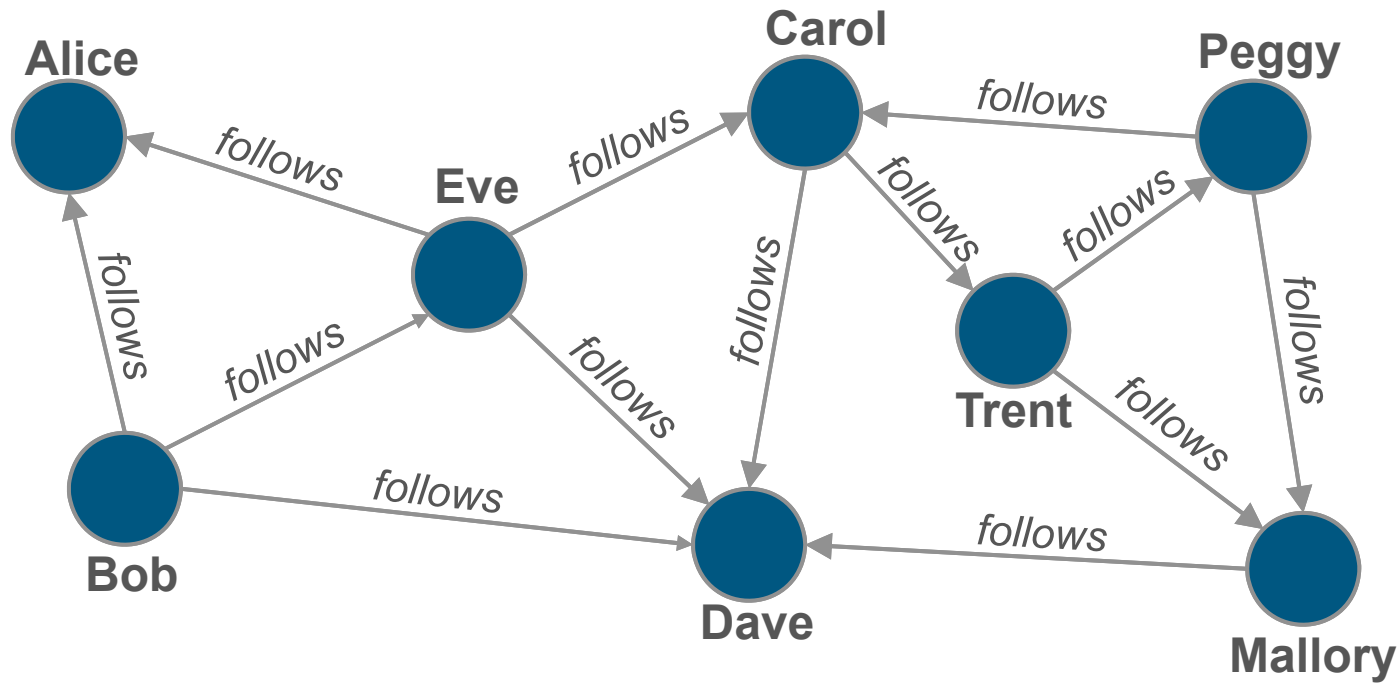
“GRAPHS ARE EVERYWHERE”



Graph = (Vertices, Edges)



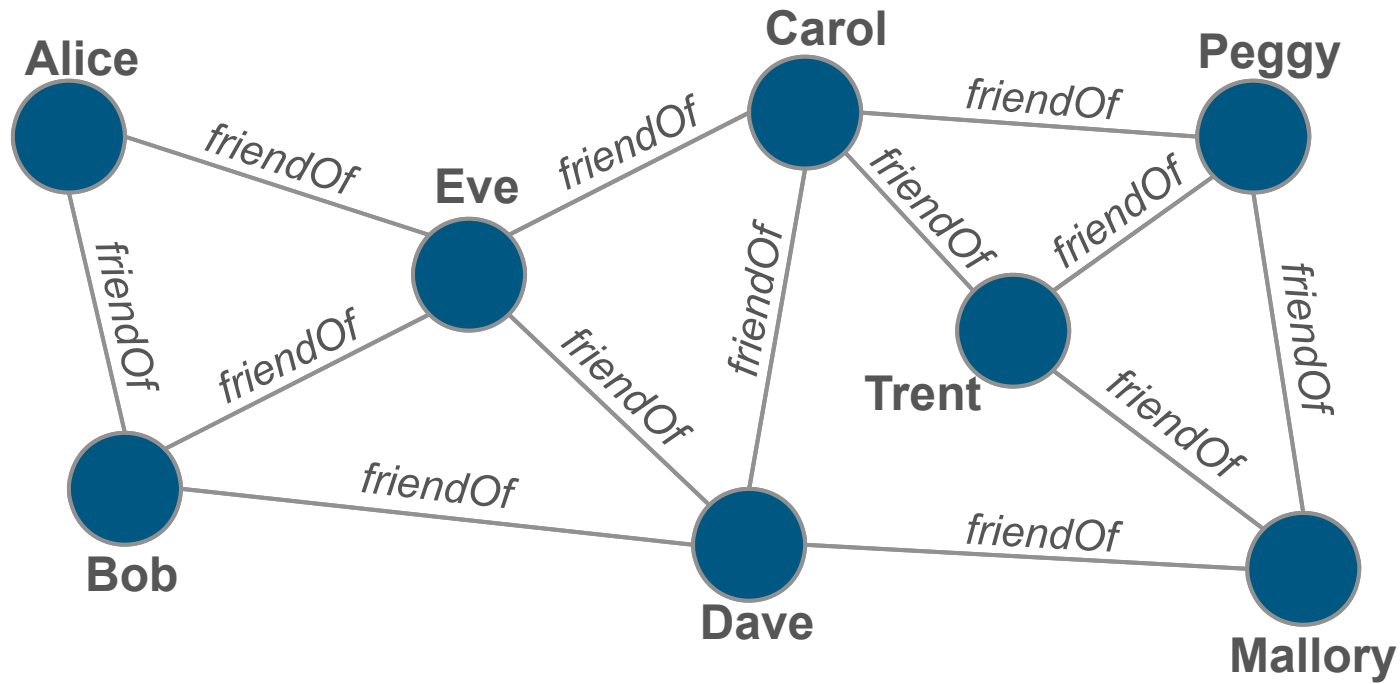
“GRAPHS ARE EVERYWHERE”



$Graph = (Users, Followers)$



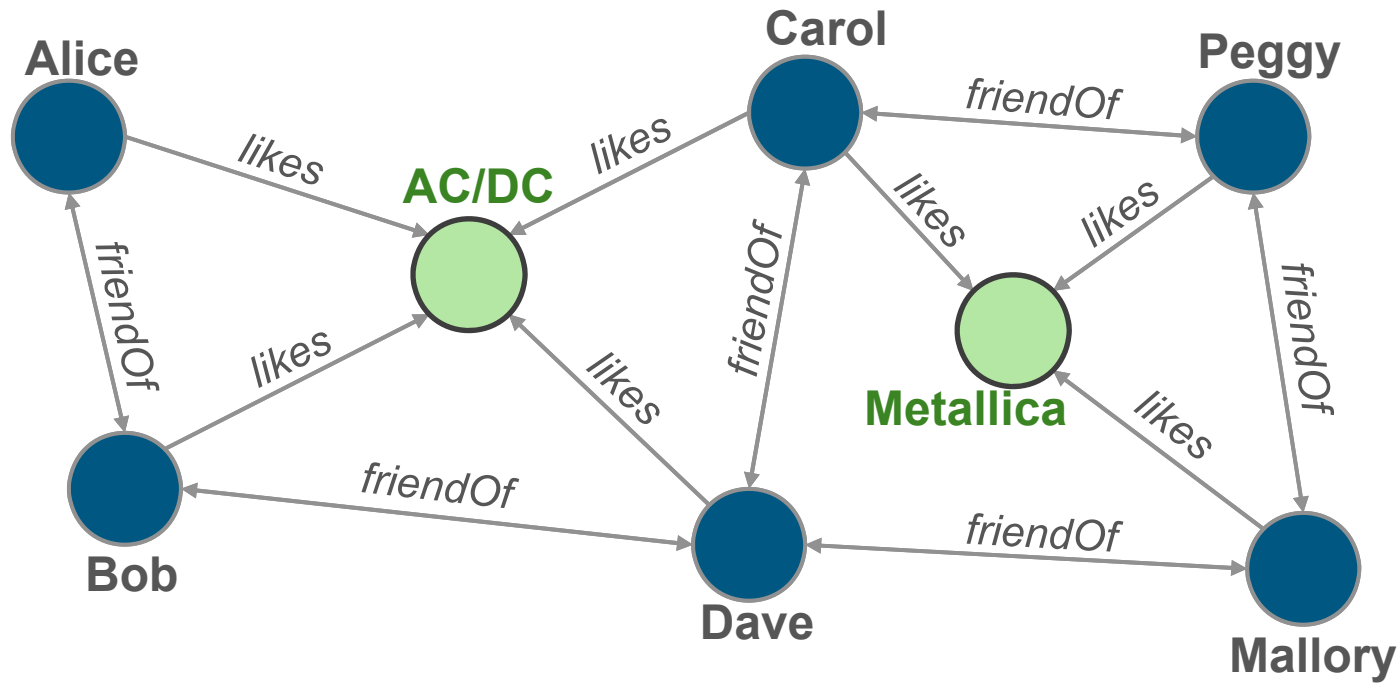
“GRAPHS ARE EVERYWHERE”



$Graph = (Users, Friendships)$

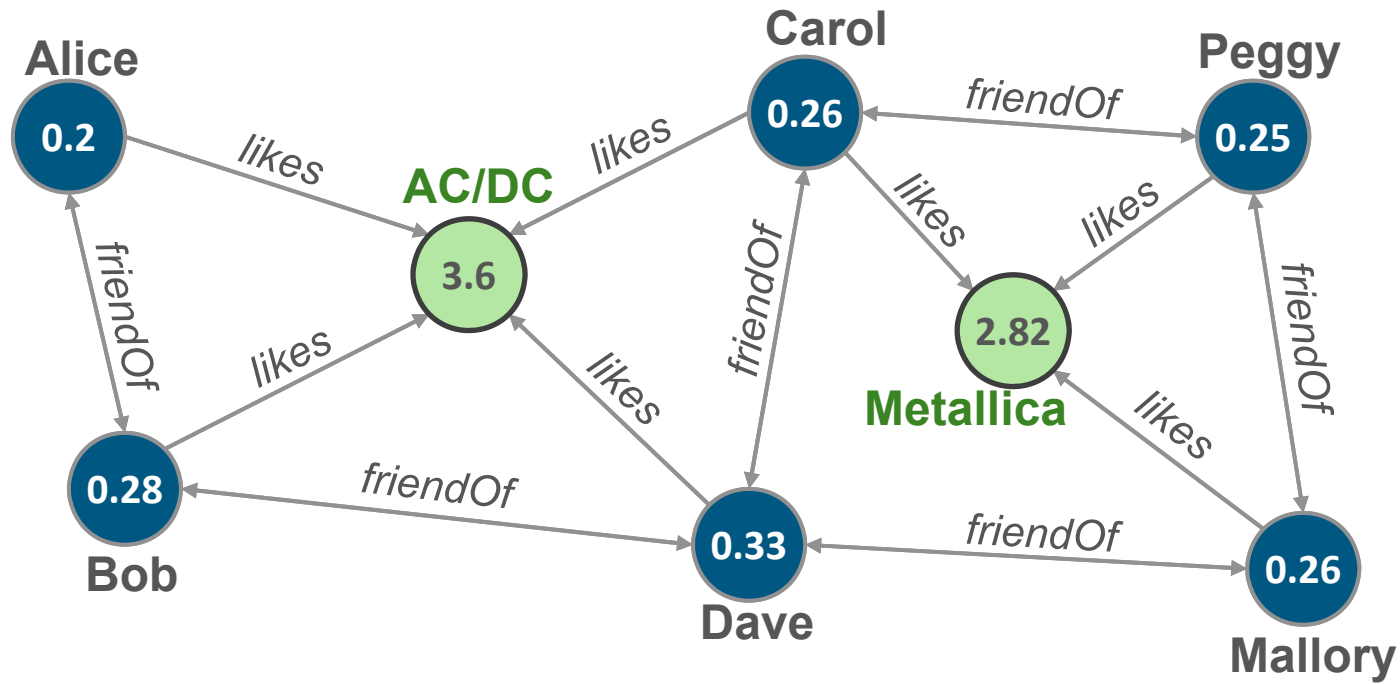


“GRAPHS ARE HETEROGENEOUS”



$$\text{Graph} = (\text{Users} \cup \text{Bands}, \text{Friendships} \cup \text{Likes})$$

“GRAPHS CAN BE ANALYZED”



$Graph = (Users \cup Bands, Friendships \cup Likes)$

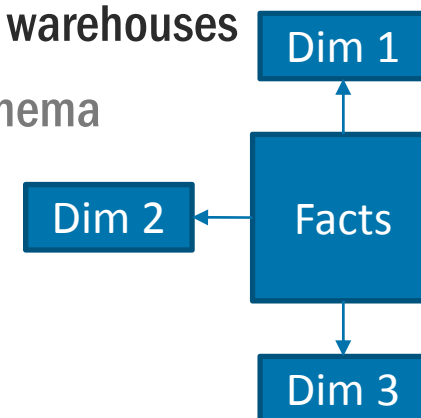
GRAPH DATA MODEL: REQUIREMENTS

- Support for heterogeneous vertices and edges
- Good semantic expressiveness
 - Typed (labeled) vertices and edges
 - Properties for vertices and edges
 - Support for collections of graphs (not only 1 graph)
- Flexibility: (semi-) structured data without strict schema
- Analysis support
 - powerful graph operators / queries
- Easy usability

Property Graph Data Model (PGM) supports most requirements

GRAPH-BASED DATA INTEGRATION

- Graphs are not only useful to represent and analyze existing networks / graph data
- Support easy linking / integration of existing data sources
 - utilized in LOD based on semantic web technology / RDF
 - can be utilized for other data models such as PGM
 - Enables graph-based analysis on relational databases and other data sources, e.g., for business intelligence
- Business intelligence usually based on relational data warehouses
 - enterprise data is integrated within dimensional schema
 - analysis limited to predefined relationships
 - no support for relationship-oriented data mining



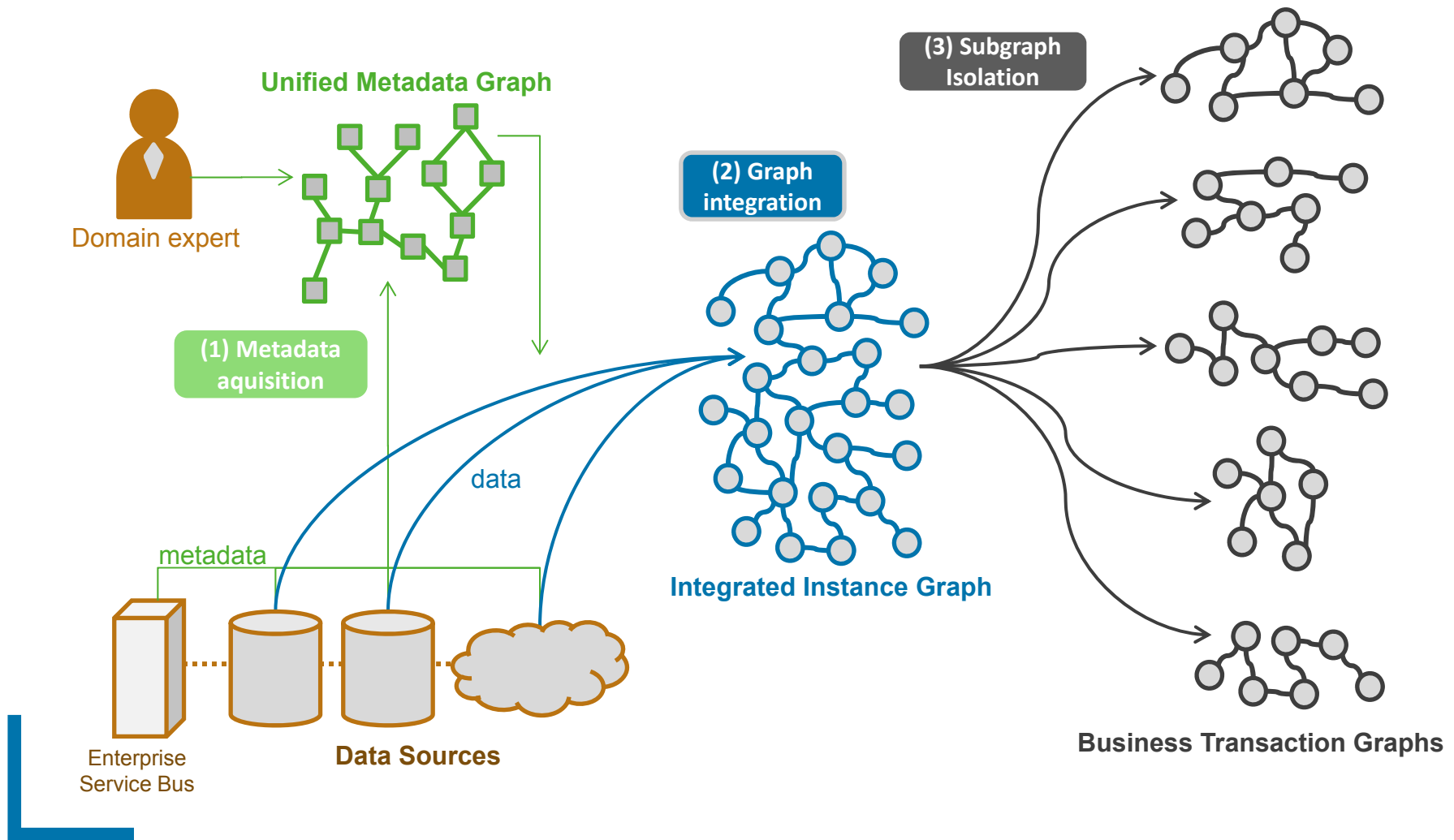
GRAPH-BASED BUSINESS INTELLIGENCE WITH BIIIG

- **BIIIG: Business Intelligence on Integrated Instance Graphs**
- **Heterogeneous data sources are integrated within an instance graph by preserving original relationships between data objects**
 - transactional and master data
- **Largely automated extraction of metadata and instance data and transformation into graphs**
 - fusion of matching entities and relations
- **Extraction of subgraphs (business transaction graphs) related to interrelated business activities**
- **Analysis of graphs/subgraphs with aggregation queries, pattern mining etc.**

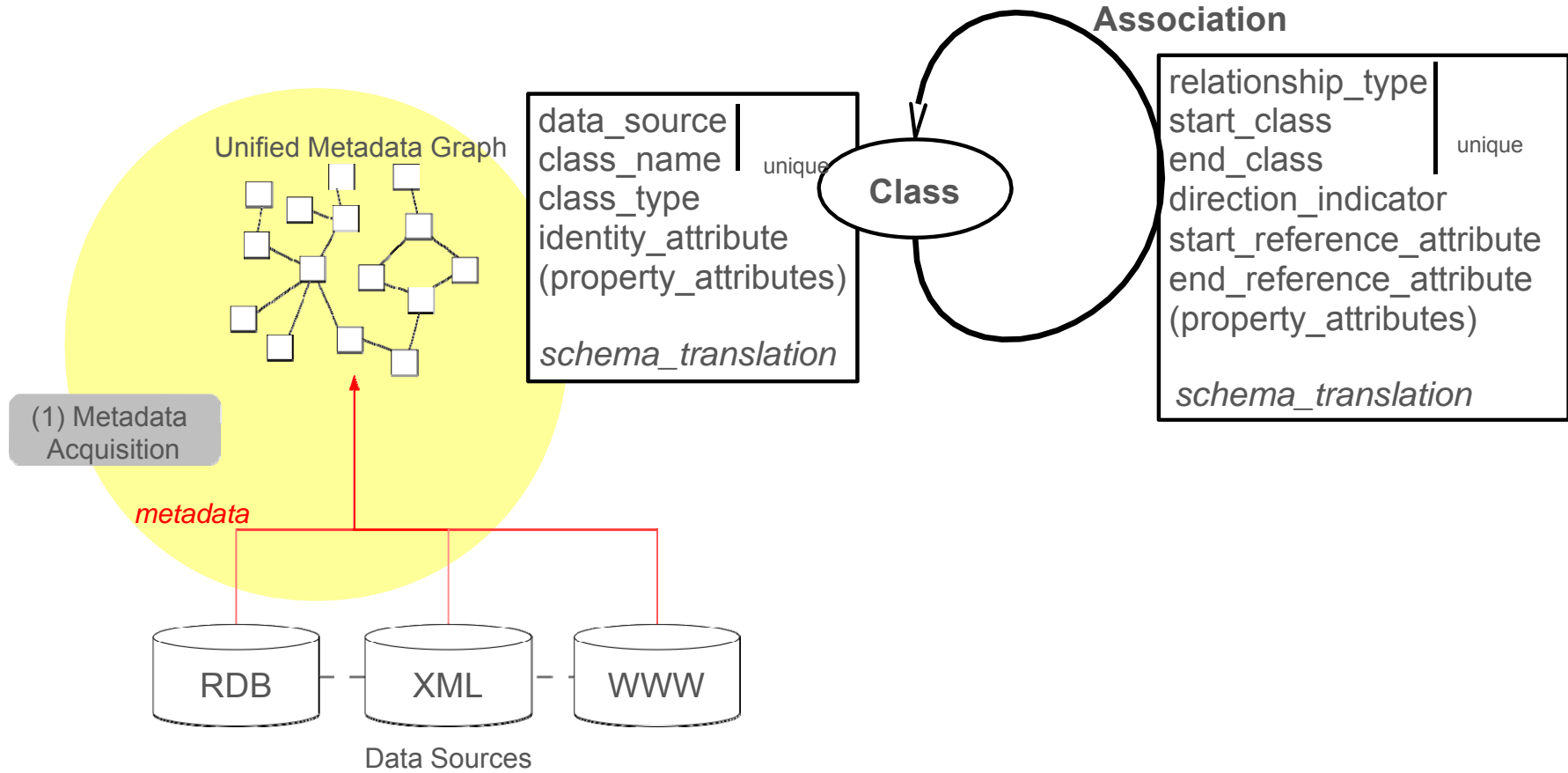


BIIG DATA INTEGRATION AND ANALYSIS WORKFLOW

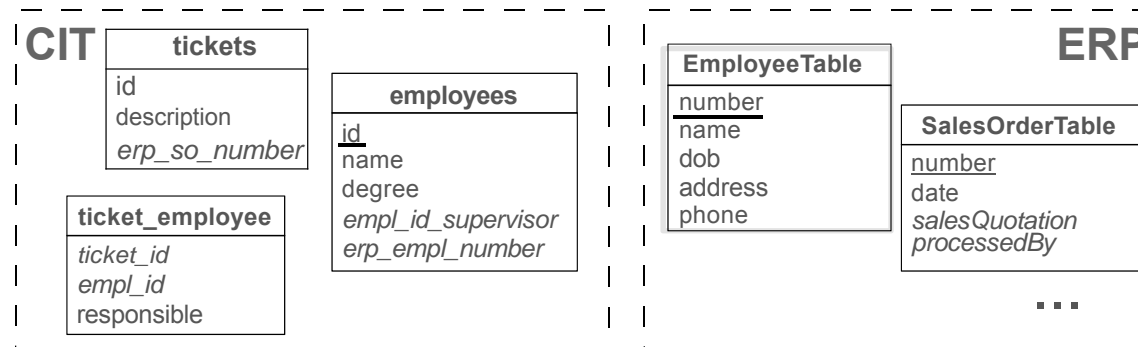
„Business Intelligence on Integrated Instance Graphs“ (PVLDB 2014)



METADATA REPRESENTATION WITH UMG



METADATA EXTRACTION (2)



Schema translations:

Class **Employee**

```
SELECT number, name,
       dob, address, phone
FROM CIT.employees
```

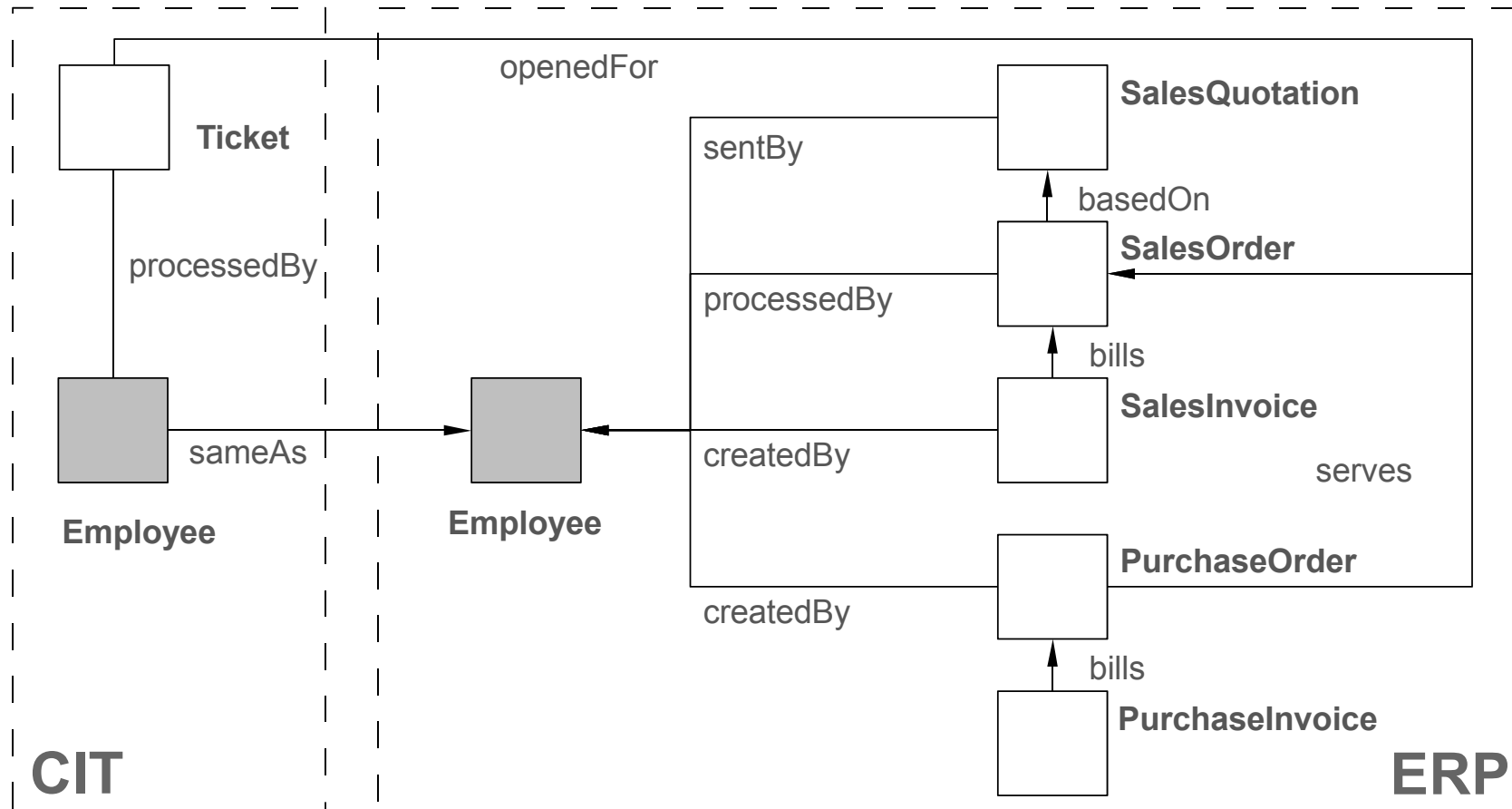
Association **processedBy** (m:n)

```
SELECT ticket_id, empl_id
FROM CIT.ticket_employee
```

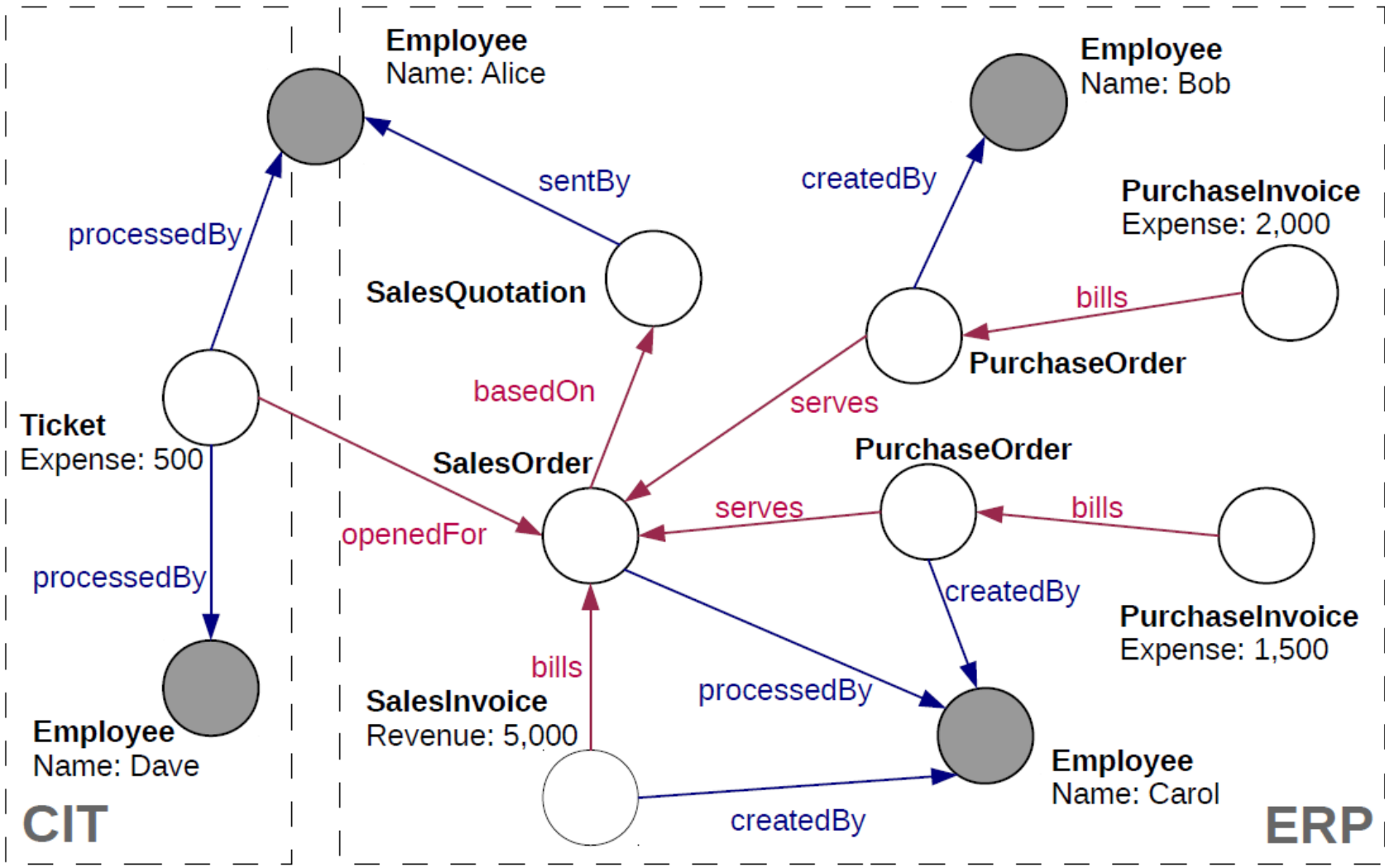
Class **Ticket**

```
SELECT id, description
FROM CIT.tickets
```

RESULTING UMG (UNIFIED METADATA GRAPH)



ScaDS  **SAMPLE INSTANCE GRAPH**
 DRESDEN LEIPZIG



SCREENSHOT OF NEO4J IMPLEMENTATION

Neo4j data/btg.db - Mozilla Firefox

Firefox Neo4j data/btg.db
localhost:7474/browser/

CYPHER MATCH (v) WHERE v.BTG_ID = 73 RETURN v

Employee [511]

Properties

- name: Leota Albery
- email: leota.albery@foodbroker.org
- gender: f
- BTG_ID: 73
- SOURCE_IDS: ["ERP_Employee_EMP1551714", "CIT_User_leota.albery@foodbroker.org"]
- TYPE: MasterData
- CLASS: Employee

Displaying 29 nodes, 41 relationships

- Relational database systems
 - can be used to implement a graph store
 - SQL alone insufficient for graph processing (need for graph operators and graph mining)
- RDF data management
 - flexible management of semantic web data
 - Data integration support (linking of entities / concepts)
 - SPARQL query processing
 - insufficient scalability of triple stores
 - insufficient support for graph mining
- Graph database system, e.g. Neo4J
 - use of *property graph data model*: vertices and edges have arbitrary set of properties (represented as key-value pairs)
 - focus on simple transactions and queries
 - insufficient scalability
 - insufficient support for graph mining



- Parallel *graph processing* systems, e.g., Google Pregel, Apache Giraph
 - in-memory storage of graphs in shared nothing clusters
 - parallel processing of general graph algorithms, e.g., page rank, connected components, ...
 - little support for semantically expressive graphs
 - no end-to-end approach with data integration and persistent graph storage
- Newer approaches (Apache Spark, Apache Flink):
 - analysis workflow with graph operators
 - no end-to-end solution with data integration and persistent graph storage



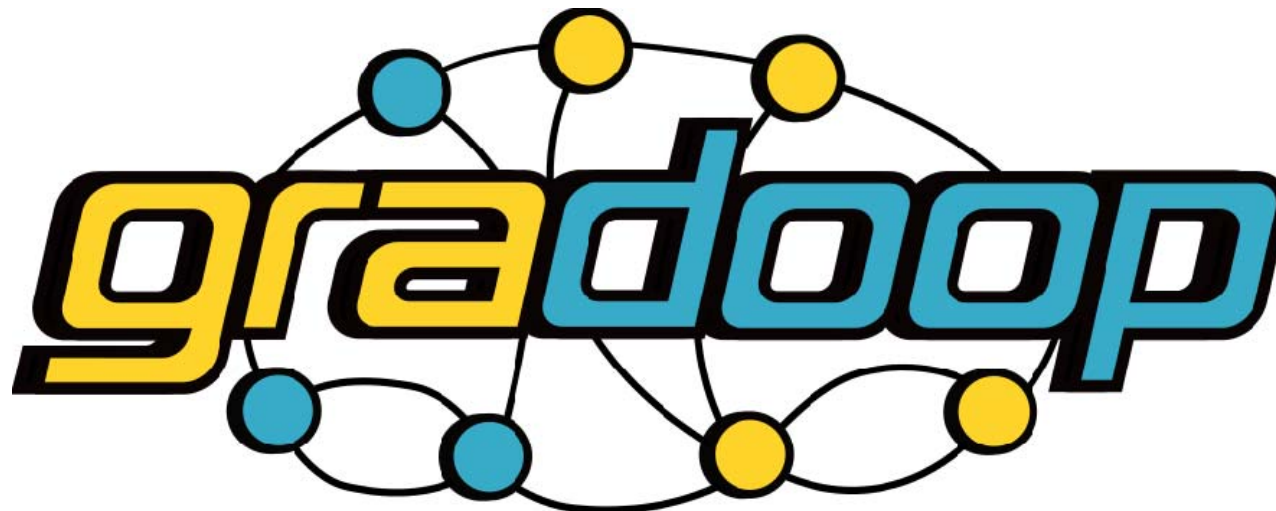
ScaDS COMPARISON

DRESDEN LEIPZIG

	Graph Database Systems Neo4j, OrientDB	Graph Processing Systems Pregel, Giraph	Distributed Dataflow Systems Flink Gelly, Spark GraphX
Data Model	Rich Graph Models (PGM)	Generic Graph Models	Generic Graph Models
Focus	transactional	analytic	analytic
Query Language	Yes	No	No
Persistency	Yes	No	No
Scalability	Vertical	Horizontal	Horizontal
Workflows	No	No	Yes
Data Integration	No	No	No
Graph Analytics	No	Yes	Yes
Visualization	Yes	No	No

WHAT'S MISSING?

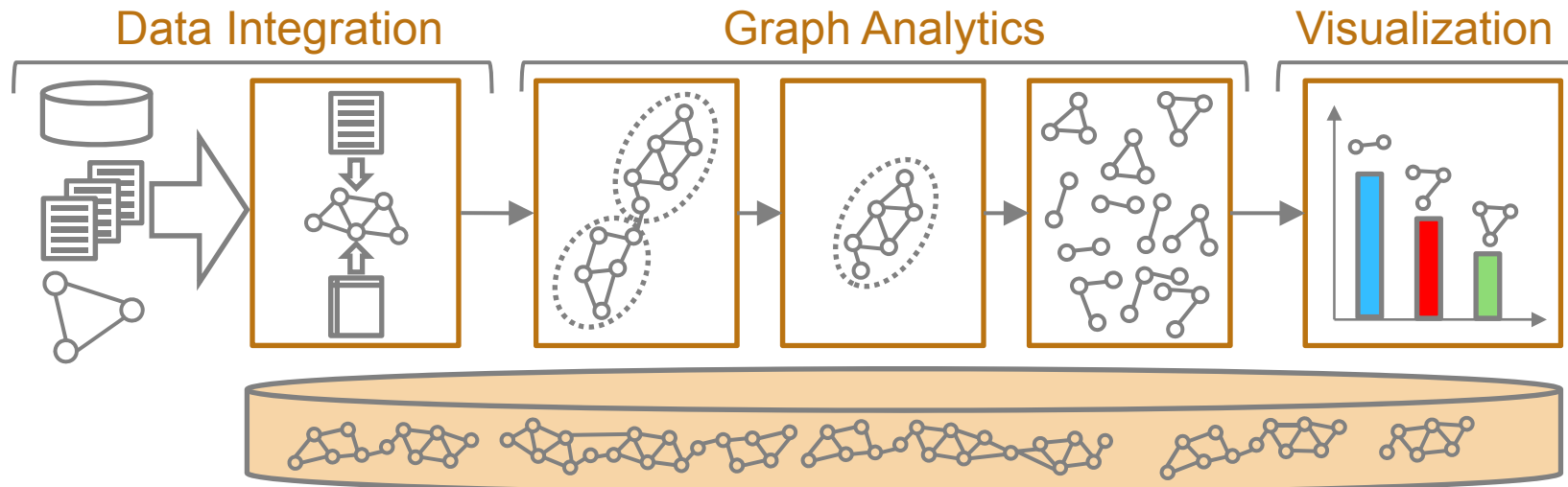
An end-to-end framework and research platform for efficient, distributed and domain independent graph data management and analytics.



GRADOOP CHARACTERISTICS

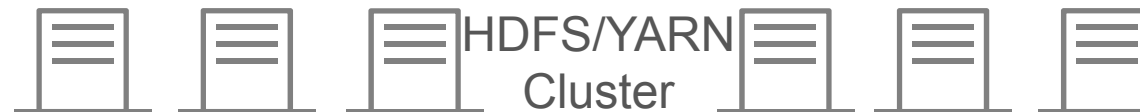
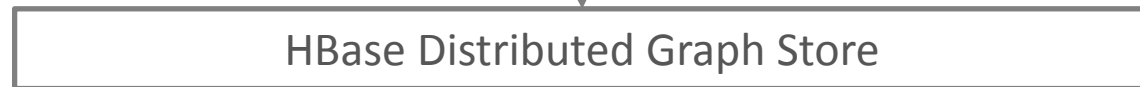
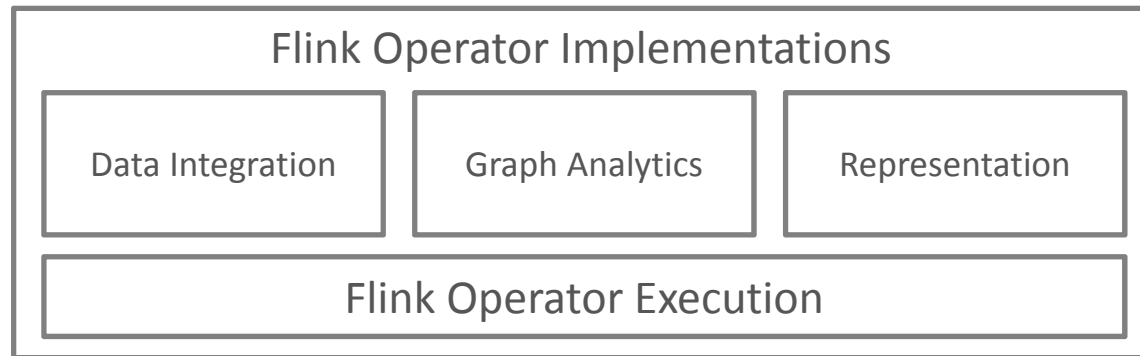
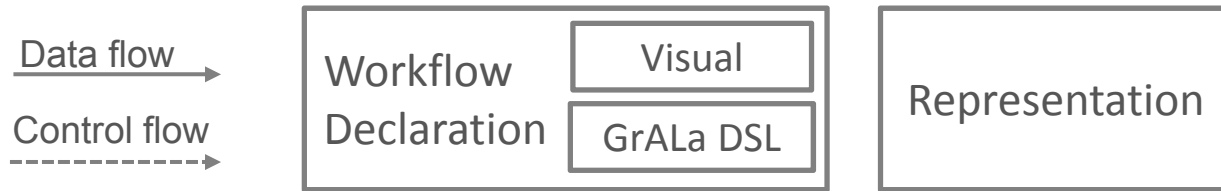
- Hadoop-based framework for graph data management and analysis
- Graph storage in scalable distributed store, e.g., HBase
- Extended property graph data model (EPGM)
 - operators on graphs and sets of (sub) graphs
 - support for semantic graph queries and mining
- Leverages powerful components of Hadoop ecosystem
 - initially: Apache HBase, MapReduce, Apache Giraph
 - now: Apache HBase, Apache Flink, ...
- New functionality for graph-based processing workflows and graph mining

END-TO-END GRAPH ANALYTICS



- **integrate data from one or more sources into a dedicated graph store with common graph data model**
- **definition of analytical workflows from operator algebra**
- **result representation in meaningful way**

ScaDS  HIGH LEVEL ARCHITECTURE
DRESDEN LEIPZIG



EXTENDED PROPERTY GRAPH MODEL (EPGM)

- Includes PGM as special case
- Support for collections of logical graphs / subgraphs
 - can be defined explicitly
 - can be result of graph algorithms / operators
- Support for graph properties
- Operators on both graphs and graph collections



EPGM – GRAPH REPRESENTATION

[0] Tag
name : Databases

[1] Tag
name : Graphs

[2] Tag
name : Hadoop

[3] Forum
title : Graph Databases

[4] Forum
title : Graph Processing

[5] Person
name : Alice
gender : f
city : Leipzig
age : 23

[6] Person
name : Bob
gender : m
city : Leipzig
age : 30

[7] Person
name : Carol
gender : f
city : Dresden
age : 30

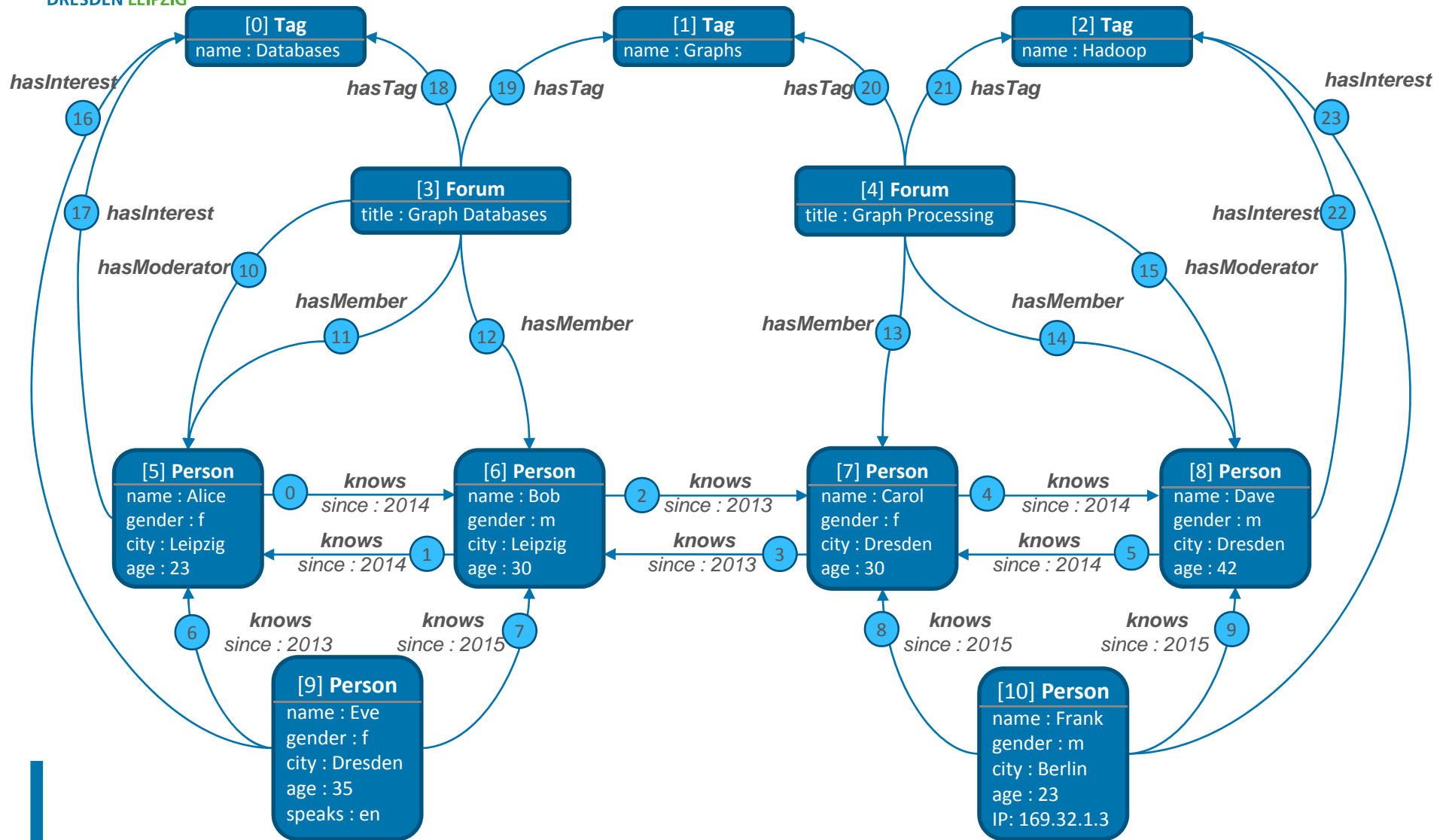
[8] Person
name : Dave
gender : m
city : Dresden
age : 42

[9] Person
name : Eve
gender : f
city : Dresden
age : 35
speaks : en

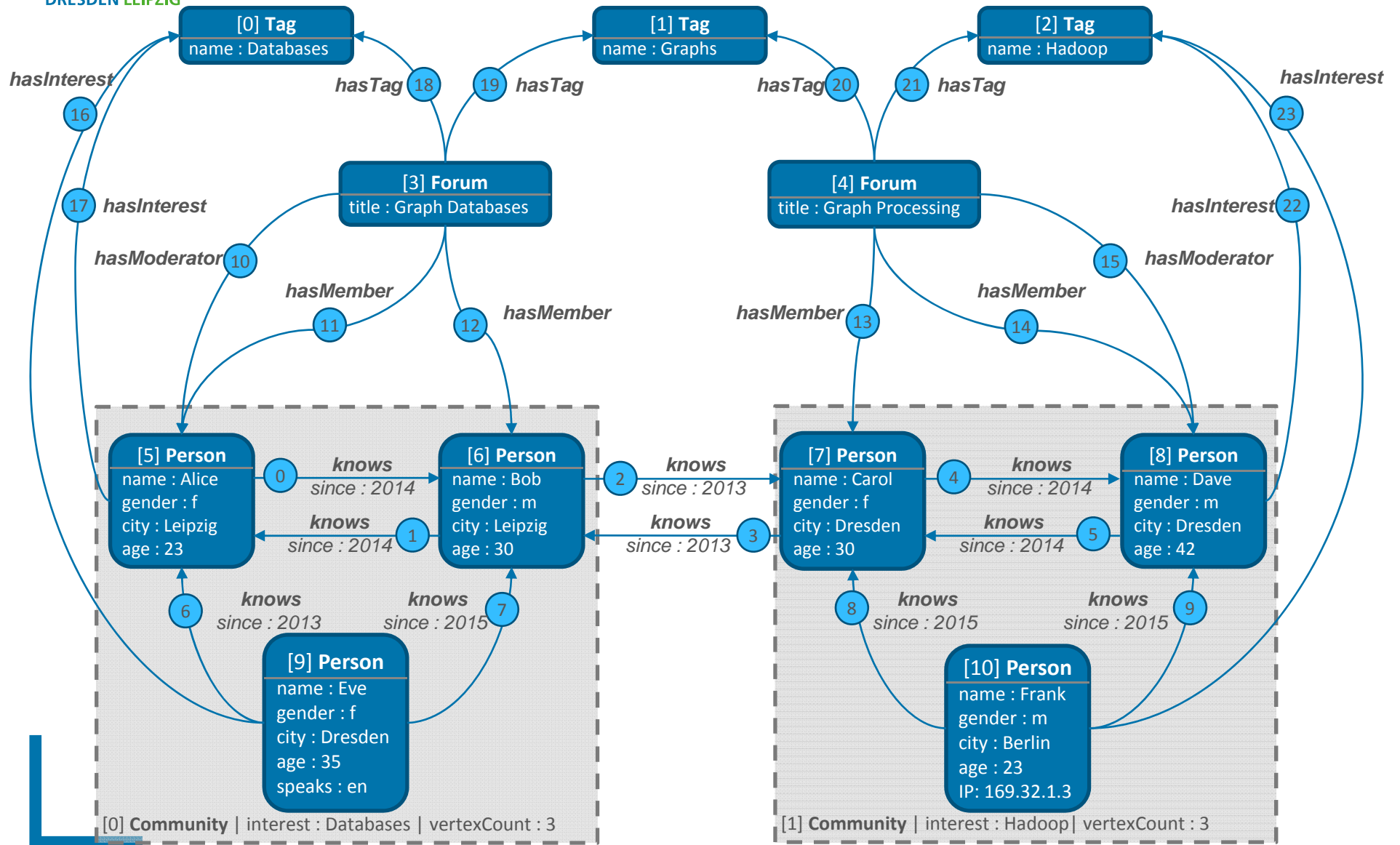
[10] Person
name : Frank
gender : m
city : Berlin
age : 23
IP: 169.32.1.3



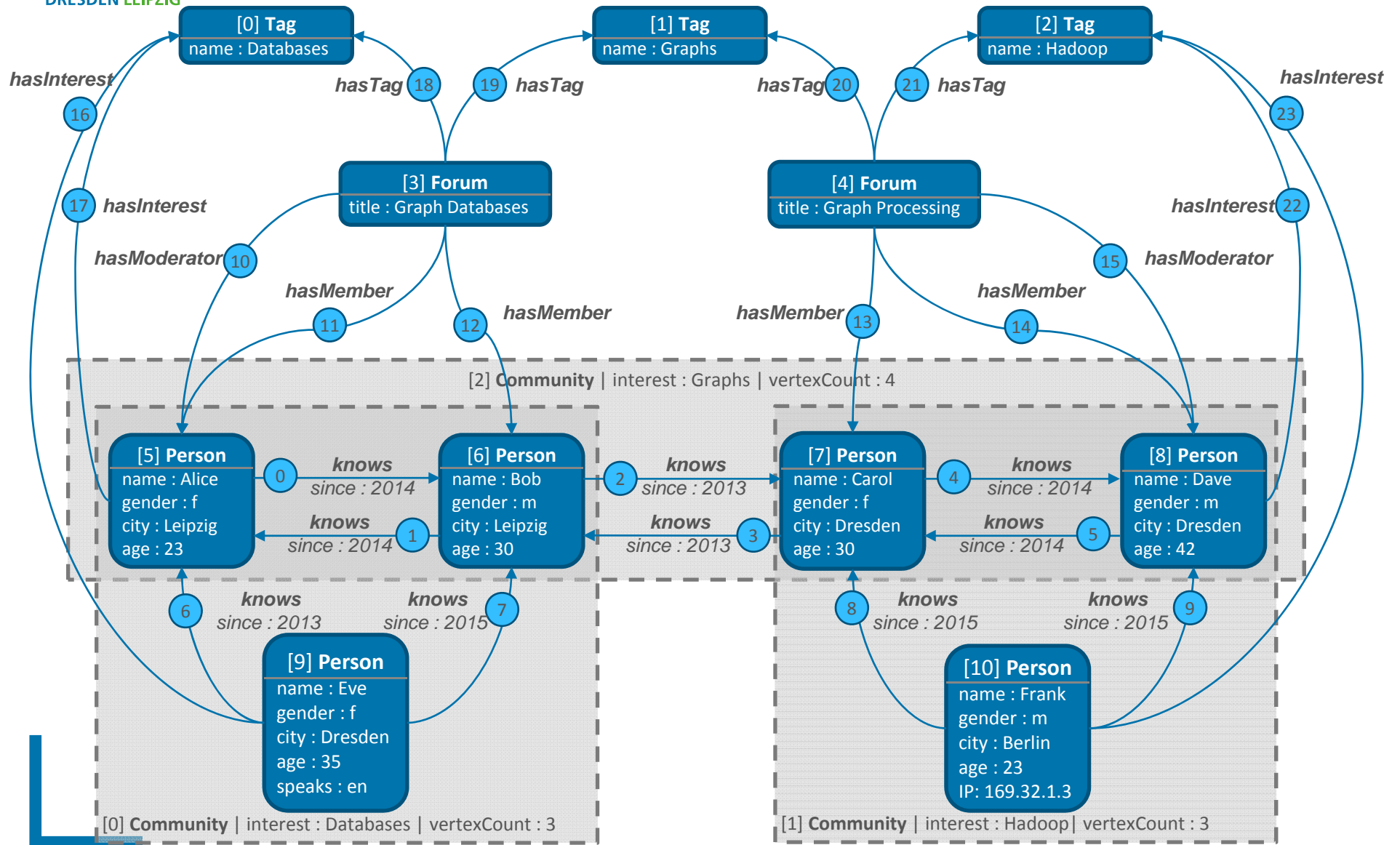
EPGM – GRAPH REPRESENTATION



EPGM – GRAPH REPRESENTATION

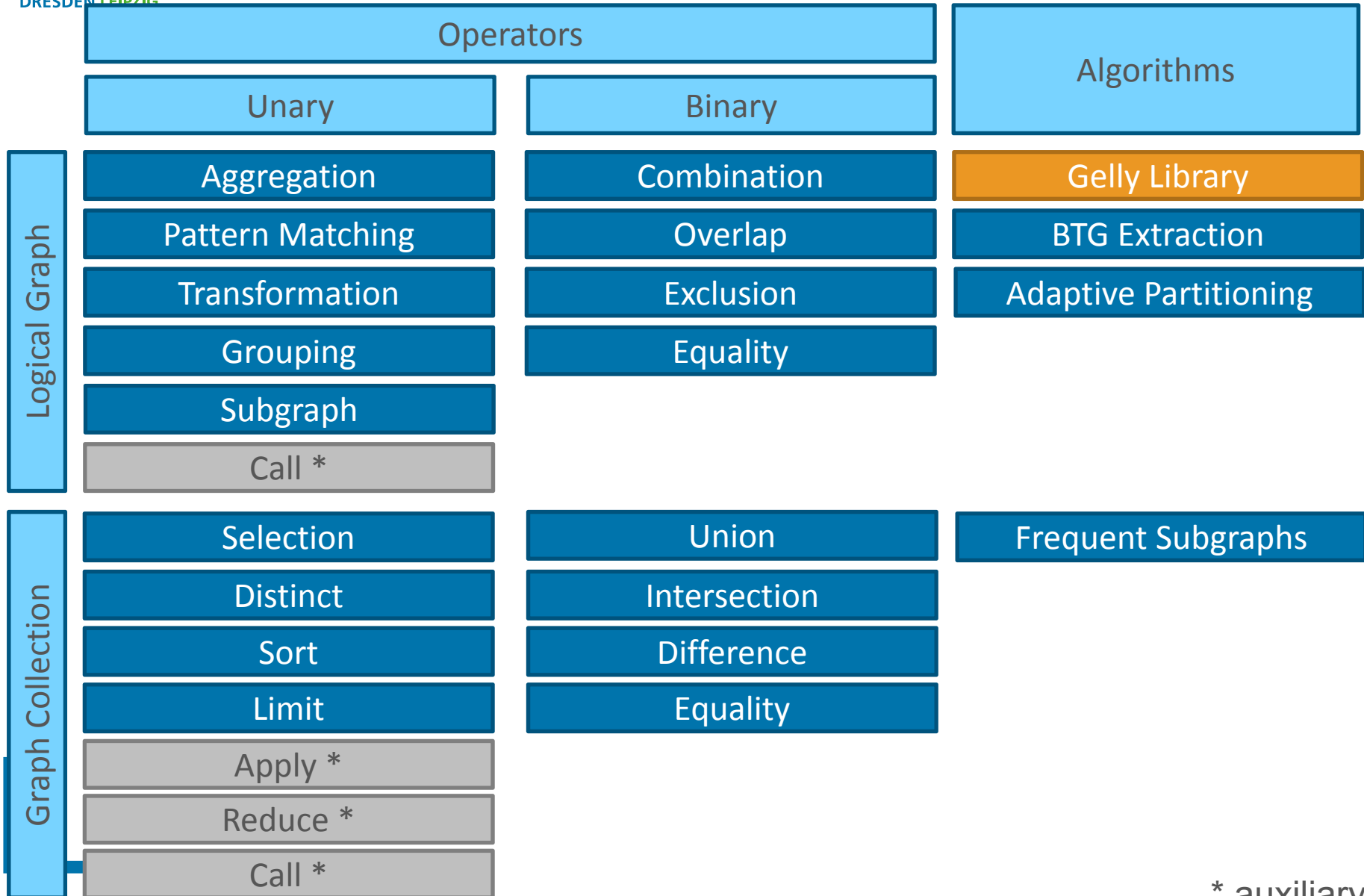


EPGM – GRAPH REPRESENTATION



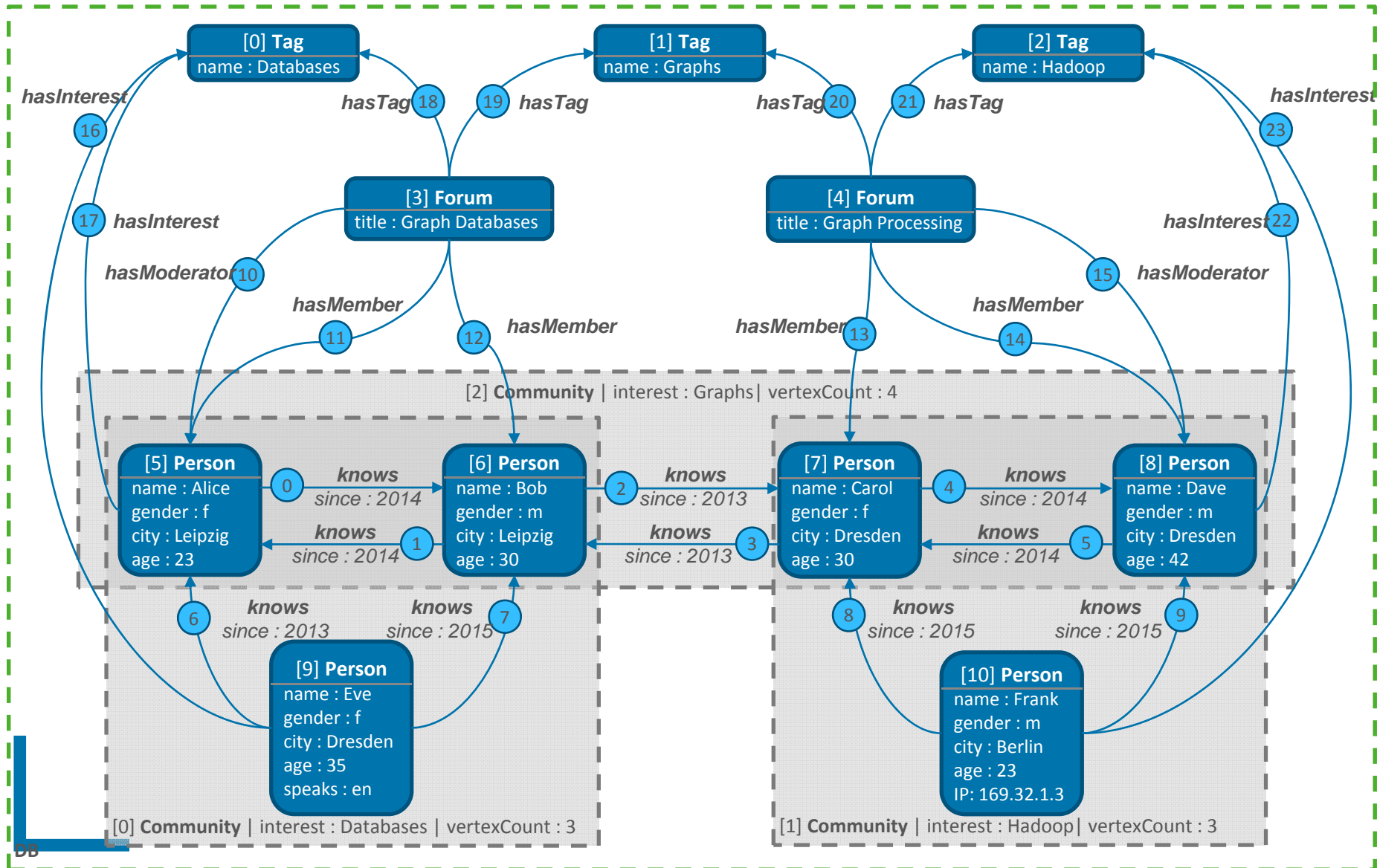


EPGM – OPERATORS AND ALGORITHMS

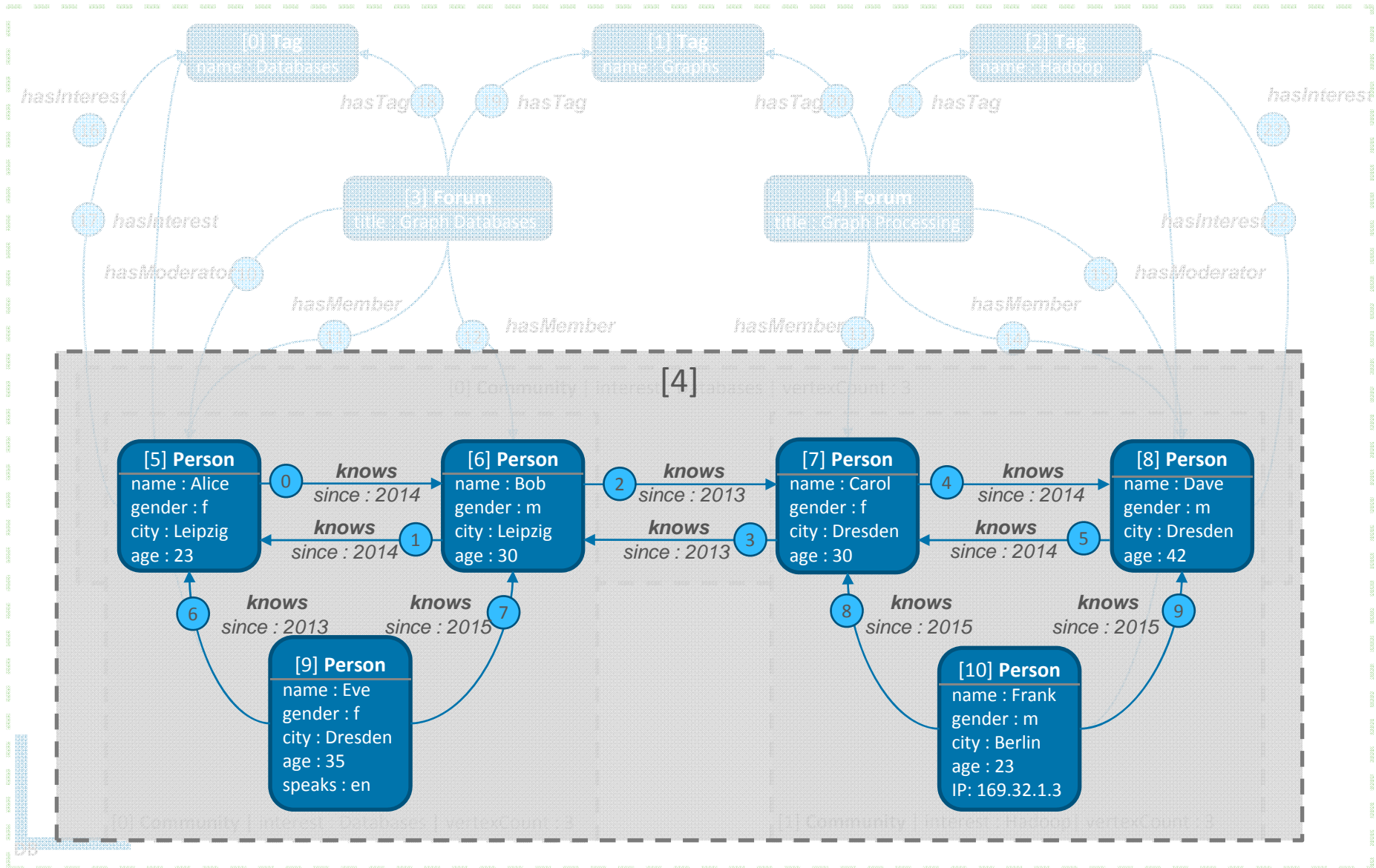


* auxiliary

1: personGraph =
db.G[0].combine(db.G[1]).combine(db.G[2])



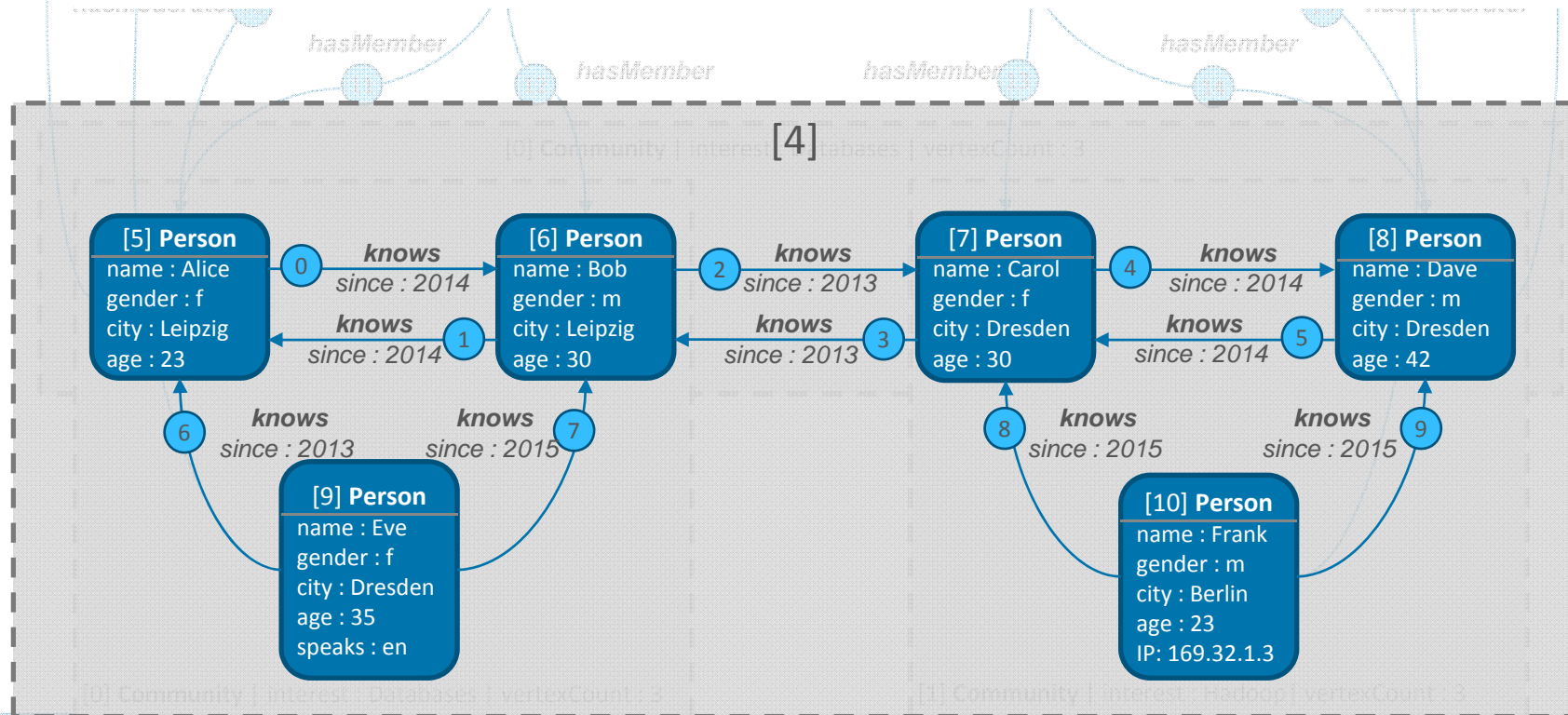
1: personGraph =
db.G[0].combine(db.G[1]).combine(db.G[2])



GROUPING (GRAPH SUMMARIZATION)

```

1: personGraph = db.G[0].combine(db.G[1]).combine(db.G[2])
2: vertexGroupingKeys = [:label, "city"]
3: edgeGroupingKeys = [:label]
4: vertexAggFunc = (superVertex, vertices => superVertex["count"] = |vertices|)
5: edgeAggFunc = (superEdge, edges => superEdge["count"] = |edges|)
6: sumGraph = personGraph.groupBy(vertexGroupingKeys, vertexAggFunc,
    edgeGroupingKeys, edgeAggFunc)
    
```



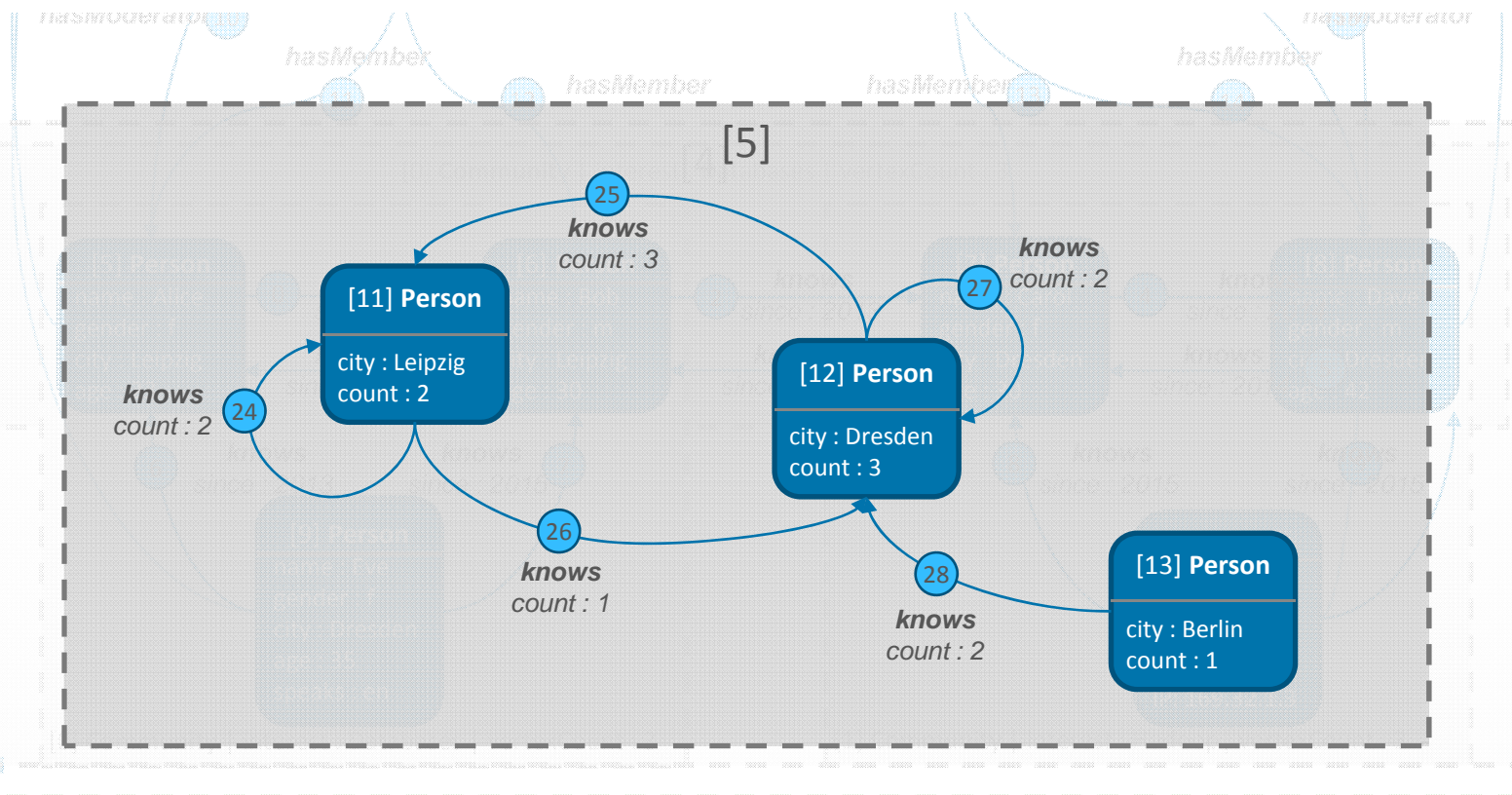
ScaDS GROUPING (2)

DRESDEN LEIPZIG

```

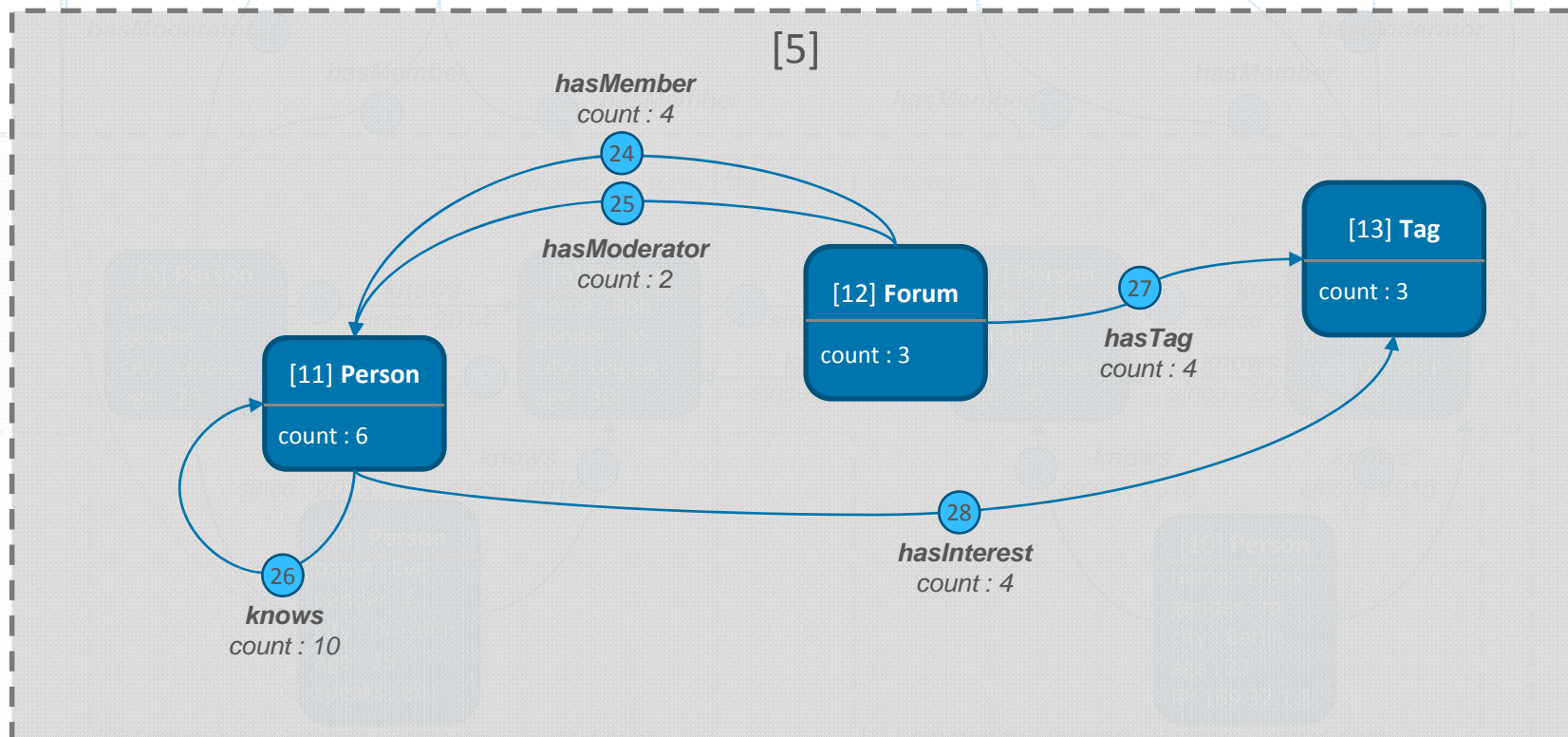
1: personGraph = db.G[0].combine(db.G[1]).combine(db.G[2])
2: vertexGroupingKeys = [:label, "city"]
3: edgeGroupingKeys = [:label]
4: vertexAggFunc = (superVertex, vertices => superVertex["count"] = |vertices|)
5: edgeAggFunc = (superEdge, edges => superEdge["count"] = |edges|)
6: sumGraph = personGraph.groupBy(vertexGroupingKeys, vertexAggFunc,
edgeGroupingKeys, edgeAggFunc)

```

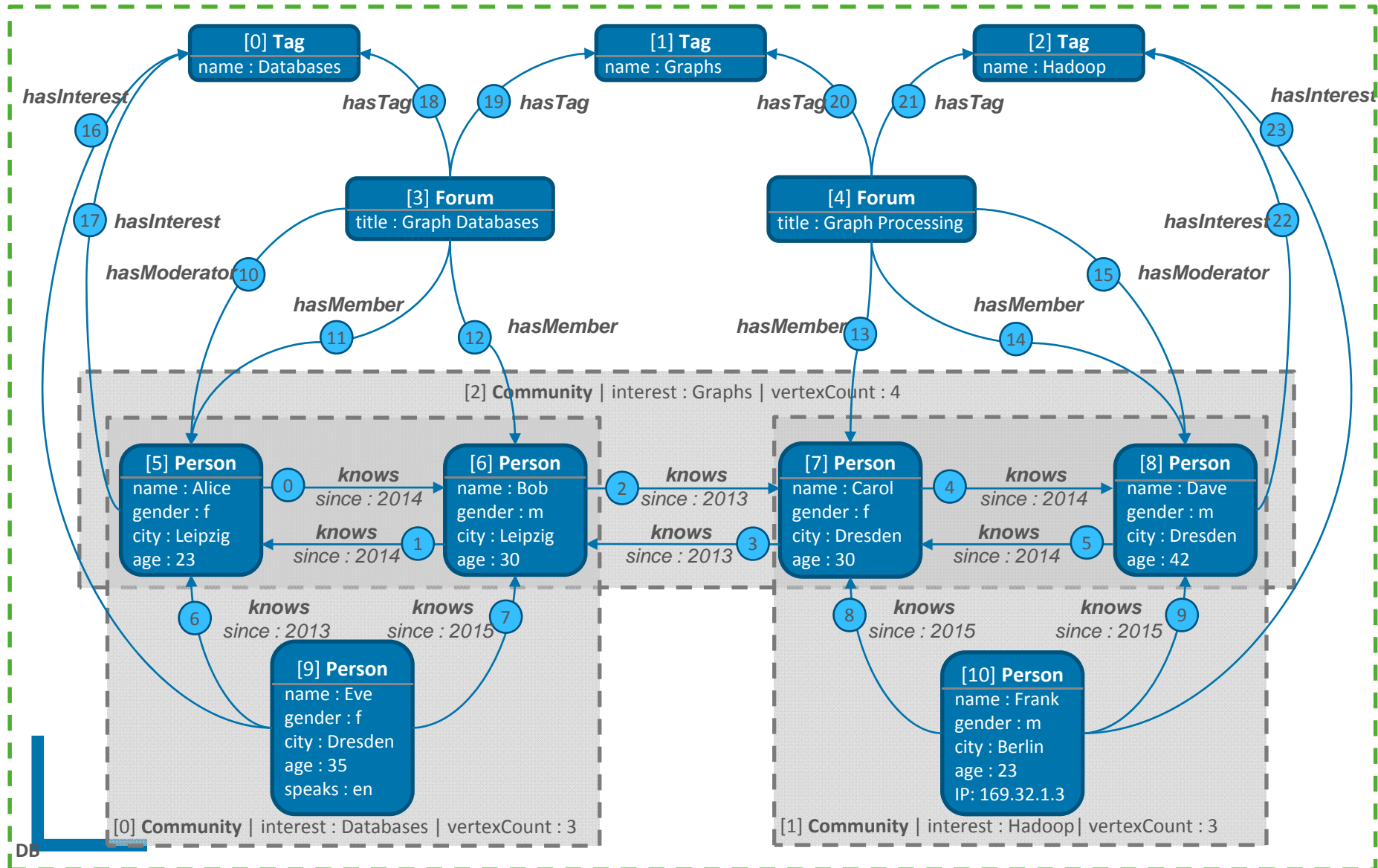


GROUPING: TYPE LEVEL (*SCHEMA GRAPH*)

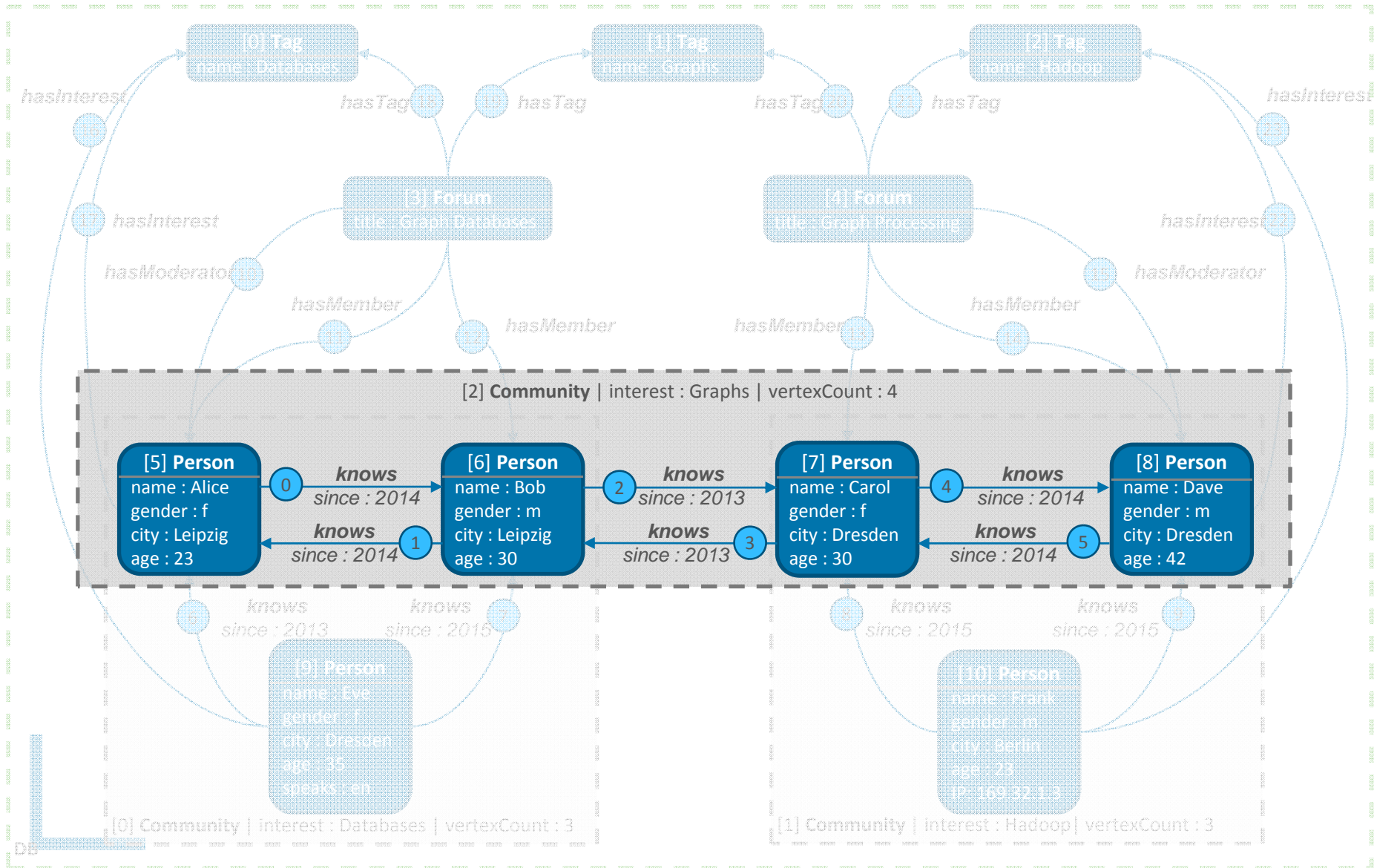
- 1: vertexGroupingKeys = [:label]
- 2: edgeGroupingKeys = [:label]
- 3: vertexAggFunc = (superVertex, vertices => superVertex["count"] = |vertices|)
- 4: edgeAggFunc = (superEdge, edges => superEdge["count"] = |edges|)
- 5: sumGraph = databaseGraph.**groupBy**(vertexGroupingKeys, vertexAggFunc, edgeGroupingKeys, edgeAggFunc)



```
1: resultColl =
db.G[0,1,2].select((g => g["vertexCount"] > 3))
```



```
1: resultColl =
db.G[0,1,2].select((g => g["vertexCount"] > 3))
```



- 1. Large-scale graphs**
 - Support for real-world graphs with millions of vertices and billions of edges
- 2. Graph partitioning**
 - Efficient data distribution to balance load and minimize communication during computation
- 3. Data versioning**
 - Enable time-based graph analytics on properties and graph structure
- 4. Fault tolerance**
 - Prevent data loss in case of cluster failures



- Open Source implementation of Google BigTable
- **Distributed**, persistent, **sparse**, **multidimensional** sorted map based on HDFS
- Data distribution based on row key (i.e., horizontal **partitioning**)
- **Flexible** storage layout (handles only byte[], no types, no schema)
- **Fault tolerancy** through data replication (HDFS)
- **Data versioning** on cell level

HTable

Sorted	Column family 1		Column family 2	
	Column identifier		C. identifier	C. identifier
	versioned value		v. value	v. value
Sorted	Column family 1		Column family 2	
	C. Identifier	C. identifier	Column identifier	
	v. value	v. value	versioned value	

Cell: <rowkey>.<column_family>.<column_identifier>[.<version>]

ScaDS  VERTEX TABLE
DRESDEN LEIPZIG

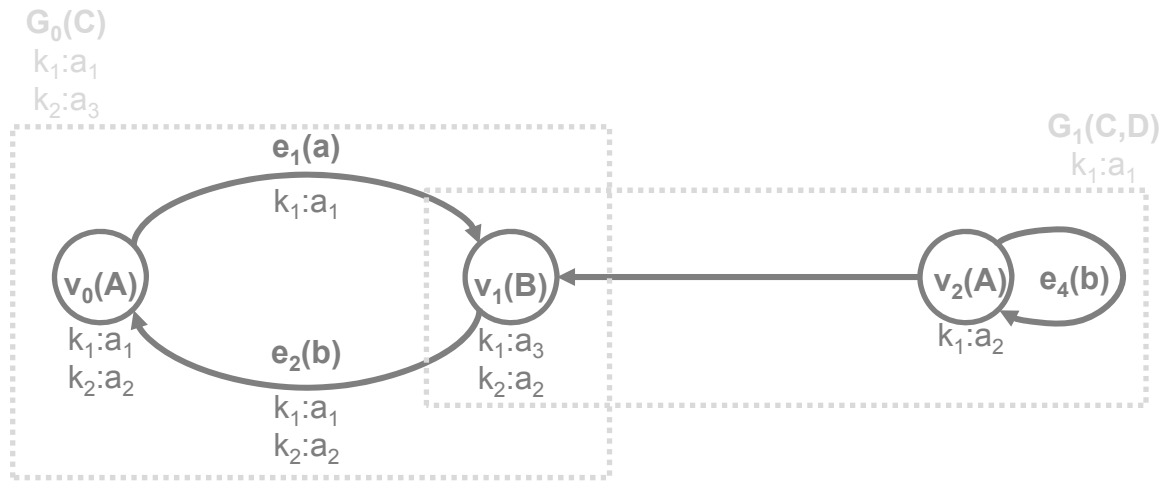


Table 'vertices'

0-0	meta			properties		out edges		in edges	
	type	idx	graphs	k_1	k_2	$\langle a, 0 - 1, 0 \rangle$		$\langle b, 0 - 1, 0 \rangle$	
	A	1	[0]	$\langle t_1, a_1 \rangle$	$\langle t_2, a_2 \rangle$	[[$(k_1, \langle t_1, a_1 \rangle)$]]		[[$(k_1, \langle t_1, a_1 \rangle)$], [$(k_2, \langle t_2, a_2 \rangle)$]]	
0-1	meta			properties		out edges		in edges	
	type	idx	graphs	k_1	k_2	$\langle b, 0 - 0, 0 \rangle$		$\langle a, 0 - 0, 0 \rangle$	$\langle a, 0 - 2, 0 \rangle$
	B	1	[0,1]	$\langle t_2, a_3 \rangle$	$\langle t_2, a_2 \rangle$	[[$(k_1, \langle t_1, a_1 \rangle)$], [$(k_2, \langle t_2, a_2 \rangle)$]]		[[$(k_1, \langle t_1, a_1 \rangle)$]]	[]
0-2	meta			properties		out edges		in edges	
	type	idx	graphs	k_1		$\langle a, 0 - 1, 0 \rangle$	$\langle b, 0 - 2, 1 \rangle$	$\langle b, 0 - 2, 1 \rangle$	
	A	2	[1]	$\langle t_2, a_2 \rangle$		[]		[]	[]

ScaDS  VERTEX TABLE
DRESDEN LEIPZIG

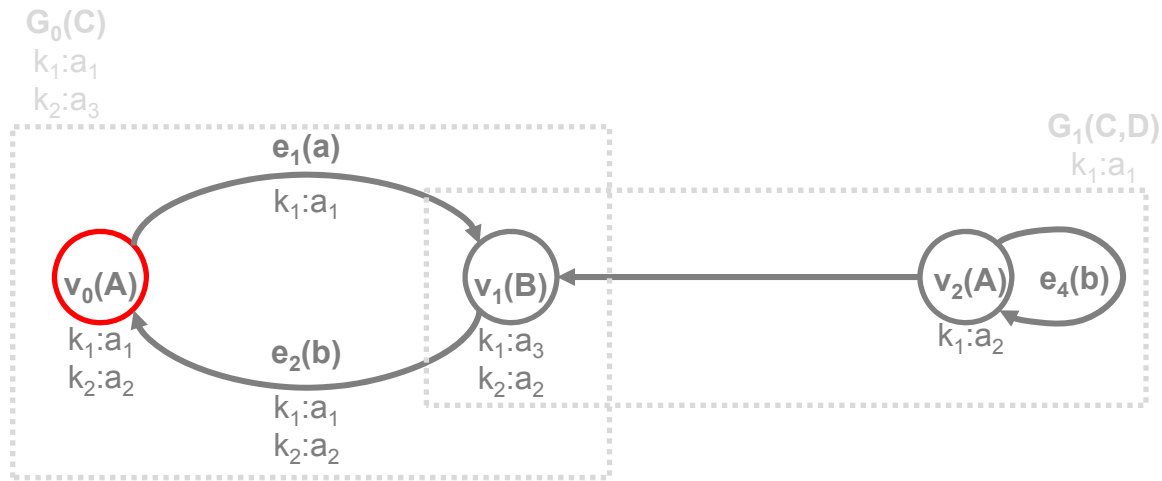


Table 'vertices'

0-0	meta			properties		out edges		in edges	
	type	idx	graphs	k_1	k_2	$\langle a, 0 - 1, 0 \rangle$		$\langle b, 0 - 1, 0 \rangle$	
	A	1	[0]	$\langle t_1, a_1 \rangle$	$\langle t_2, a_2 \rangle$	[[$(k_1, \langle t_1, a_1 \rangle)$]]		[[$(k_1, \langle t_1, a_1 \rangle)$], [$(k_2, \langle t_2, a_2 \rangle)$]]	
0-1	meta			properties		out edges		in edges	
	type	idx	graphs	k_1	k_2	$\langle b, 0 - 0, 0 \rangle$		$\langle a, 0 - 0, 0 \rangle$	$\langle a, 0 - 2, 0 \rangle$
	B	1	[0,1]	$\langle t_2, a_3 \rangle$	$\langle t_2, a_2 \rangle$	[[$(k_1, \langle t_1, a_1 \rangle)$], [$(k_2, \langle t_2, a_2 \rangle)$]]		[[$(k_1, \langle t_1, a_1 \rangle)$]]	[]
0-2	meta			properties		out edges		in edges	
	type	idx	graphs	k_1		$\langle a, 0 - 1, 0 \rangle$	$\langle b, 0 - 2, 1 \rangle$	$\langle b, 0 - 2, 1 \rangle$	
	A	2	[1]	$\langle t_2, a_2 \rangle$		[]		[]	[]

ScaDS  VERTEX TABLE
DRESDEN LEIPZIG

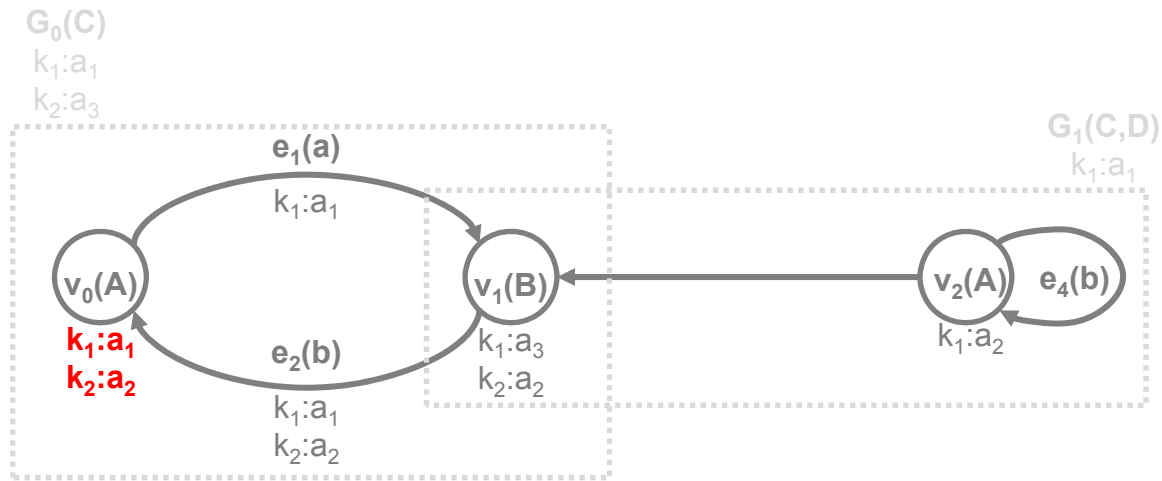


Table 'vertices'

id	meta			properties		out edges		in edges	
	type	idx	graphs	k_1	k_2				
0-0	A	1	[0]	$\langle t_1, a_1 \rangle$	$\langle t_2, a_2 \rangle$	[[$(k_1, \langle t_1, a_1 \rangle)$]]		[[$(k_1, \langle t_1, a_1 \rangle)$], [$(k_2, \langle t_2, a_2 \rangle)$]]	
0-1	B	1	[0,1]	$\langle t_2, a_3 \rangle$	$\langle t_2, a_2 \rangle$	[[$(k_1, \langle t_1, a_1 \rangle)$], [$(k_2, \langle t_2, a_2 \rangle)$]]		[[$(k_1, \langle t_1, a_1 \rangle)$]]	[]
0-2	A	2	[1]	$\langle t_2, a_2 \rangle$		[[$(a, 0 - 1, 0)$]]		[[$(b, 0 - 2, 1)$]]	[]
						[]		[]	[]

ScaDS  VERTEX TABLE
DRESDEN LEIPZIG

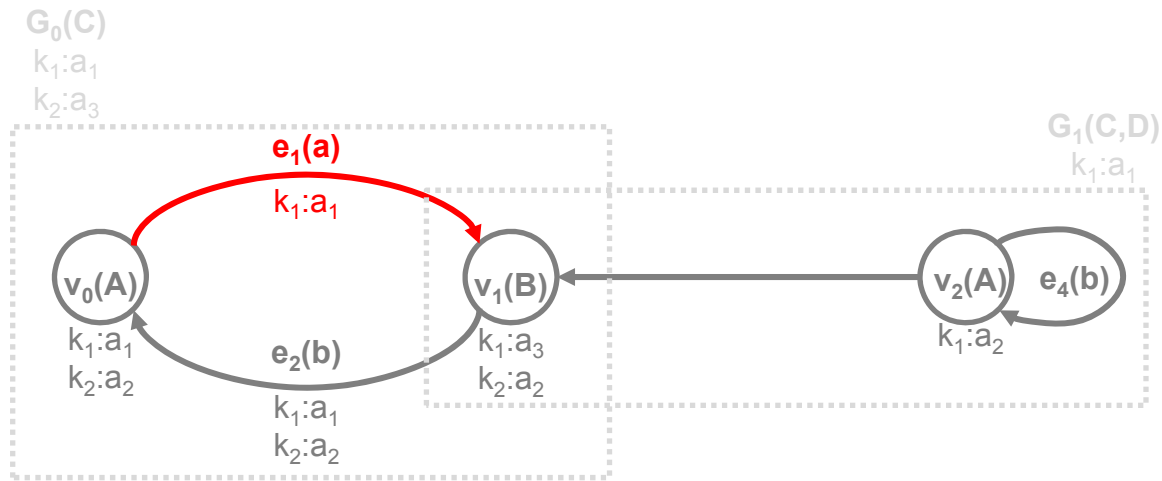


Table 'vertices'

id	meta			properties		out edges		in edges	
	type	idx	graphs	k_1	k_2				
0-0	A	1	[0]	$\langle t_1, a_1 \rangle$	$\langle t_2, a_2 \rangle$	$\langle a, 0 - 1, 0 \rangle$		$\langle b, 0 - 1, 0 \rangle$	
						$[(k_1, \langle t_1, a_1 \rangle)]$		$[(k_1, \langle t_1, a_1 \rangle), (k_2, \langle t_2, a_2 \rangle)]$	
0-1	B	1	[0,1]	$\langle t_2, a_3 \rangle$	$\langle t_2, a_2 \rangle$	$\langle b, 0 - 0, 0 \rangle$		$\langle a, 0 - 0, 0 \rangle$	$\langle a, 0 - 2, 0 \rangle$
						$[(k_1, \langle t_1, a_1 \rangle), (k_2, \langle t_2, a_2 \rangle)]$		$[(k_1, \langle t_1, a_1 \rangle)]$	$[\]$
0-2	A	2	[1]	$\langle t_2, a_2 \rangle$		$\langle a, 0 - 1, 0 \rangle$	$\langle b, 0 - 2, 1 \rangle$	$\langle b, 0 - 2, 1 \rangle$	
						$[\]$		$[\]$	$[\]$

ScaDS  VERTEX TABLE
DRESDEN LEIPZIG

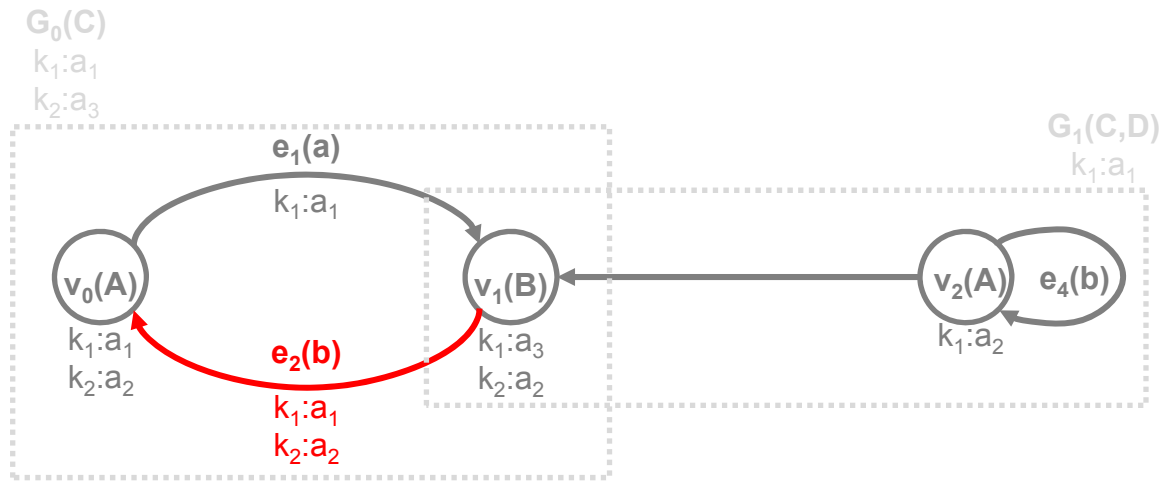


Table 'vertices'

id	meta			properties		out edges		in edges	
	type	idx	graphs	k_1	k_2				
0-0	A	1	[0]	$\langle t_1, a_1 \rangle$	$\langle t_2, a_2 \rangle$	[[$\langle k_1, \langle t_1, a_1 \rangle \rangle$]]		[[$\langle k_1, \langle t_1, a_1 \rangle \rangle, \langle k_2, \langle t_2, a_2 \rangle \rangle$]]	
								$\langle b, 0 - 1, 0 \rangle$	
0-1	B	1	[0,1]	$\langle t_2, a_3 \rangle$	$\langle t_2, a_2 \rangle$	[[$\langle k_1, \langle t_1, a_1 \rangle \rangle, \langle k_2, \langle t_2, a_2 \rangle \rangle$]]		[[$\langle k_1, \langle t_1, a_1 \rangle \rangle$]]	
								$\langle b, 0 - 0, 0 \rangle$	
0-2	A	2	[1]	$\langle t_2, a_2 \rangle$		[[$\langle a, 0 - 1, 0 \rangle, \langle b, 0 - 2, 1 \rangle$]]		$\langle b, 0 - 2, 1 \rangle$	
						[]		[]	

1. Business Intelligence

- Top Revenue Subgraph
- Find the common subgraph of the top 100 revenue business transaction graphs

2. Social Network Analysis

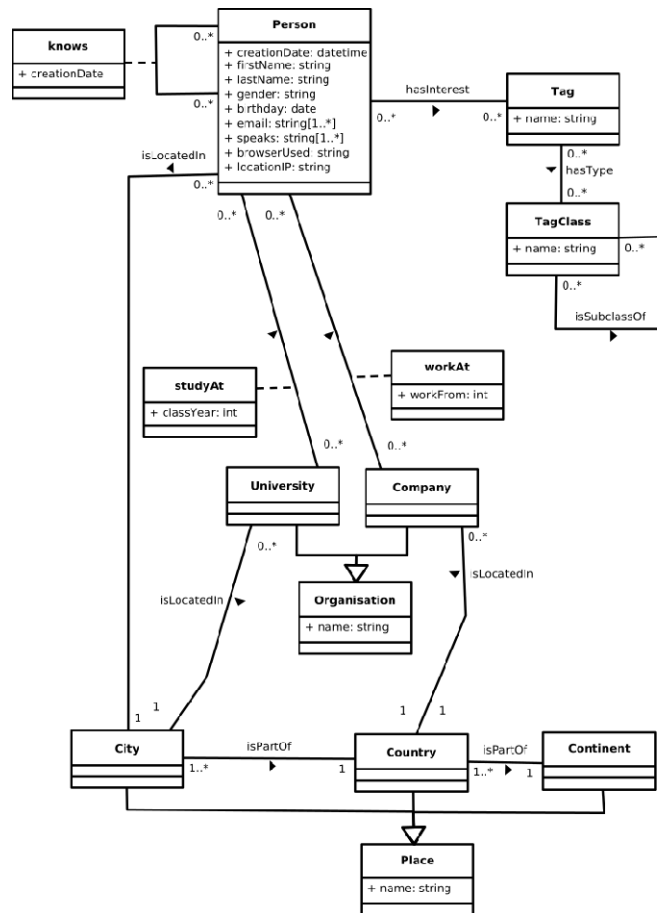
- “Summarized Communities”
- Find communities by label propagation
- Filter communities by number of users
- Summarize vertices per community and edges between community members



GRALA EXAMPLE : TOP REVENUE SUBGRAPH

```
// compute logical graphs
1: btgs = db.callForCollection( :BusinessTransactionGraphs , {} )
// define and apply aggregate function (number of invoices per graph)
2: aggFuncInvoiceCount = ( Graph g =>
    |g.V.filter( Vertex v => v[:type] == "Invoice")|)
3: btgs = btgs.apply(
    Graph g => g.aggregate( "invoiceCount",aggFuncInvoiceCount) )
// select logical graphs with at least one invoice
4: invBtgs = btgs.select(
    Graph g => g["invoiceCount"] > 0)
// define and apply aggregate function (revenue per graph)
5: aggFuncRevenue = ( Graph g =>
    g.V.values("revenue").sum())
6: invBtgs = invBtgs.apply(
    Graph g => g.aggregate( "revenue",aggFuncRevenue) )
// sort graphs by revenue and return top 100
7: topBtgs = invBtgs.sortBy( "revenue" , :desc ).top( 100 )
// compute overlap to find master data objects (e.g., Employees)
8: topBtgOverlap = invBtgs.reduce(
    Graph g, Graph h => g.overlap(h))
```





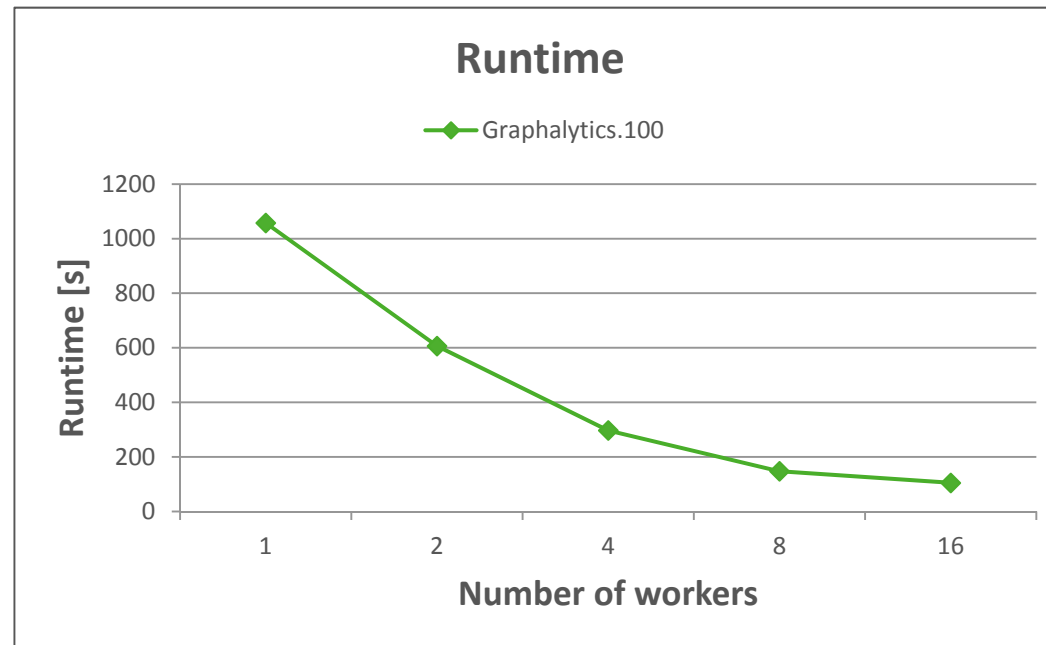
socialNetwork

```

.subgraph(
  (v => v.label == 'Person'),
  (e => e.label == 'knows'))
.transform(
  (gIn, gOut => gOut = gIn),
  (vIn, vOut => {
    vOut.label      = vIn.label,
    vOut['city']    = vIn['city']
    vOut['gender']  = vIn['gender']
    vOut['key']     = vIn['birthday']}),
  (eIn, eOut) => eOut.label = eIn.label))
.callForCollection(:LabelPropagation, ['key', 4])
.apply(g =>
  g.aggregate('vertexCount', (h => |h.V|))
  .select(g => g['vertexCount'] > 50000)
  .reduce(g, h => g.combine(h))
  .groupBy(
    ['city', 'gender'], (superVertex, vertices =>
      superVertex['count'] = |vertices|),
    [], (superEdge, edges =>
      superEdge['count'] = |edges|)
  )
  .aggregate('vCount', (g => |g.V|))
  .aggregate('eCount', (g => |g.E|))

```

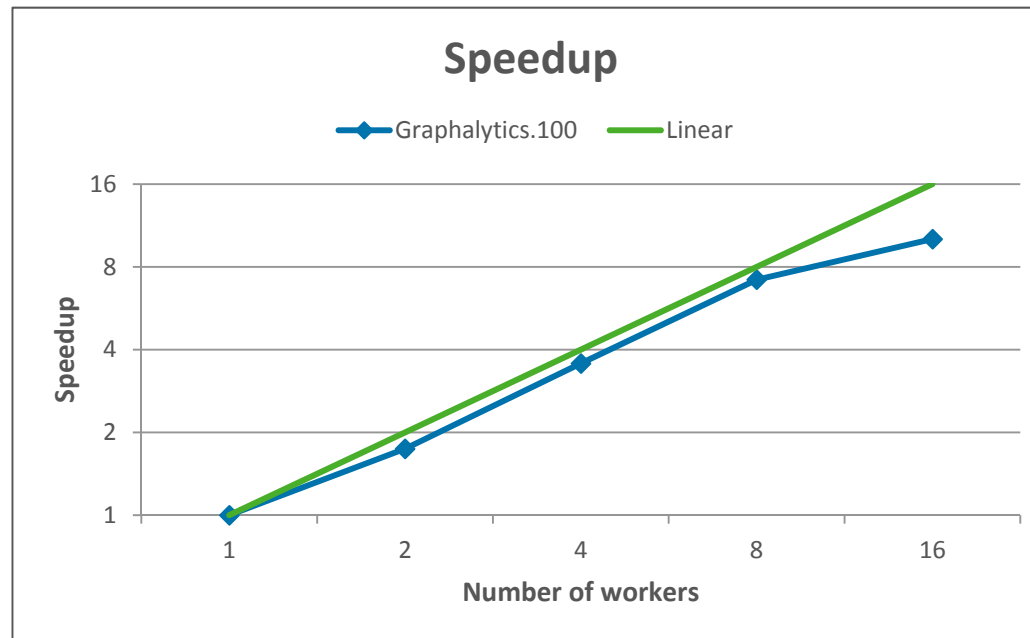
BENCHMARK RESULTS



Dataset	# Vertices	# Edges
Graphalytics.1	61,613	2,026,082
Graphalytics.10	260,613	16,600,778
Graphalytics.100	1,695,613	147,437,275
Graphalytics.1000	12,775,613	1,363,747,260
Graphalytics.10000	90,025,613	10,872,109,028

- 16x Intel(R) Xeon(R) 2.50GHz (6 Cores)
- 16x 48 GB RAM
- 1 Gigabit Ethernet
- Hadoop 2.6.0
- Flink 1.0-SNAPSHOT

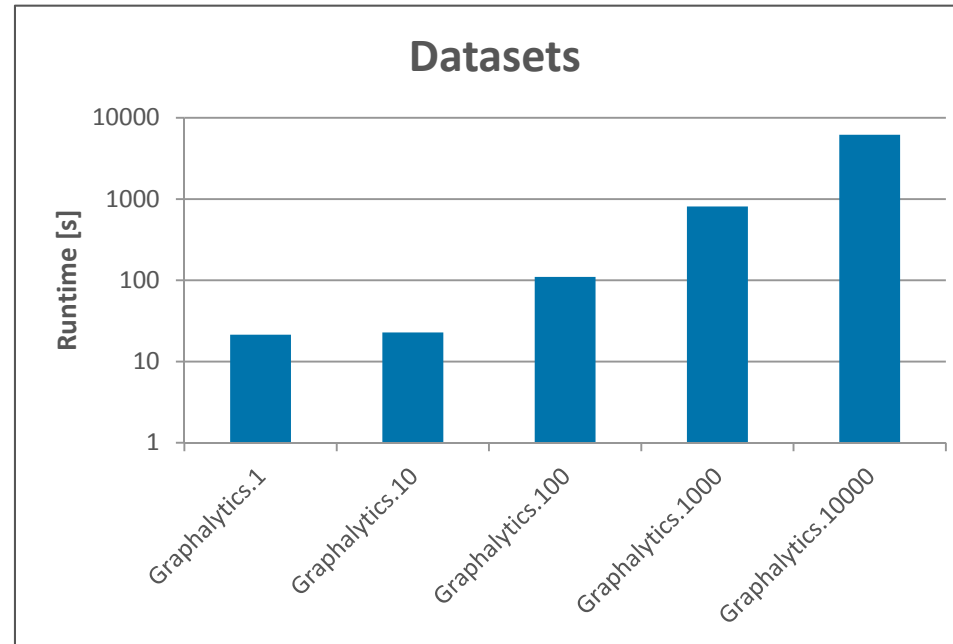
BENCHMARK RESULTS 2



Dataset	# Vertices	# Edges
Graphalytics.1	61,613	2,026,082
Graphalytics.10	260,613	16,600,778
Graphalytics.100	1,695,613	147,437,275
Graphalytics.1000	12,775,613	1,363,747,260
Graphalytics.10000	90,025,613	10,872,109,028

- 16x Intel(R) Xeon(R) 2.50GHz (6 Cores)
- 16x 48 GB RAM
- 1 Gigabit Ethernet
- Hadoop 2.6.0
- Flink 1.0-SNAPSHOT

BENCHMARK RESULTS 3



Dataset	# Vertices	# Edges
Graphalytics.1	61,613	2,026,082
Graphalytics.10	260,613	16,600,778
Graphalytics.100	1,695,613	147,437,275
Graphalytics.1000	12,775,613	1,363,747,260
Graphalytics.10000	90,025,613	10,872,109,028

- 16x Intel(R) Xeon(R) 2.50GHz (6 Cores)
- 16x 48 GB RAM
- 1 Gigabit Ethernet
- Hadoop 2.6.0
- Flink 1.0-SNAPSHOT

- **Graph-based data integration**
 - centralized „linked data“ using PGM rather than RDF
 - data/metadata extraction and transformation into graphs
 - linking / matching + fusion
- **Big Graph Analytics**
 - high potential even for business intelligence (BIIG)
 - Hadoop-based graph processing frameworks based on generic graphs
 - Spark/Flink: batch-oriented workflows (rather than OLAP)
 - Graph collections not generally supported



- **GraDoop**
 - infrastructure for entire processing pipeline: graph acquisition, storage, integration, transformation, analysis (queries + graph mining), visualization
 - extended property graph model (EPGM) with powerful operators (e.g. grouping) and support for graph collections
 - leverages Hadoop ecosystem
 - **Apache HBase for permanent graph storage**
 - **Apache Flink to implement operators**
 - ongoing implementation



OUTLOOK / CHALLENGES

- **Graph-based data integration**
 - unified approach for knowledge graphs and regular data graphs
 - evaluate/improve scalability and data quality

- **Graph analytics**
 - automatic optimization of analysis workflows
 - optimized graph partitioning approaches
 - load balancing
 - interactive graph analytics
 - visualization of graphs and analysis results



- R. Angles, C. Guitierrez. *Survey of Graph Database Models*. ACM Comput. Surv., 2008
- M. Junghanns, A. Petermann, K. Gomez, E. Rahm: *GRADOOP - Scalable Graph Data Management and Analytics with Hadoop*. Tech. report, Univ. of Leipzig, June 2015
- M. Junghanns, A. Petermann, N. Teichmann, K. Gomez, E. Rahm: *Analyzing Extended Property Graphs with Apache Flink*. Tech. report, Univ. of Leipzig, Feb. 2016
- A. Petermann, M. Junghanns, R. Müller, E. Rahm: *BIIG : Enabling Business Intelligence with Integrated Instance Graphs*. Proc. 5th Int. Workshop on Graph Data Management (GDM 2014)
- A. Petermann, M. Junghanns, R. Müller, E. Rahm: *Graph-based Data Integration and Business Intelligence with BIIG*. Proc. VLDB Conf., 2014
- Petermann, A.; Junghanns, M.; Müller, R.; Rahm, E.: *FoodBroker - Generating Synthetic Datasets for Graph-Based Business Analytics*. Proc. 5th Int. Workshop on Big Data Benchmarking (WBDB), 2014
- R.S. Xin, J.E. Gonzales, M.J. Franklin, I. Stoica: *GraphX: A resilient distributed graph system on Spark*. GRADES, 2013
- Z. Wang et. al. *Pagrol: Parallel Graph OLAP over Large-scale Attributed Graphs*. ICDE, 2014
- P. Zhao et. al. *Graph Cube: On Warehousing and OLAP Multidimensional Networks*. Sigmod, 2011

