



Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

FAKULTÄT
FÜR MATHEMATIK, INFORMATIK
UND NATURWISSENSCHAFTEN

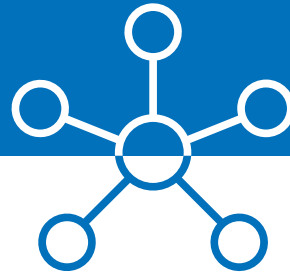
Databases and Information Systems (DIS) – Quiz 2

Universität Hamburg



Foto: UHH/Esfandiari

Modern DBS: NoSQL Systems



Which graph models include edge and vertex properties?

Which model is used in neo4j?

RDF Graph

Weighted Graph

Property Graph

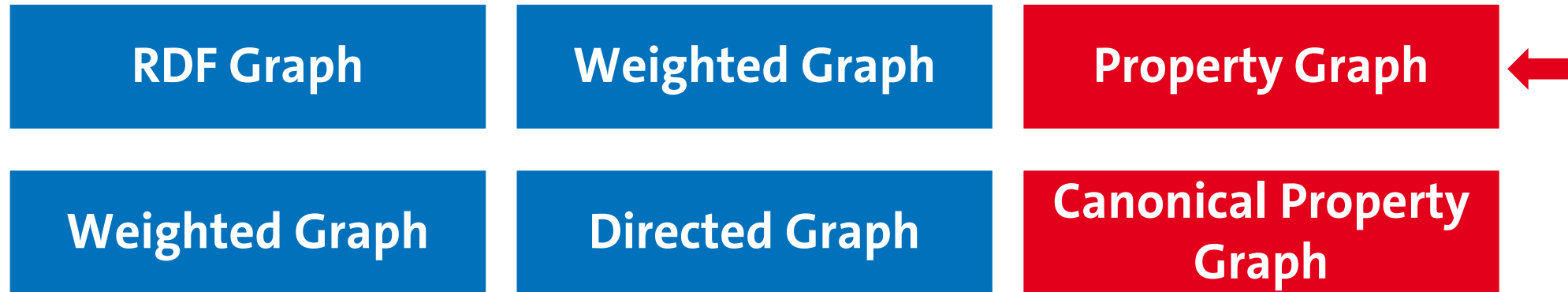
Weighted Graph

Directed Graph

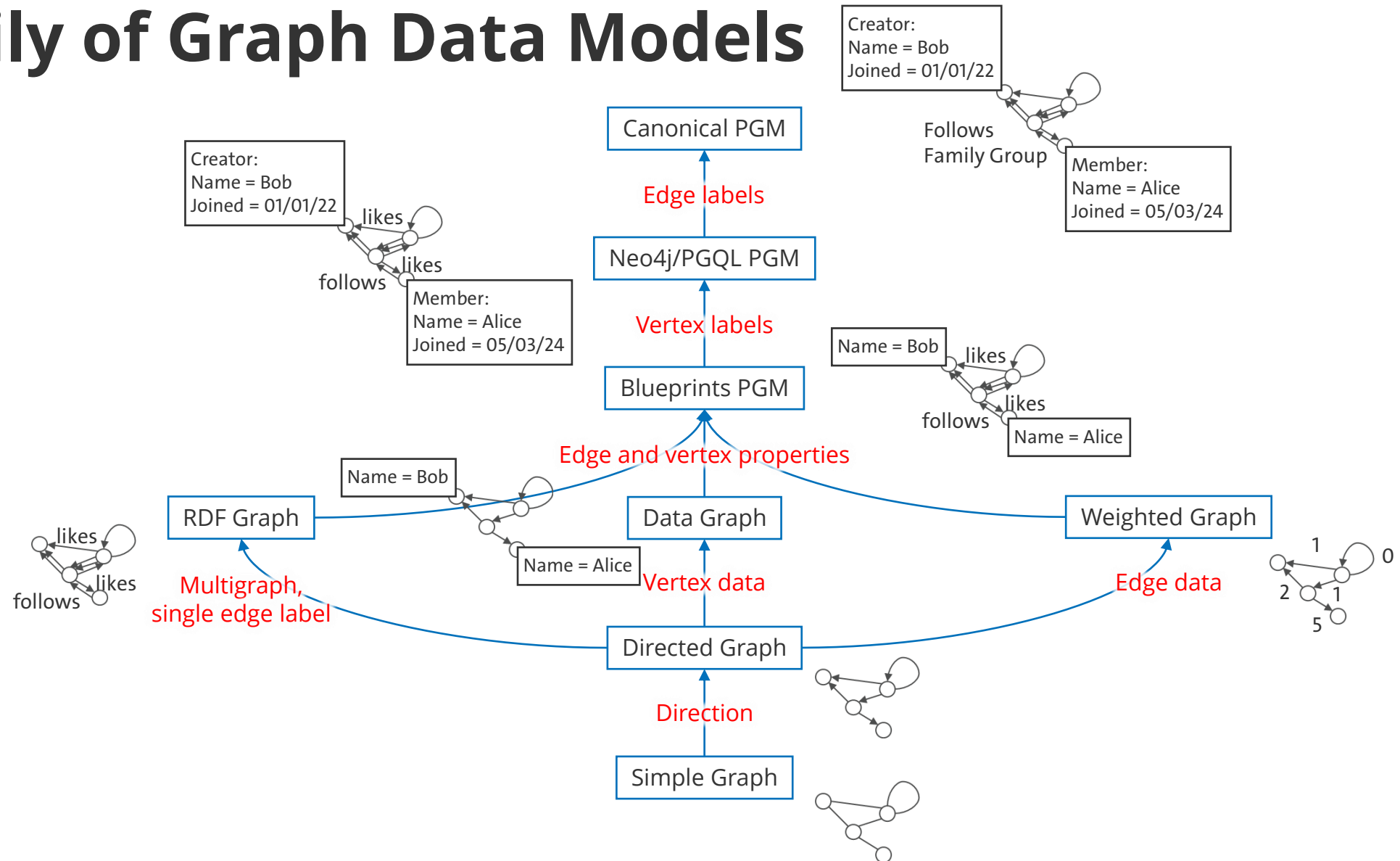
Canonical Property
Graph

Which graph models include edge and vertex properties?

Which model is used in neo4j?



Family of Graph Data Models



Which of the following elements can directly be modelled using neo4j?

Entities with
properties

Edge properties

Multivalued
properties

N-ary relationships

N:M relationships

Which of the following elements can directly be modelled using neo4j?

**Entities with
properties**

Edge properties

**Multivalued
properties**

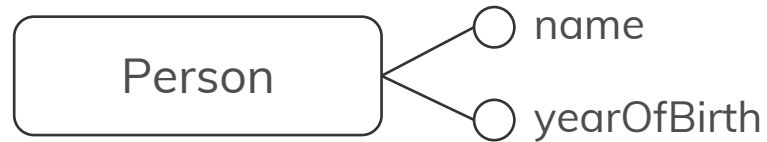
N-ary relationships

N:M relationships

Property Graph Modelling

Entity with Properties

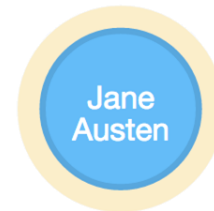
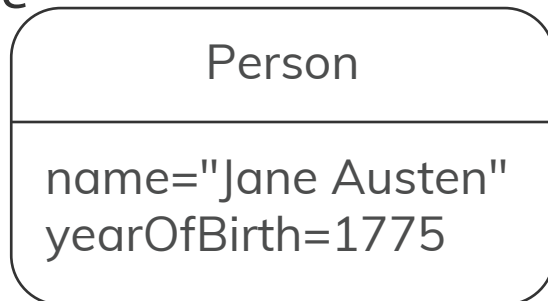
- Entity relationship model



- Schema typically implicit, i.e. given with instances

```
(ja:Person { name: 'Jane Austen', yearOfBirth: 1775 })
```

- Instance



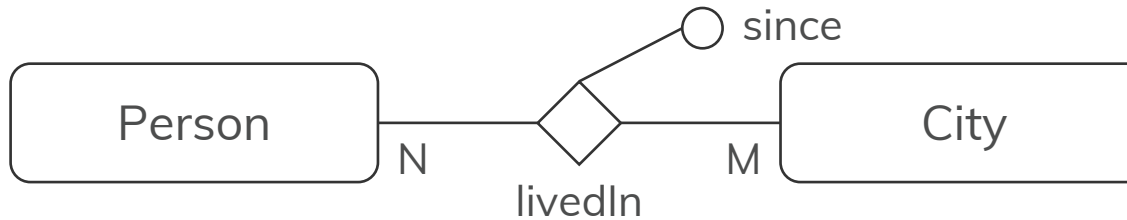
Person

<id>: 175 yearOfBirth: 1775 name: Jane Austen

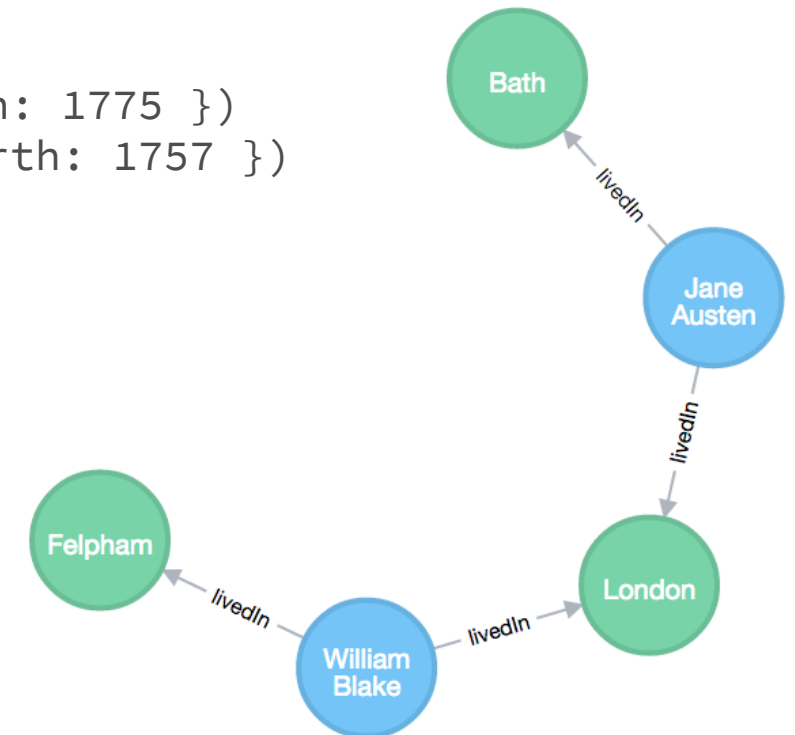
Property Graph Modelling

Relationships (N:M)

- Entity relationship model



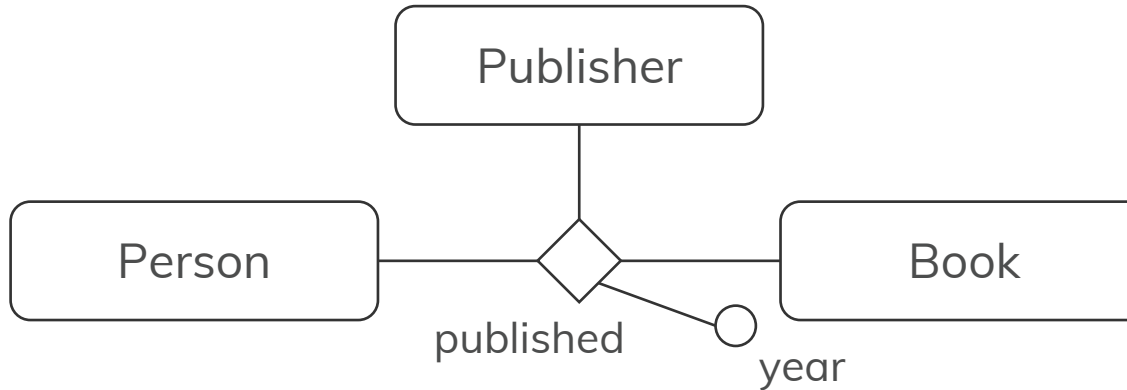
- Vertices
 - (**ja**:**Person** { name: 'Jane Austen', yearOfBirth: 1775 })
 - (**wb**:**Person** { name: 'William Blake', yearOfBirth: 1757 })
 - (**lo**:**City** {name: 'London'})
 - (**ba**:**City** {name: 'Bath'})
 - (**fe**:**City** {name: 'Felpham'})
- Edges
 - (**ja**)-[:**livedIn** {since: 1775}]->(lo)
 - (**ja**)-[:**livedIn** {since: 1800}]->(ba)
 - (**wb**)-[:**livedIn** {since: 1757}]->(lo)
 - (**wb**)-[:**livedIn** {since: 1800}]->(fe)



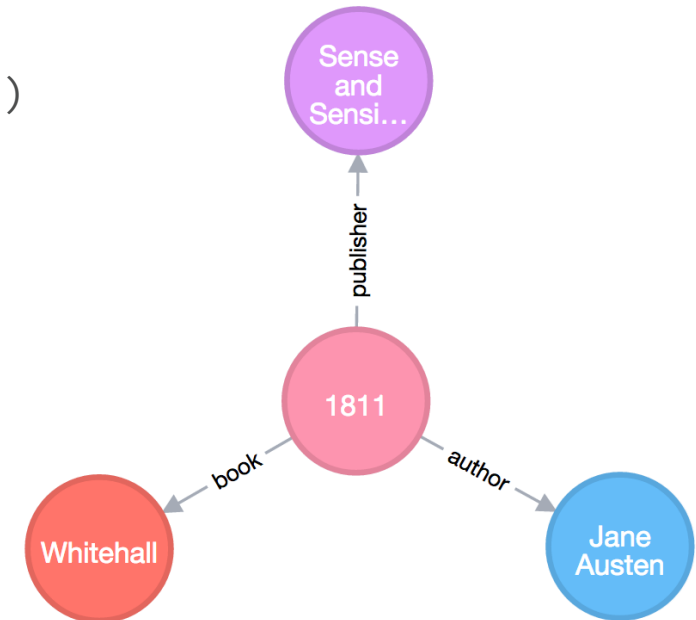
Property Graph Modelling

n-ary relationships (with $n > 2$)

- Entity relationship model

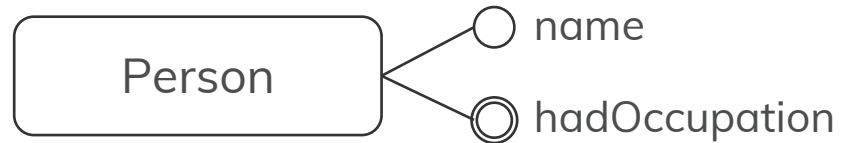


```
(ja:Person { name: 'Jane Austen', yearOfBirth: 1775 })  
(wh:Publisher { name: 'Whitehall' })  
(sas:Book {title: 'Sense and Sensibility' })  
(pub:Publication {year: 1811 })  
(pub)-[:author]->(ja)  
(pub)-[:book]->(wh)  
(pub)-[:publisher]->(sas)
```



Property Graph Modelling

Multivalued properties



```
(wb:Person {name: 'William Blake', hadOccupation: ['Poet','Painter','Printmaker']})
```



Person

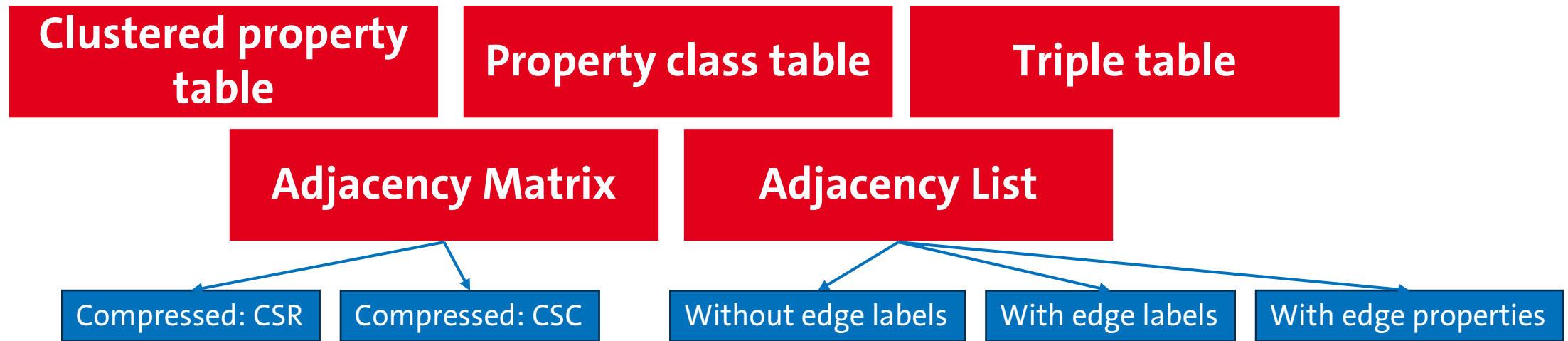
<id>: 183

hadOccupation: Poet,Painter,Printmaker

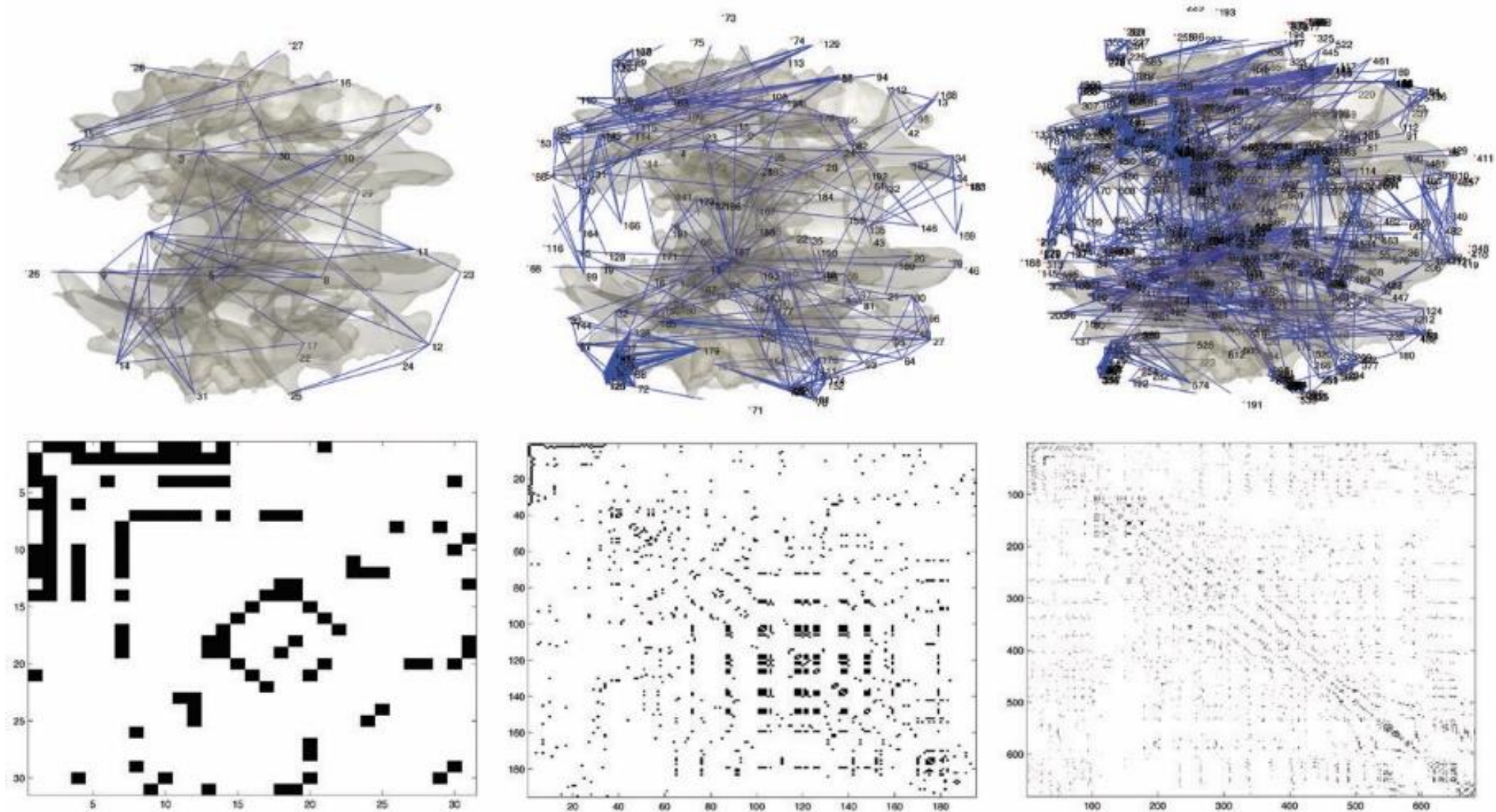
name: William Blake

Name 3 different ways of representing graph data in a database system.

Name 3 different ways of representing graph data in a database system.



Adjacency Matrix



[<http://brainimaging.waisman.wisc.edu/~chung/graph/admatrix.jpg>]

Compress Sparse Row (CSR)

$$\begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix} \begin{pmatrix} \overset{0}{a} & \overset{1}{b} & \overset{2}{c} & \overset{3}{e} \\ d & i & f & \\ h & g & & \end{pmatrix}$$

Position for row # in other two arrays:

#

0

1

2

3

Row position array:

0

2

4

7

#

0

1

2

3

4

5

6

7

8

Column index array:

0

2

1

3

0

1

2

1

2

Cell value array:

a

c

b

e

d

i

f

h

g

Representations: Adjacency List

Source vertex with outgoing edges...

$$\begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix} & \begin{pmatrix} a & & c & \\ & b & & e \\ d & i & f & \\ & h & g & \end{pmatrix} \end{matrix}$$

...without edge labels

0 \rightarrow (0,2)
1 \rightarrow (1,3)
2 \rightarrow (0,1,2)
3 \rightarrow (1,2)

...with edge labels

0 -> ([0,a],[2,c])
 1 -> ([1,b],[3,e])
 2 -> ([0,d],[1,i],[2,f])
 3 -> ([1,h],[2,g])

...with edge properties

0 -> ([0,a,(weight=4)],[2,c,(weight=3)])
 1 -> ([1,b,(weight=3)],[3,e,(weight=2)])
 2 -> ([0,d,(weight=5)],[1,i,(weight=2)],...)
 3 -> ([1,h,(weight=9)],[2,g,(weight=7)])

The same!

Almost the same!

Compressed Sparse Row

Diagram illustrating the insertion of the element 7 into the sorted array [0, 2, 4]. The array is shown as a 2x2 grid. Arrows indicate the shifting of elements to the right to make space for the new element.

source-oriented

Compressed Sparse Column

Diagram illustrating a sorting step in a bubble sort algorithm. The top row shows the numbers 0, 2, 5, 8. The bottom row shows the numbers 0, 2, 1, 2, 3, 0, 2, 3, 1. Arrows indicate comparisons between adjacent elements in the top row and the bottom row. The elements 1, 2, 3, 0, 2, 3, 1 are highlighted in blue, and the elements a, d, b, i, h, c, f, g, e are highlighted in red.

target-oriented

Property Table Approaches

Triple Table

Subj.	Prop.	Obj.
ID1	type	BookType
ID1	title	"XYZ"
ID1	author	"Fox, Joe"
ID1	copyright	"2001"
ID2	type	CDType
ID2	title	"ABC"
ID2	artist	"Orr, Tim"
ID2	copyright	"1985"
ID2	language	"French"
ID3	type	BookType
ID3	title	"MNO"
ID3	language	"English"
ID4	type	DVDType
ID4	title	"DEF"
ID5	type	CDType
ID5	title	"GHI"
ID5	copyright	"1995"
ID6	type	BookType
ID6	copyright	"2004"

(Clustered) property table

Property Table

Subj.	Type	Title	copyright
ID1	BookType	"XYZ"	"2001"
ID2	CDType	"ABC"	"1985"
ID3	BookType	"MNP"	NULL
ID4	DVDType	"DEF"	NULL
ID5	CDType	"GHI"	"1995"
ID6	BookType	NULL	"2004"

Left-Over Triples

Subj.	Prop.	Obj.
ID1	author	"Fox, Joe"
ID2	artist	"Orr, Tim"
ID2	language	"French"
ID3	language	"English"

Property-Class Table

Class: BookType

Subj.	Title	Author	copyright
ID1	"XYZ"	"Fox, Joe"	"2001"
ID3	"MNP"	NULL	NULL
ID6	NULL	NULL	"2004"

Class: CDType

Subj.	Title	Artist	copyright
ID2	"ABC"	"Orr, Tim"	"1985"
ID5	"GHI"	NULL	"1995"

Left-Over Triples

Subj.	Prop.	Obj.
ID2	language	"French"
ID3	language	"English"
ID4	type	DVDType
ID4	title	"DEF"

Reduce numbers of subject-subject self joins necessary to reconstruct entities

Which of the following are distance measures in graphs?

Degree

Eccentricity

Diameter

Radius

Which of the following are distance measures in graphs?

Degree

Eccentricity

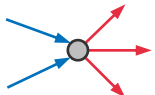
Diameter

Radius

Basic Concepts: Measures

Degree (Valency)

Degree of a vertex v

$\text{deg}_{in} = 2$  $\text{deg}_{out} = 3$

$$\text{deg}(v) = \text{deg}_{out}(v) + \text{deg}_{in}(v)$$

Degree of a graph G

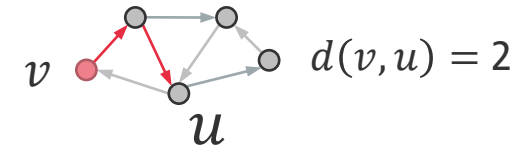
$$\text{deg}(G) = \max_{v \in V} \text{deg}(v)$$

Degree distribution

→ Probability distribution of vertex degree in a graph

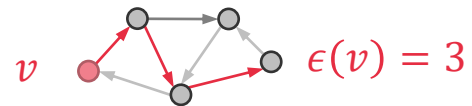


Distance



→ Number of edges in a shortest path connecting two vertices

Eccentricity



→ Longest shortest path starting from v

Average Distance

→ Average shortest path
→ Average eccentricity of any vertex in the graph

Radius

→ Minimum eccentricity of any vertex in the graph

$$r(G) = 2$$

Diameter

→ Maximum eccentricity of any vertex in the graph

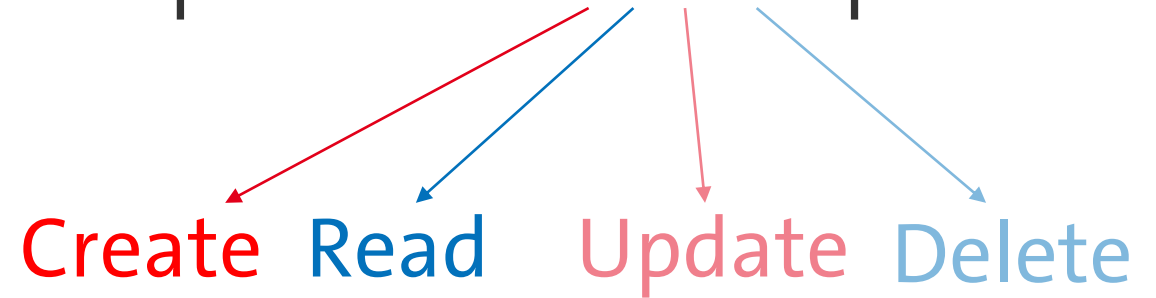


What do the letters CRUD (in relation to NoSQL DBs) stand for?

What do the letters CRUD (in relation to NoSQL DBs) stand for?

Key value store:

Basic key-value mapping with a simple API for **CRUD** operations



Which operations are offered by every key-value store?

Put

Get

Merge

Delete

Which operations are offered by every key-value store?

Put

Get

Merge

Delete

Key-Value Stores

Basic key-value mapping with a simple API for **CRUD** operations



Horizontal partitioning

users:1:a	4711
users:1:b	"[12, 34, 45, 67, 89]"
<hr/>	
users:2:a	01101010010110010101001...
users:2:b	"[12, ABC, 3212, 0xff]"

CRUD realized by at least 3 types of queries

Put: Add a new pair

Get: Retrieve a pair

Delete

Additional Operators implemented by some systems

Merge

MuliGet/MGet

MSet

...

Which are the variants for merging levels in an LSM tree?

Which are the variants for merging levels in an LSM tree?

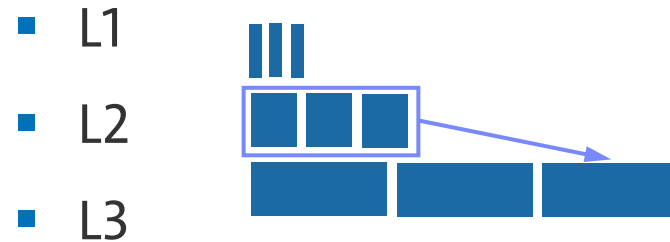
Tiering

Leveling

LSM Trees: Merging

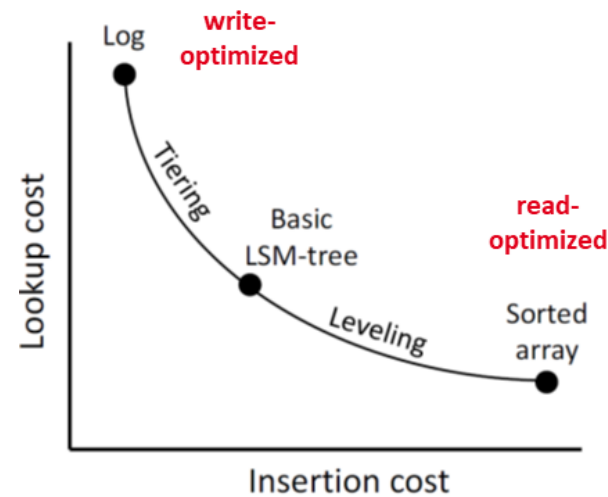
LSM Tiering

- Keep up to $T-1$ (sorted) runs per level L
- Merge all runs of L_i into 1 run of L_{i+1}

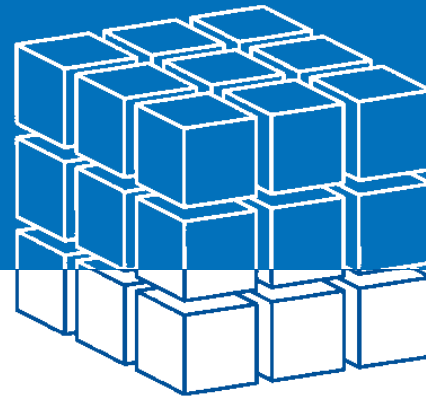


LSM Leveling

- Keep 1 (sorted) run per level L
- Sort-Merge run of L_i with L_{i+1}



Data Warehouses and OLAP



Which of the following are characteristics of Data Warehouses?

classical

Integration of data from different source systems

No user updates

Transaction-oriented organization of data

Historic data

Analysis-oriented organization of data

Real time data

Which of the following are characteristics of Data Warehouses?

classical

Integration of data from different source systems

No user updates

Transaction-oriented organization of data

Historic data

Analysis-oriented organization of data

Real time data



Characteristics of DWHs

Analysis-oriented organization of data

Domain-oriented, models a specific application goal

Integration of data from different source systems

Integrated database, integration on structural and data level of multiple databases

No user updates

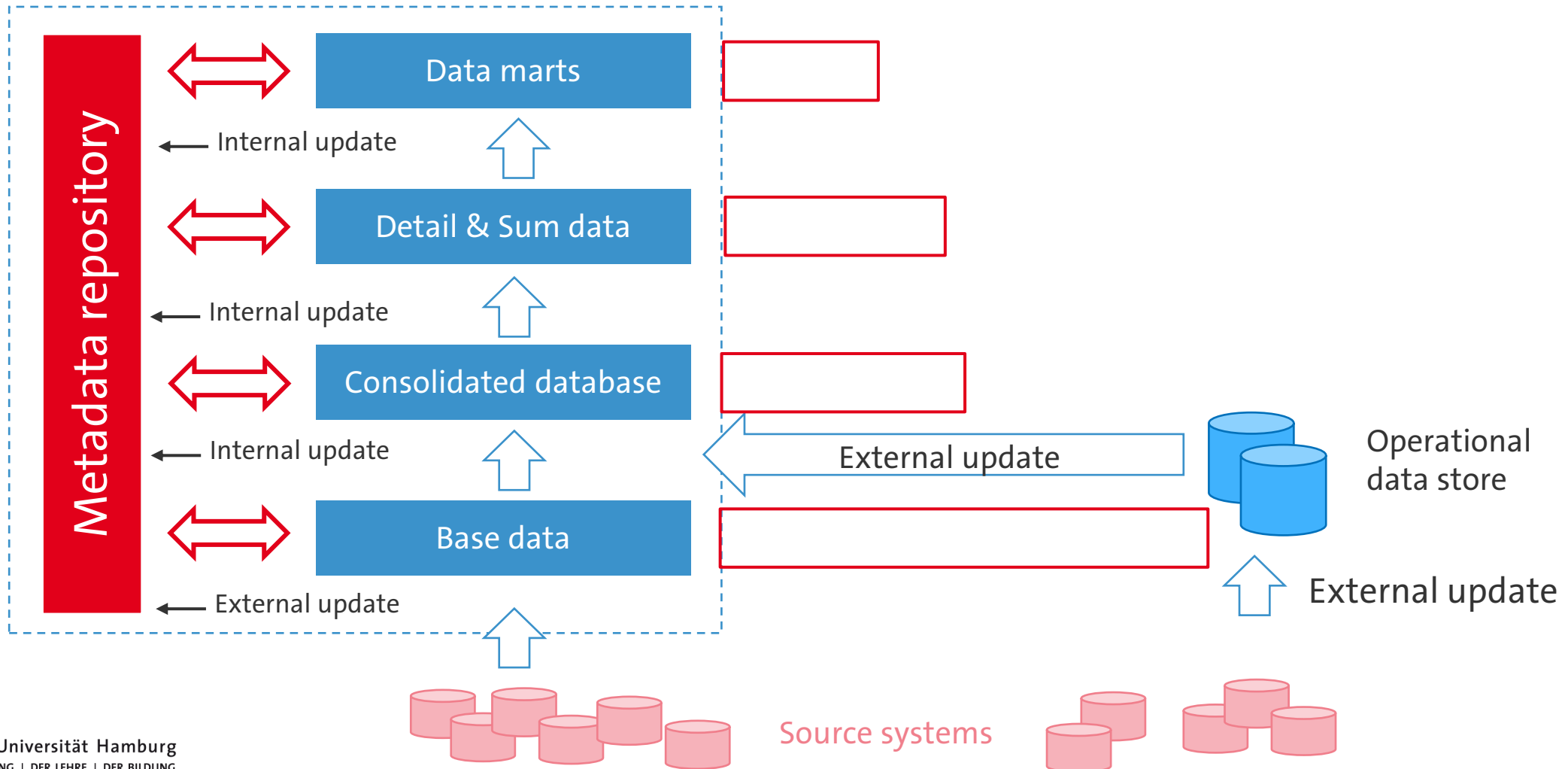
(almost) **no updates and deletes** (technical updates / deletes only, quality assurance)

Historic data

Data is kept over a long period of time

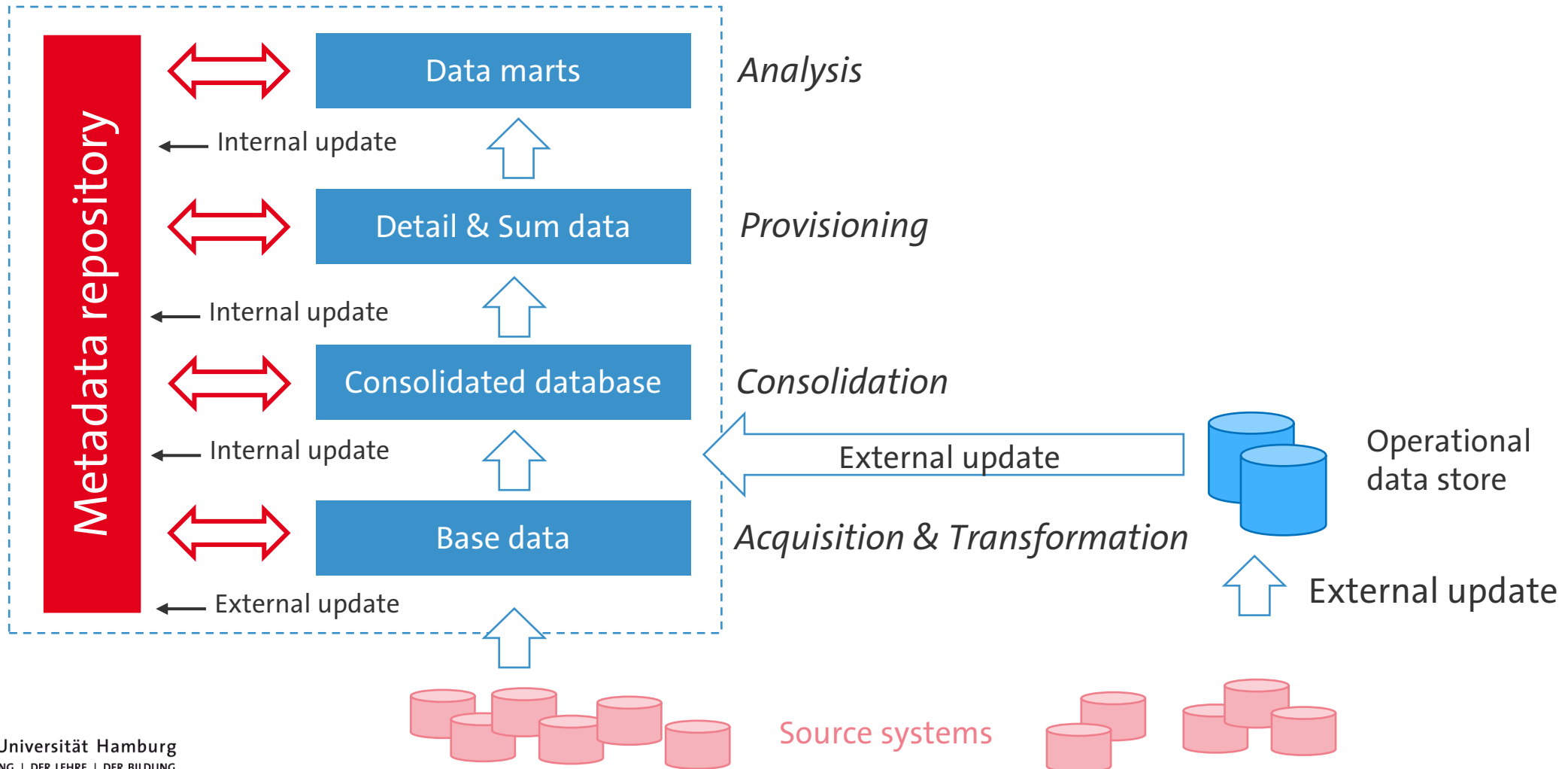
Fill in the steps involved in creating the Data Warehouse

Data-Warehouse-System



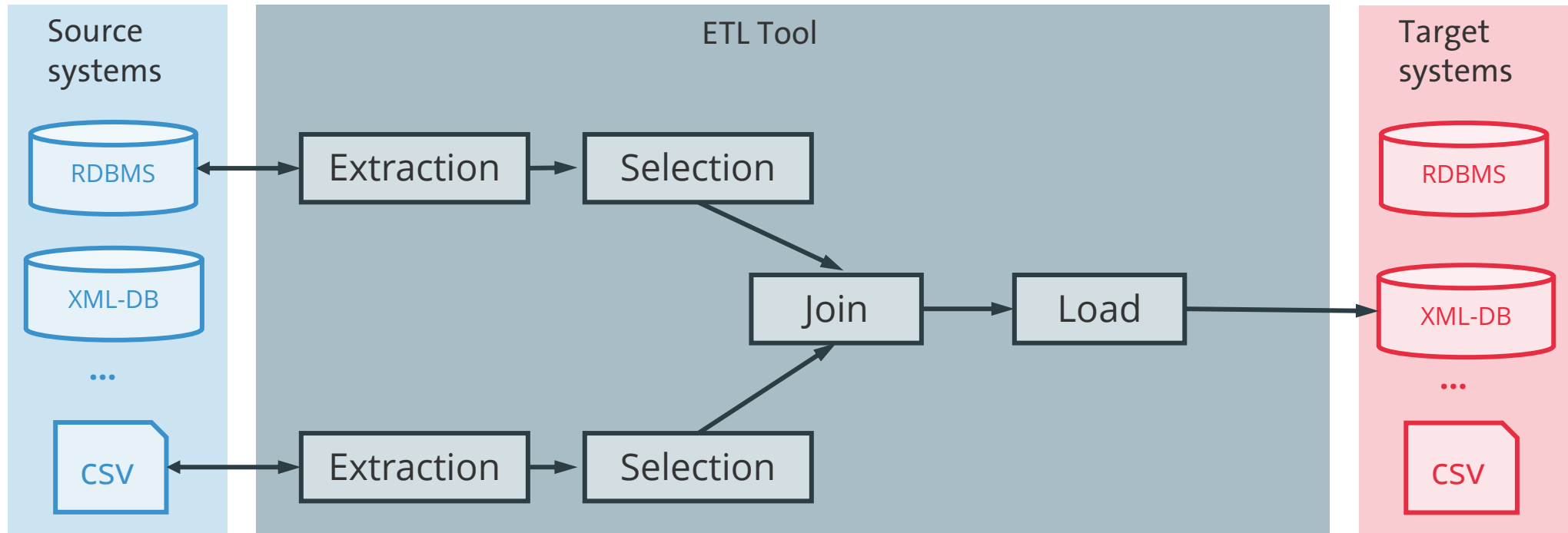
Fill in the steps involved in creating the Data Warehouse

Data-Warehouse-System



What happens during the ETL step?

Data acquisition and transformation



Which of the following statements about ADAPT are wrong?

Different hierarchies of the same dimension can have different leaf nodes.

Attributes can be assigned to a whole cube, a dimension, and individual levels.

Self-precedence can be modelled

A cube must have exactly 3 dimensions.

Which of the following statements about ADAPT are wrong?

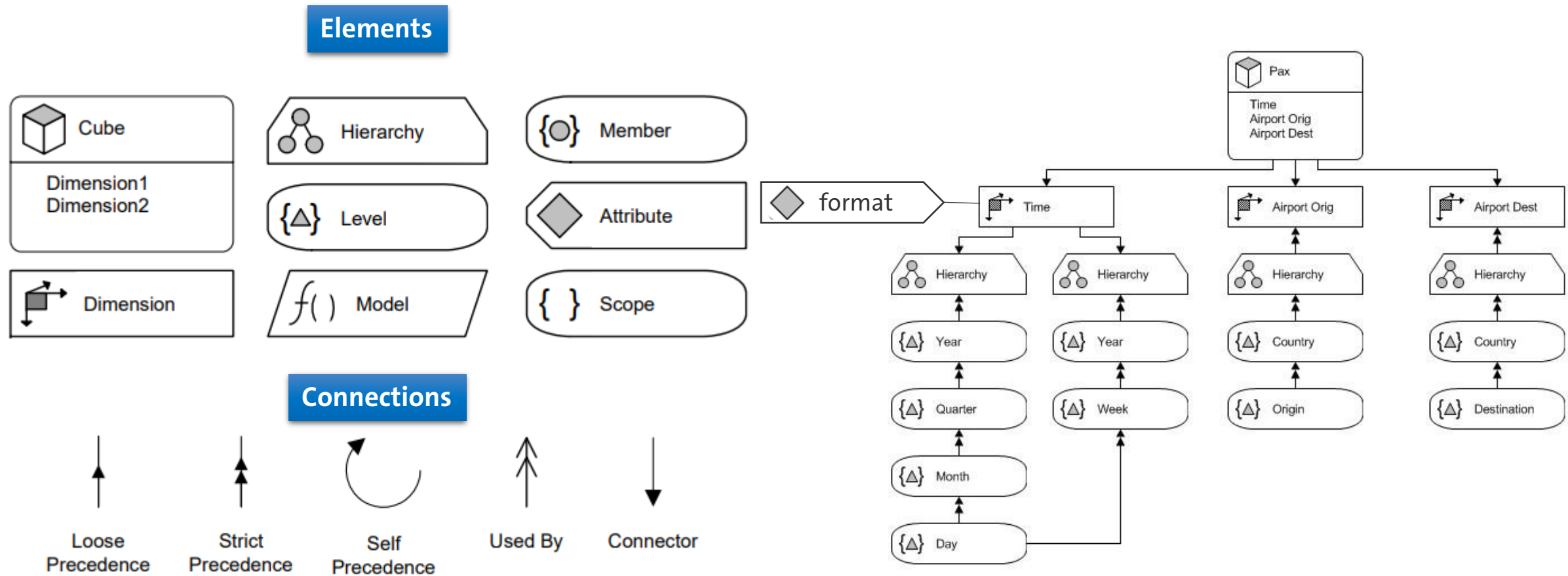
Different hierarchies of the same dimension can have different leaf nodes.

Attributes can be assigned to a whole cube, a dimension, and individual levels.

Self-precedence can be modelled

A cube must have exactly 3 dimensions.

Conceptual Models: ADAPT



Name two relational schemas typically used in OLAP and briefly explain them.

Name two relational schemas typically used in OLAP and briefly explain them.

Star schema

- One fact table and multiple dimension tables
- Dimension tables are not necessarily normalized
- Fact table is linked directly to all dimension tables via a foreign key

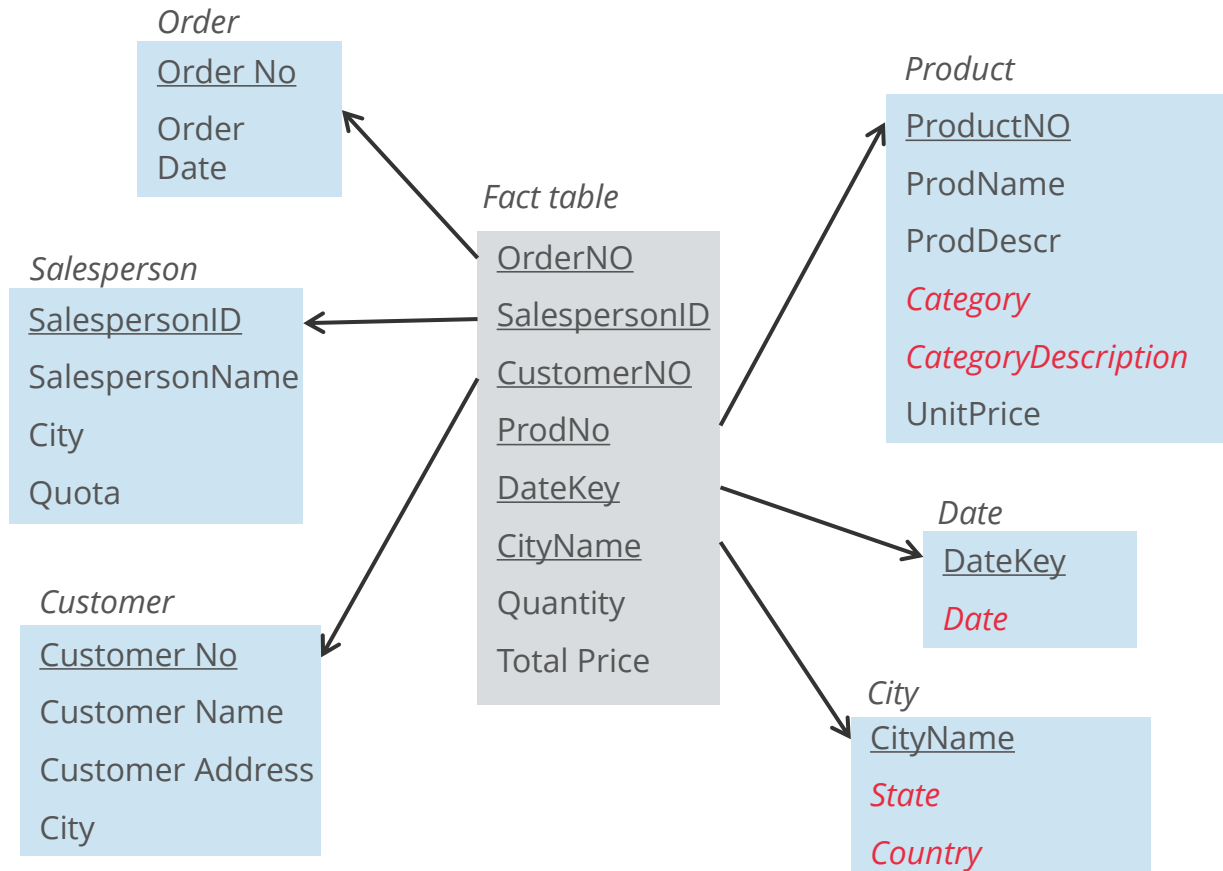
Snowflake schema

- One fact table and multiple dimension tables
 - Smaller memory footprint than star schema
- Dimension tables are normalized
 - Introduces more joins in the queries

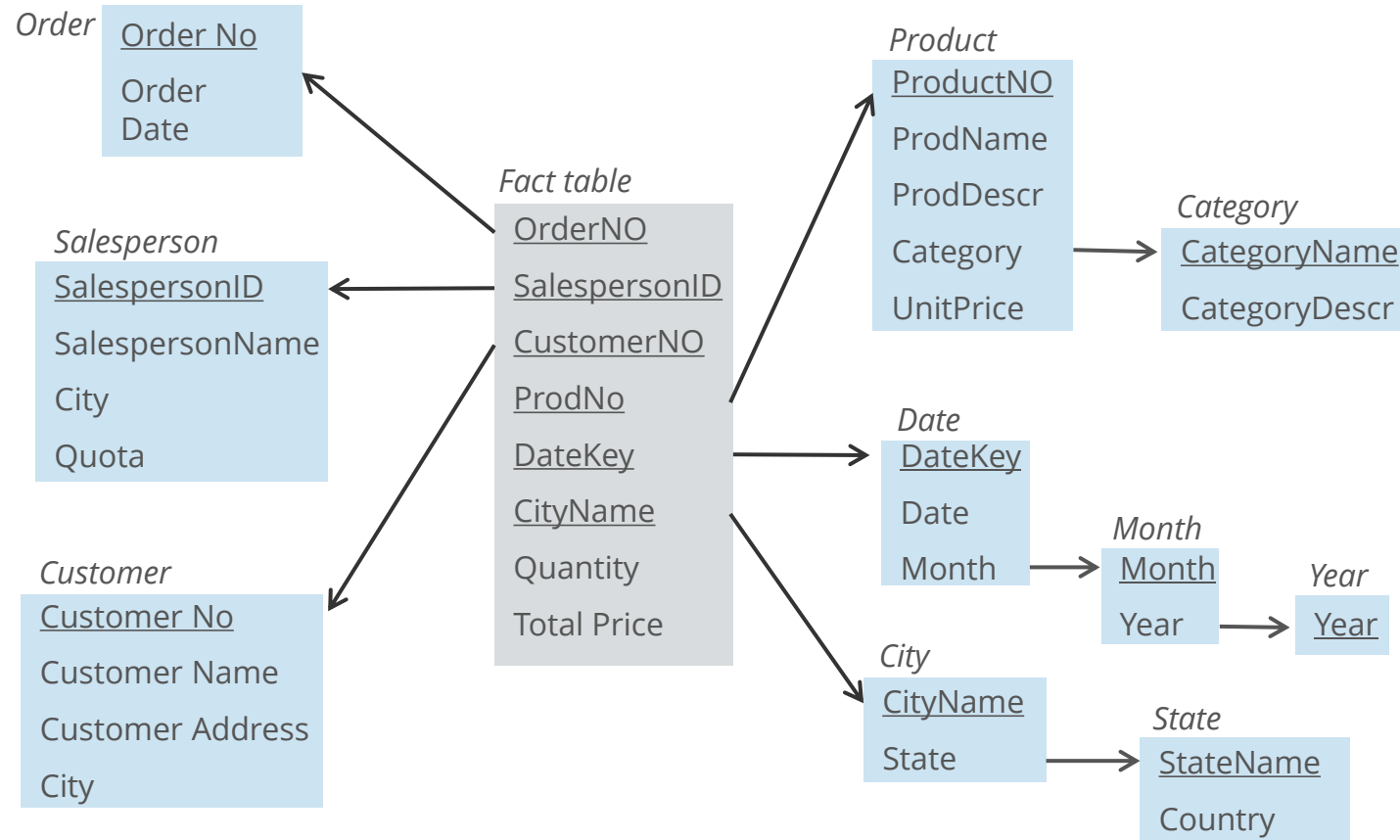
Galaxy schema

- Multiple fact tables
- Dimension tables are shared
- Joining fact tables (which are usually huge) is extremely inefficient
 - Long query execution times

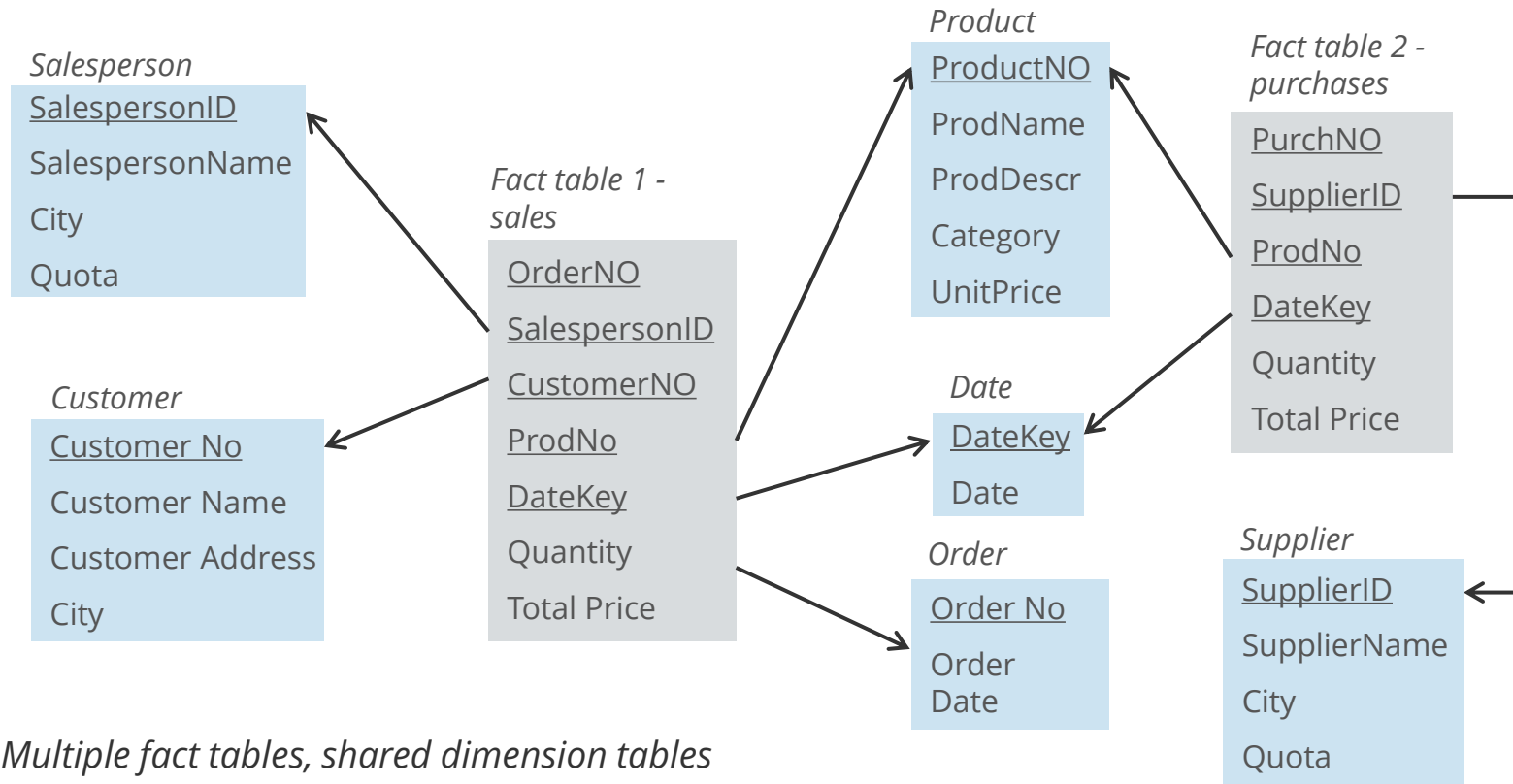
Star Schema



Snowflake Schema



Galaxy Schema



Which of the following is/are not (a) typical method(s) of relational mapping of hierarchies?

Horizontal Mapping

Recursive Vertical Mapping

Diagonal Mapping

Iterative Horizontal Mapping

Vertical Mapping

Random Mapping

Which of the following is/are not (a) typical method(s) of relational mapping of hierarchies?

Horizontal Mapping

Recursive Vertical Mapping

Diagonal Mapping

Iterative Horizontal Mapping

Vertical Mapping

Random Mapping

Relational Mapping of Hierarchies

Vertical Mapping

Product			Product group			Product family			Category	
ID	Product	ID2	ID	Product group	ID3	ID	Product family	ID4	ID	Category
1	Flour X	1	1	Flour	1	1	Bakery goods	1	1	Food
2	Sugar Y	2	2	Sugar	1	2	Beverages	1		
3	Water Z	3	3	Water	2					

Horizontal Mapping

Relation product				
ID	Product	Product group	Product family	Category
1	Flour X	Flour	Bakery goods	Food
2	Sugar Y	Sugar	Bakery goods	Food
3	Water Z	Water	Beverages	Food

Recursive Vertical Mapping

ID	Product	Product group	Product family	Category
1	Flour X	Flour	Bakery goods	Food
2	Sugar Y	Sugar	Bakery goods	Food
3	Water Z	Water	Beverages	Food



Relation product		
ID	Product	ParentID
1	Food	NULL
2	Bakery goods	1
3	Beverages	1
4	Flour	2
5	Sugar	2
6	Water	3
7	Flour X	4
8	Sugar Y	5
9	Water Z	6

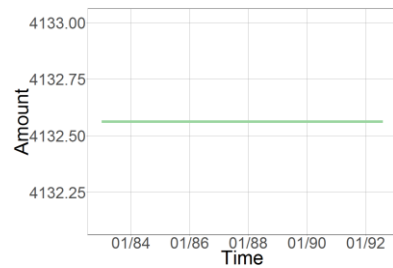


Data Mining



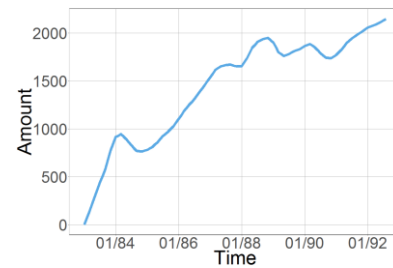
Into which components can time series typically be decomposed?

Into which components can a time series typically be decomposed?



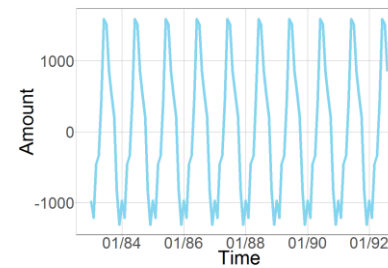
Base

+



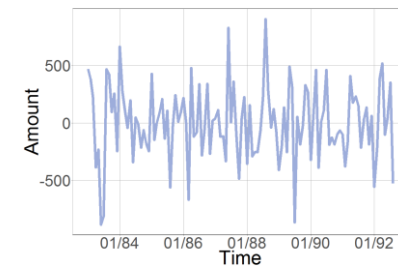
Trend

+



Season

+



Residuals

Is Piecewise aggregate approximation (PAA) a shape-based or a feature based representation for time series?

Shape-based

Feature-based

What is the difference to symbolic aggregate approximation?

Is Piecewise aggregate approximation (PAA) a shape-based or a feature based representation for time series?

Shape-based

Feature-based

What is the difference to symbolic aggregate approximation?

In SAX, values are discretized into an alphabet.

Shape-based Representation

Data Types

Metrics

Handling of Missing Data

Outlier Detection

Dimensionality Reduction

Value Count Reduction

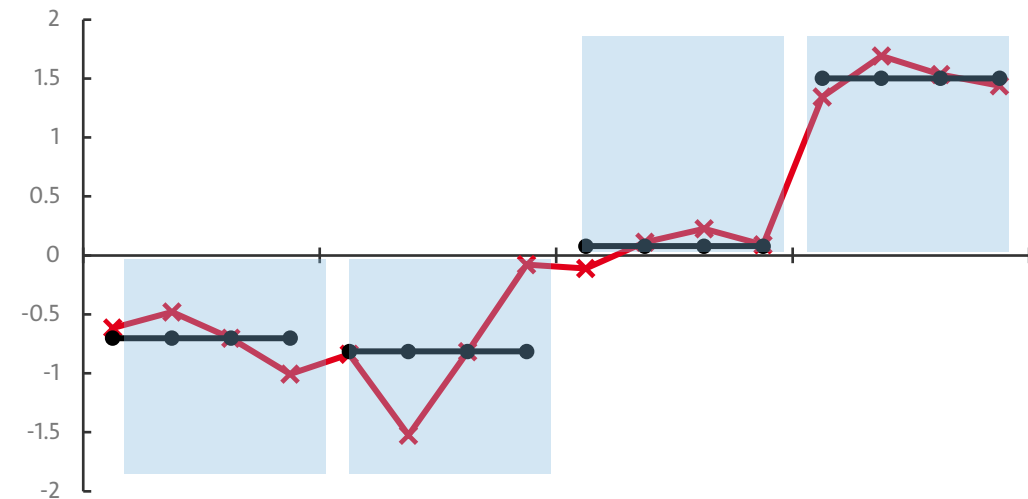
- Transform a series to segments
- Time-dependent representation
 - High compression



Piecewise aggregate approximation (PAA)
→ Represent each segment by its mean



Symbolic aggregate approximation (SAX)
→ Discretize the mean into an alphabet

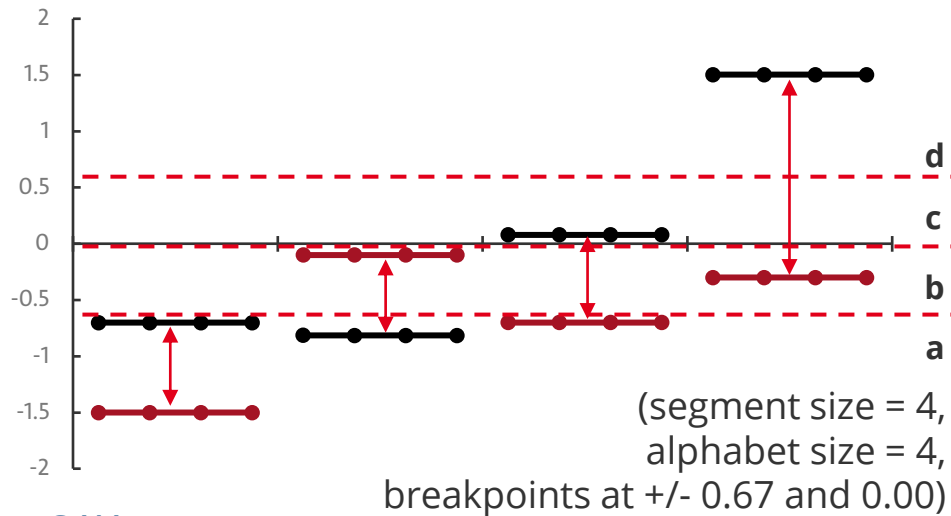


PAA: -0.70 -0.81 0.08 1.56

SAX: b b a a
(segments per series = 4, alphabet size = 2, breakpoint at 0.0)

The SAX Distance

Segments



SAX

\hat{x}^1	a	a	c	d
\hat{x}^2	a	b	a	b

$$d_{SAX}(\hat{x}^i, \hat{x}^j) = \sqrt{T/W} \cdot \sqrt{\sum_{w=1}^W cell(\hat{x}_w^i, \hat{x}_w^j)^2}$$
 $cell$ returns the minimum distance of two symbols:

	a	b	c	d
a	0	0	0.67	1.34
b	0	0	0	0.67
c	0.67	0	0	0
d	1.34	0.67	0	0

PCA

Standardize all values per dimension

value-mean/standard deviation

$$\text{std} \left[\begin{array}{|c|c|c|c|} \hline \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \hline \end{array} \right] = \begin{array}{|c|c|c|c|} \hline \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array}$$

*Standardized data set * Feature vector
= New Dataset (reduced to x dimensions)*

$$\begin{array}{|c|c|c|c|} \hline \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array} * \begin{array}{|c|c|} \hline \text{red} & \text{black} \\ \hline \end{array} = \begin{array}{|c|c|} \hline \text{pink} & \text{pink} \\ \hline \end{array}$$

Compute the covariance matrix

$\text{cov}(x,y) = (\sum (x - \text{mean}(x))(y - \text{mean}(y))) / \text{number of data points}$

$$\text{cov} \left[\begin{array}{|c|c|c|c|} \hline \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array} \right] = \begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} & \text{pink} \\ \hline \end{array}$$

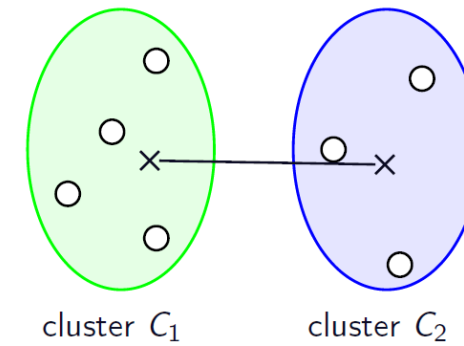
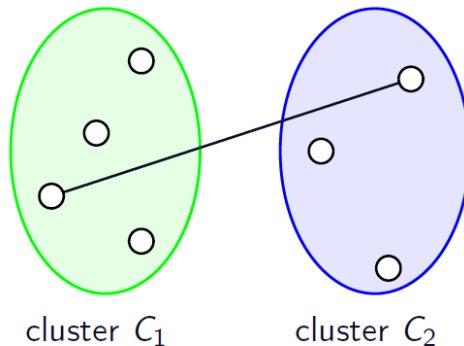
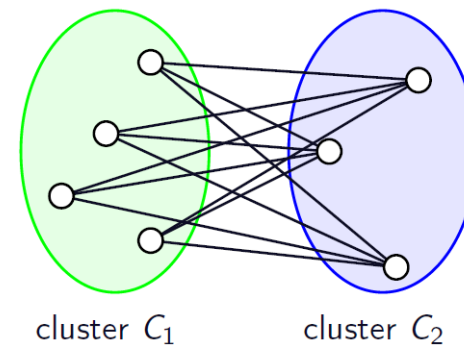
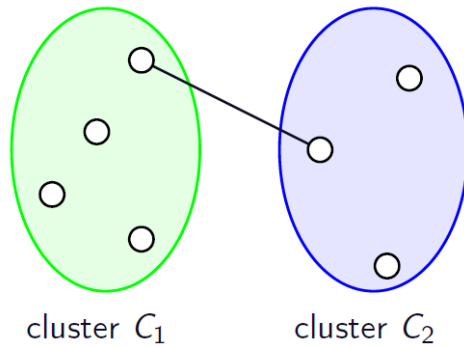
Compute Eigen values and Eigen vectors of the covariance matrix

$$\begin{array}{|c|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} & \text{pink} \\ \hline \end{array} v = \lambda v \rightarrow \begin{array}{|c|c|c|c|} \hline \lambda & \text{black} & \text{black} & \text{black} \\ \hline v & \text{black} & \text{black} & \text{black} \\ \hline \end{array}$$

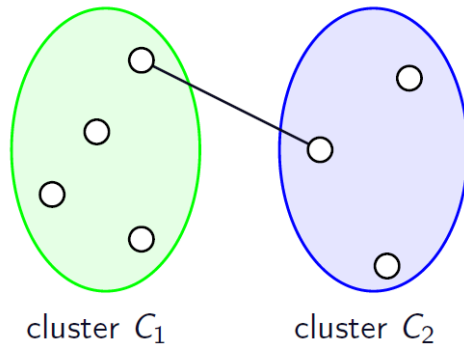
Sort elements of Eigenvectors in descending order
and pick those that belong to top x Eigen values
→ Feature Vectors

$$\text{sort} \downarrow \begin{array}{|c|c|c|c|} \hline \text{red} & \text{black} & \text{black} & \text{black} \\ \hline \end{array}$$

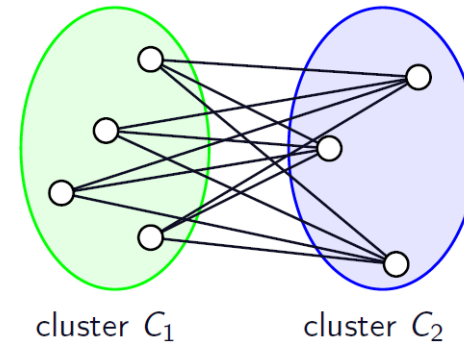
Which distance measures for clusters are shown here?



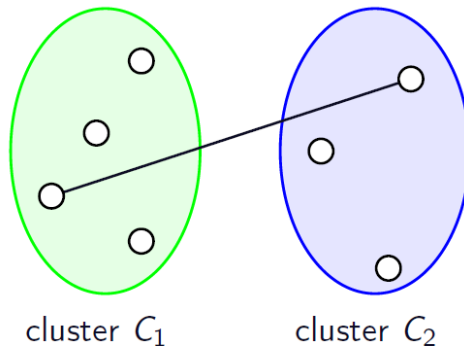
Which distance measures for clusters are shown here?



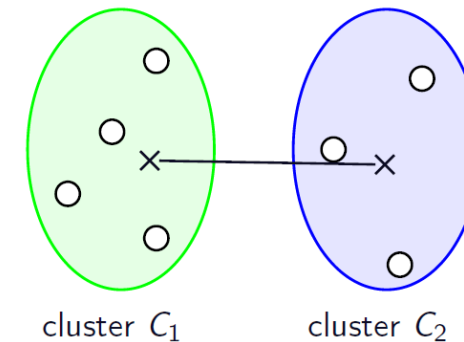
Single link



Average link



Complete link



Canonical entity

Clustering Methods

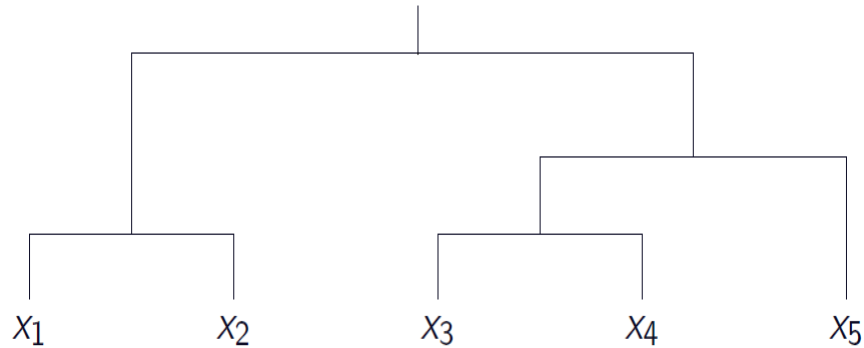
K-Means

Canopy Clustering

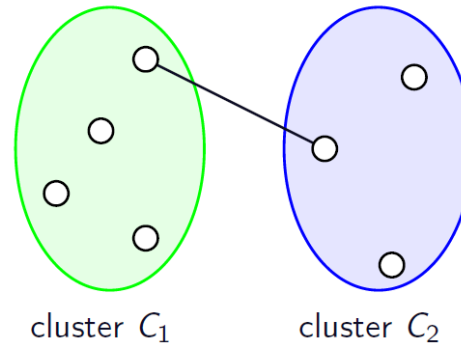
Hierarchical Clustering

Incremental Clustering

Results can be displayed as a dendrogram

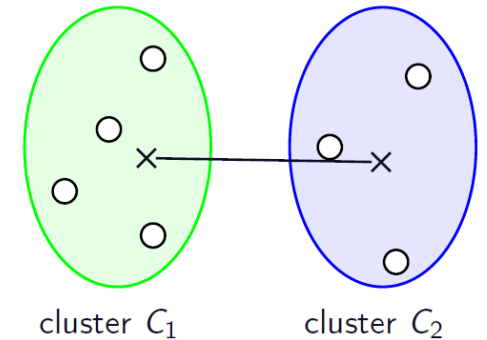


Distance measures for clusters



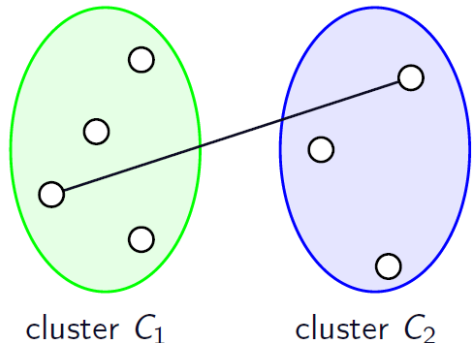
Single link

Minimal distance between two data points



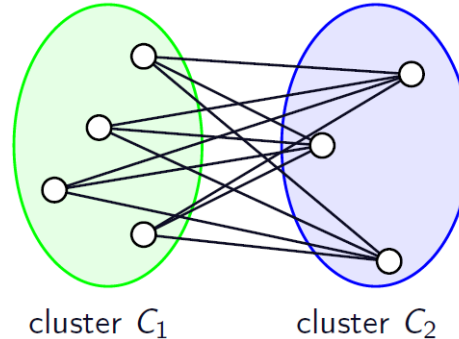
Canonical entity

Distance between two cluster representatives (e.g. the centroids)



Complete link

Maximal distance between two data points

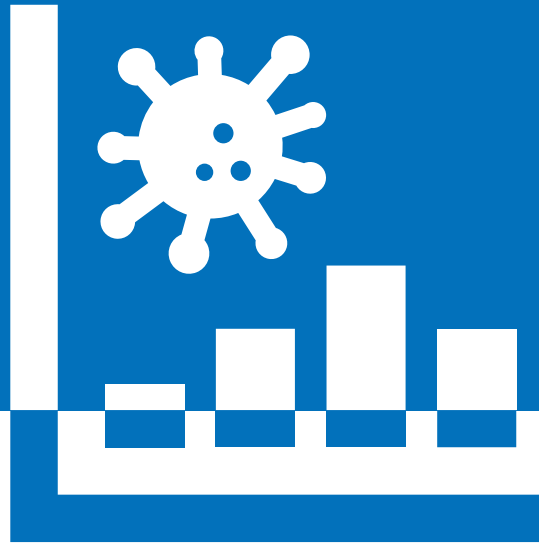


Average link

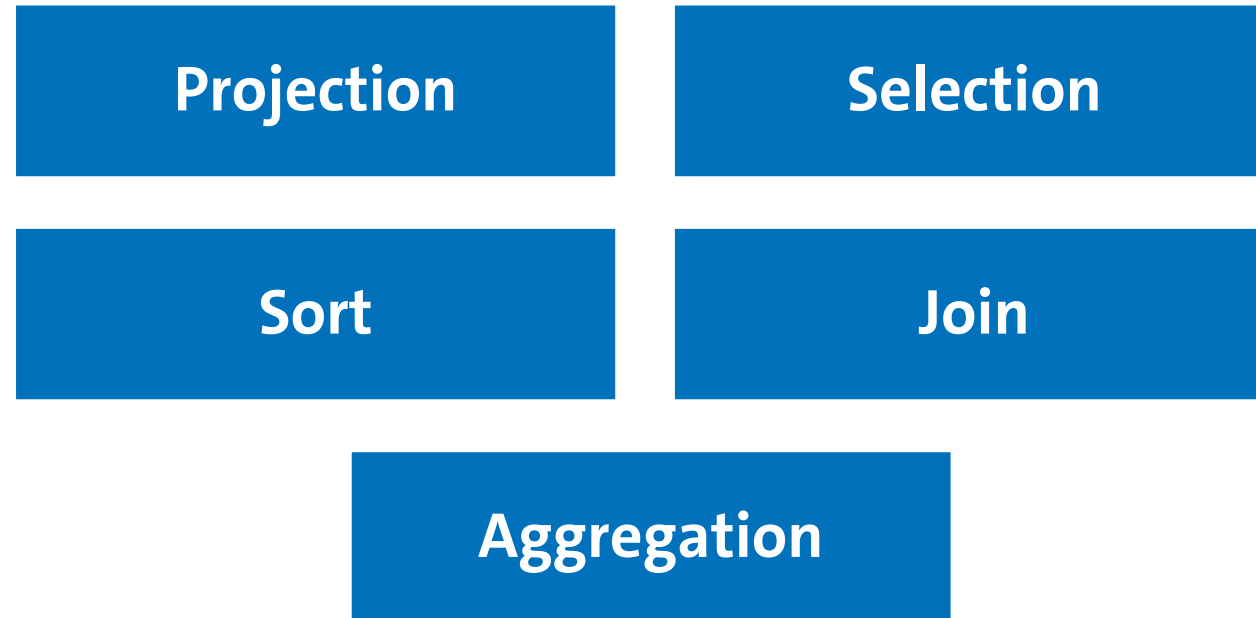
Average distance between two data points



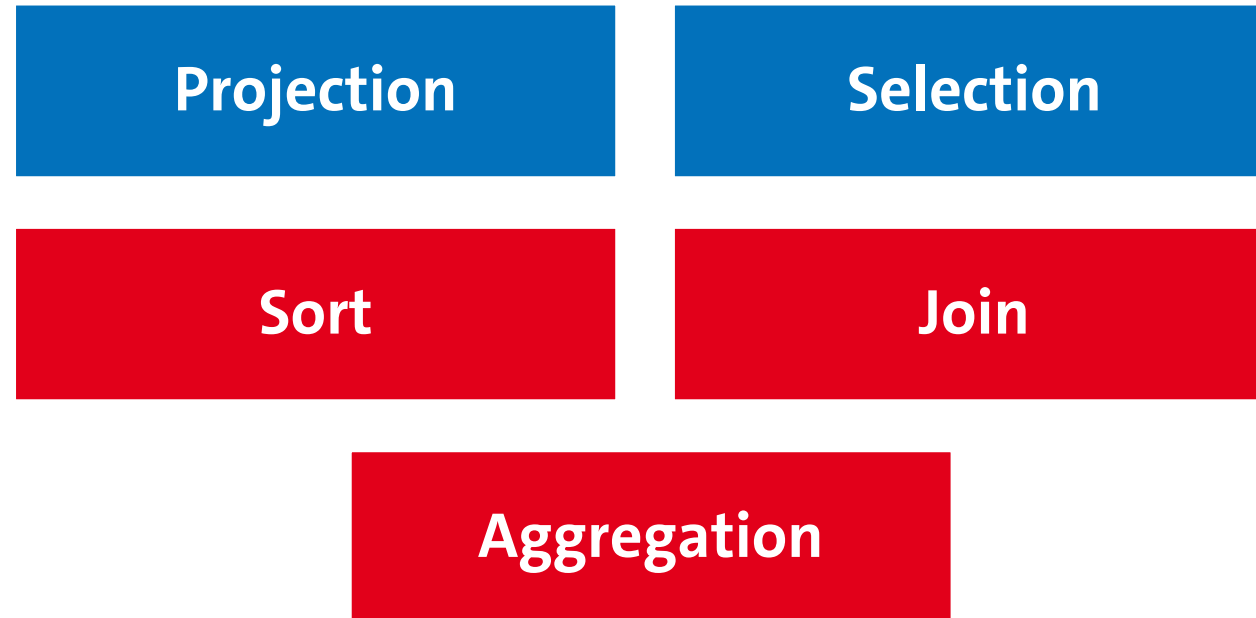
Big Data Analysis



Which operators require shuffling if they are distributed?

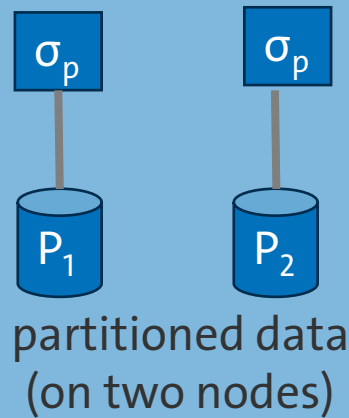


Which operators require shuffling if they are distributed?



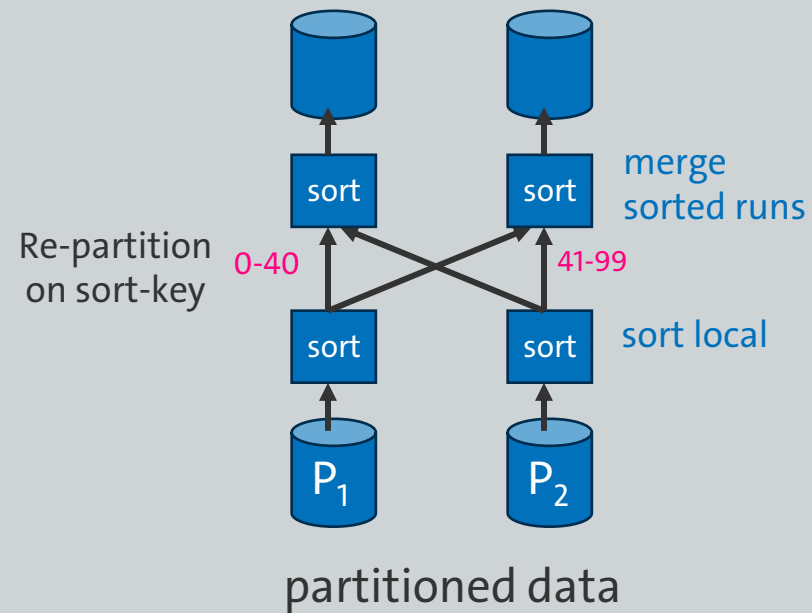
Parallel Operators

Selection

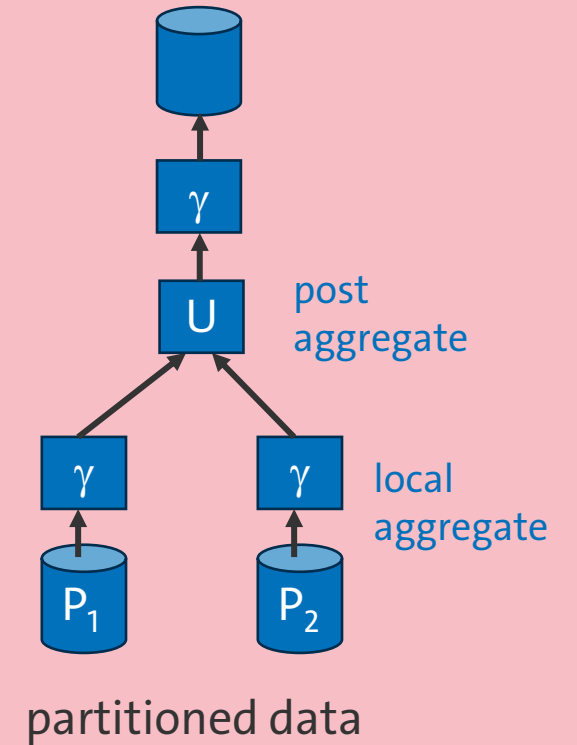


→ Projection works similarly as selection

Sort



Aggregation

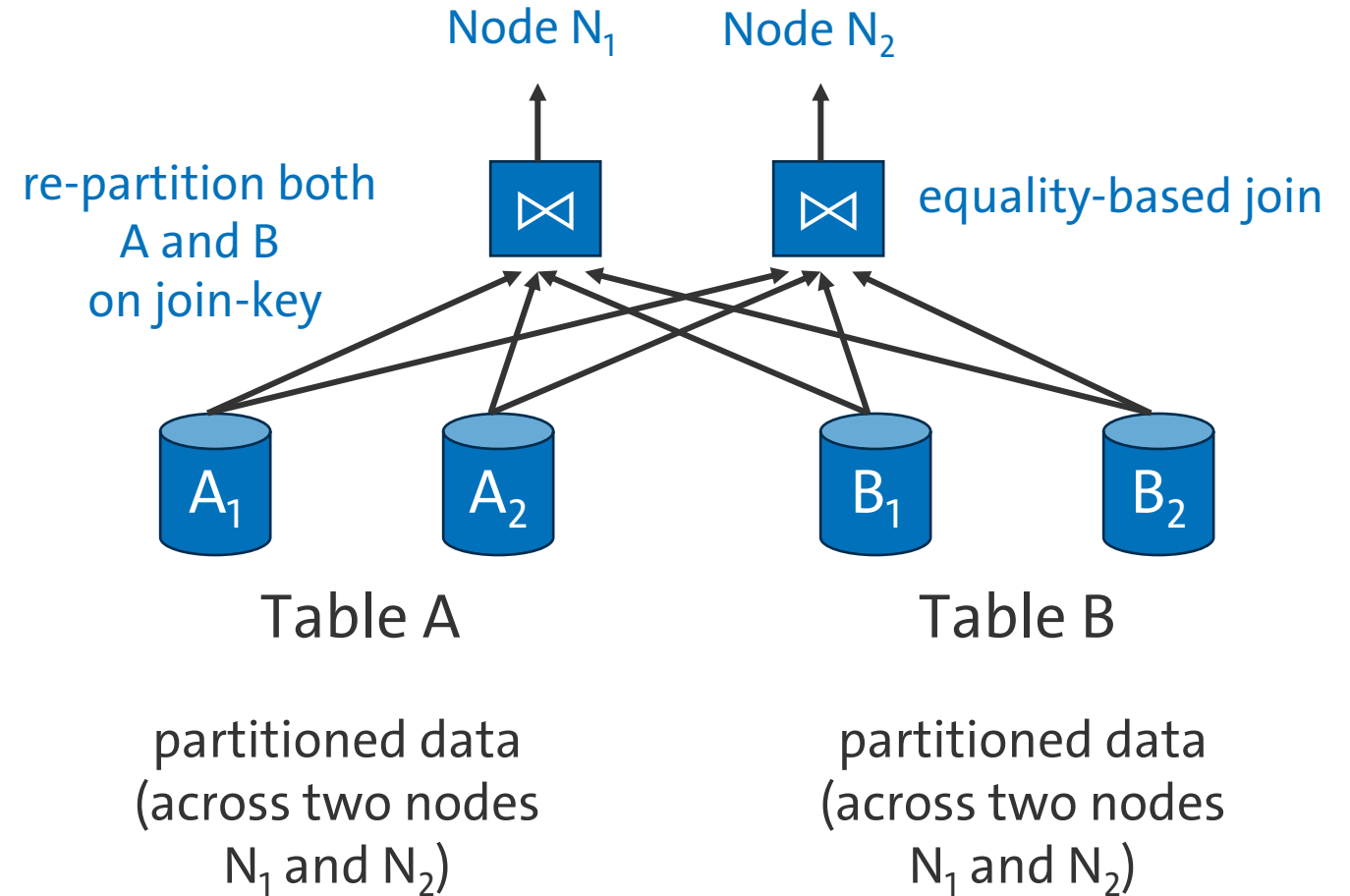


Join: Symmetric Repartitioning

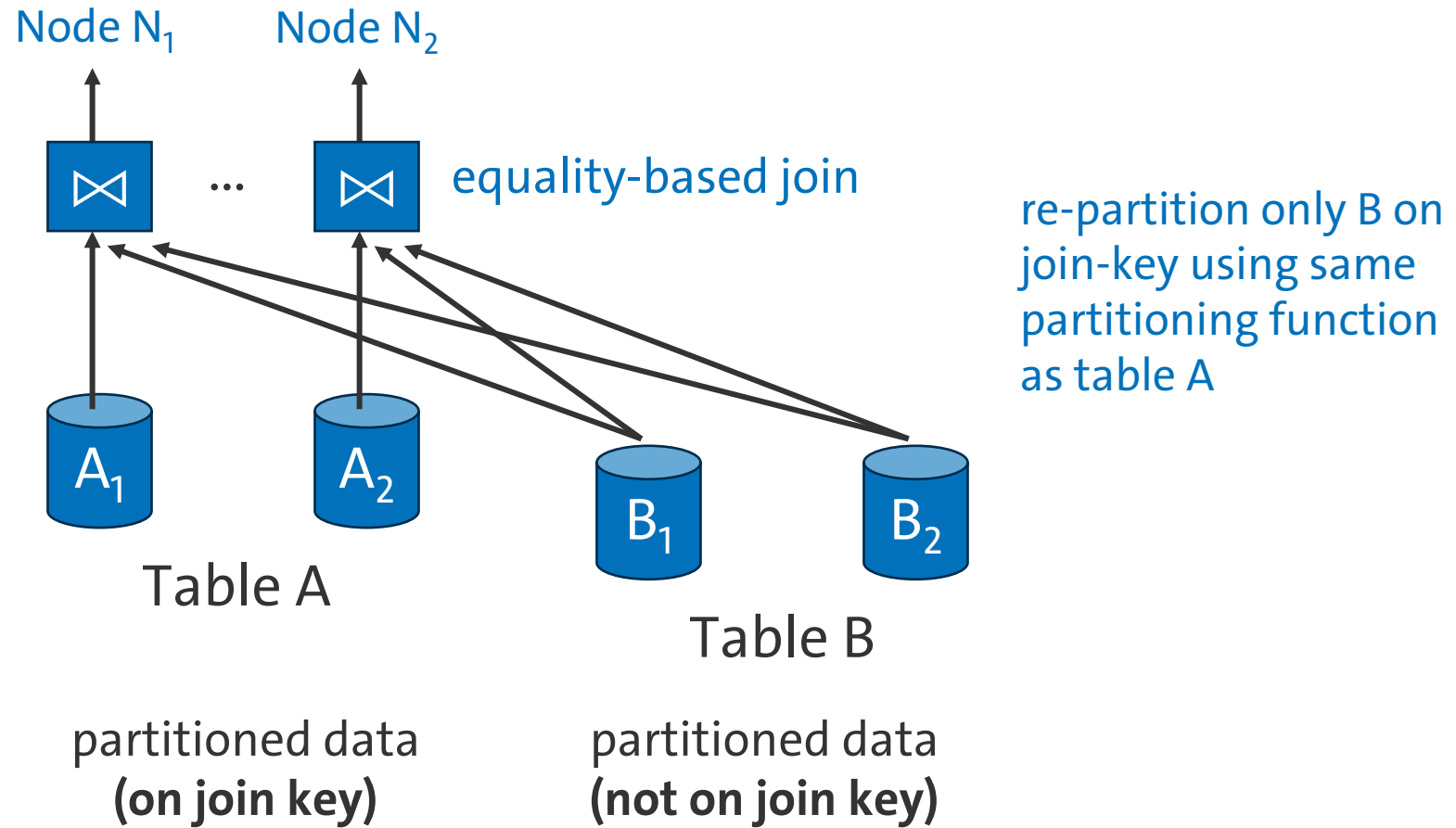
„Repartitioned join“

- Both join partners are repartitioned after the join attribute
- High communication costs

→ Avoid!

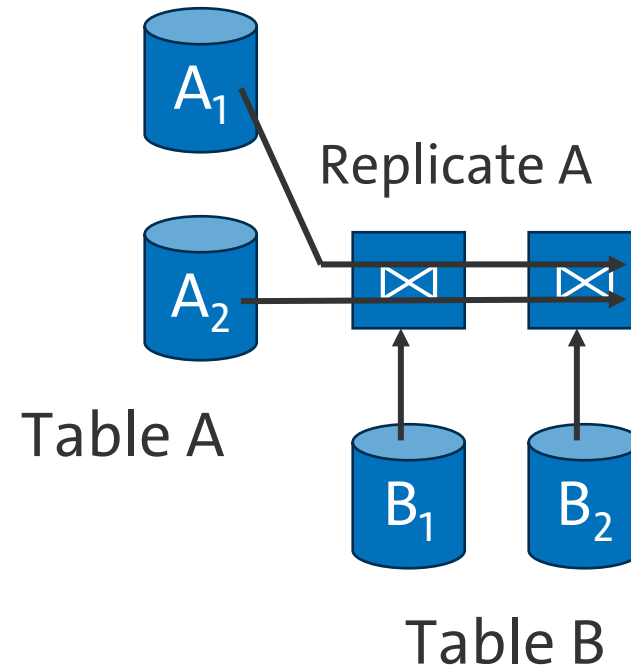


Join: Asymmetric Repartitioning



Join: Fragment and Replicate

partitioned data
(across two nodes
 A_1 and A_2)



Replicate (all fragments of) smaller table
(here: table A) to all nodes

partitioned data
(across two nodes B_1 and B_2)

Which of the following are shuffling methods for distributed operators?

Range-based N:M

Hash-based N:M

Column-based N:M

1:1

N:1

Remote N:1

Which of the following are shuffling methods for distributed operators?

Range-based N:M

Hash-based N:M

Column-based N:M

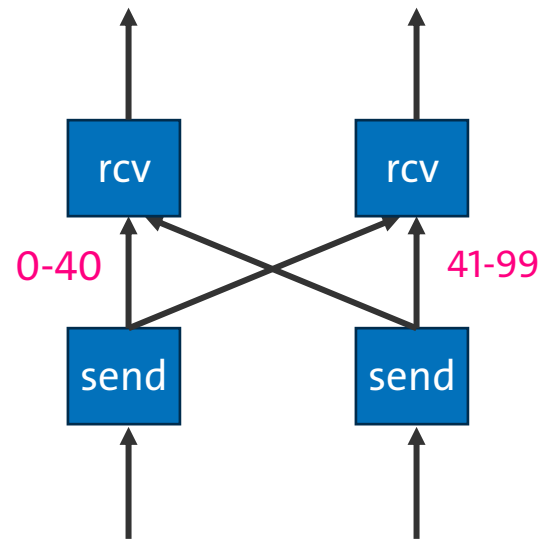
1:1

N:1

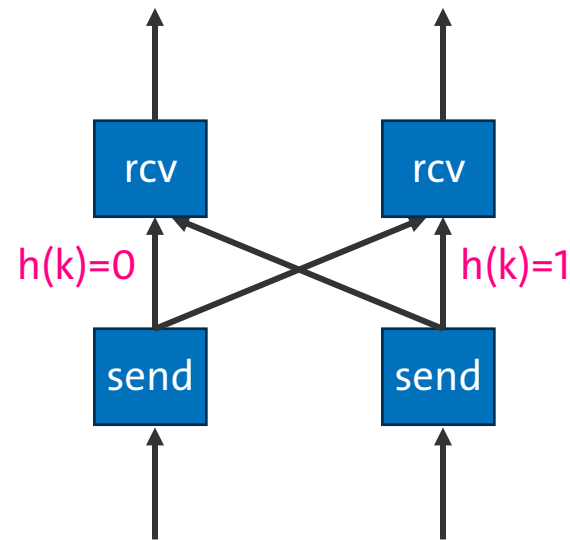
Remote N:1

Data Shuffling: Details

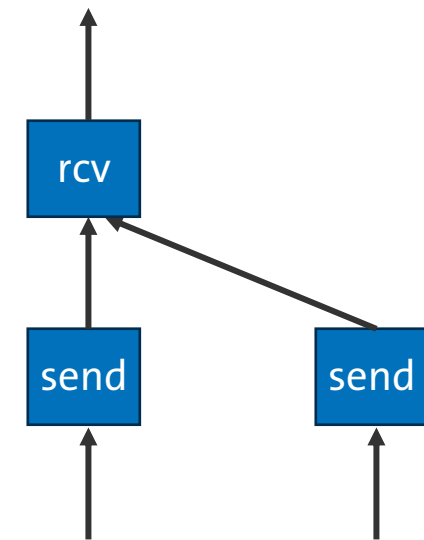
Data Shuffling can use different partitioning strategies (range vs. hash-partitioning, N:M vs. N:1) to re-partition data during query execution



Range-based N:M



Hash-based N:M



N:1 (no part. function)

Name 3 possible reasons for data skew

Name 3 possible reasons for data skew

Different partition sizes

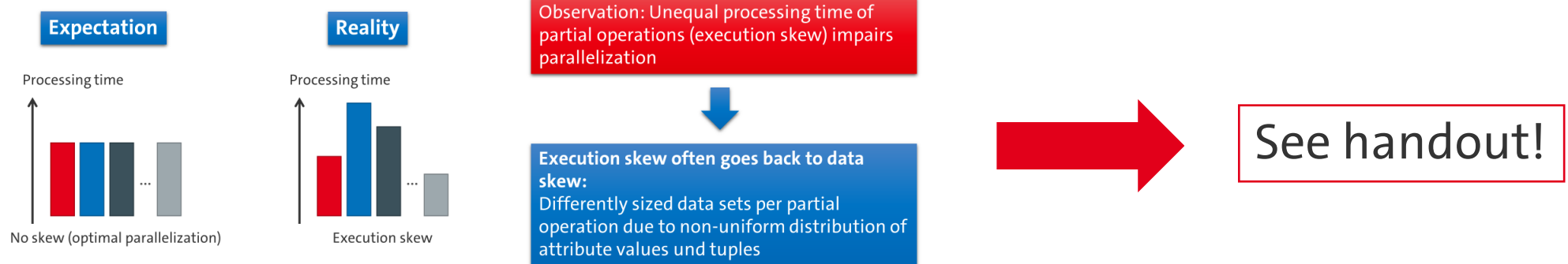
Distribution function leading to different fragment sizes

Different hit rates per partition in range queries

...

Name 3 possible reasons for data skew

The Data Skew



What is the main difference between the Kappa architecture and the Lambda architecture?

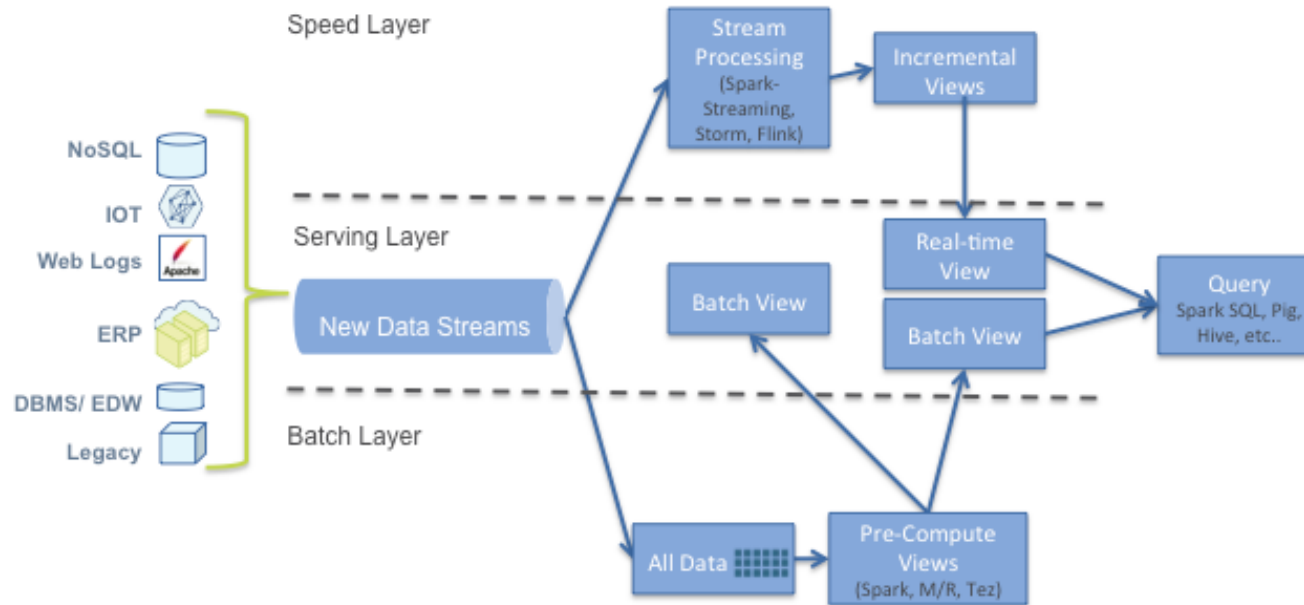
What is the main difference between the Kappa architecture and the Lambda architecture?

The Kappa architecture only has a speed layer while the Lambda architecture also has a batch layer

Lambda Architecture

→ see lecture 16

→ Batch & Stream processing

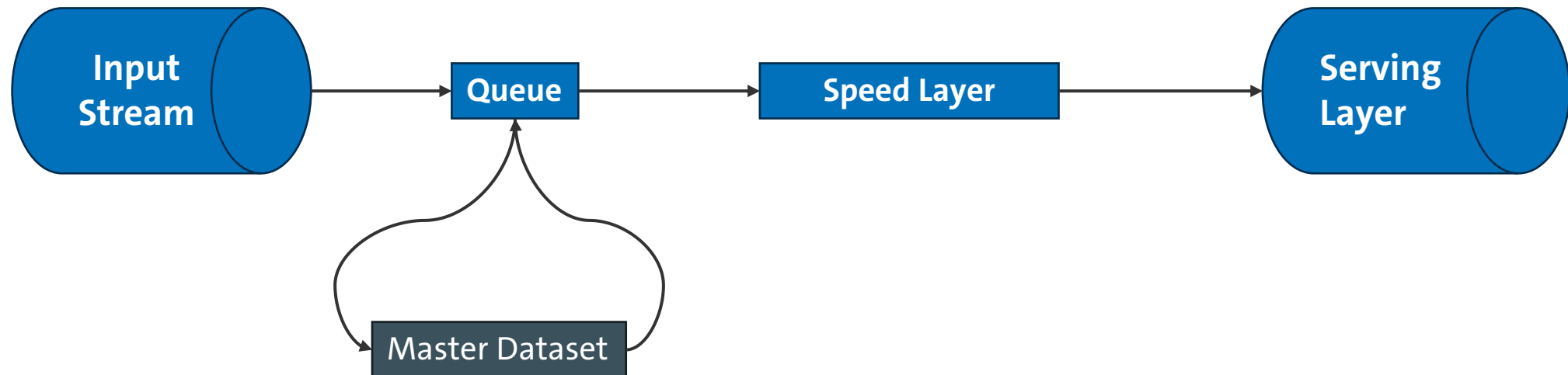


Disadvantage

2 code bases & 2 deployments, e.g. Hadoop & Storm

Kappa Architecture

→ No batch layer!



Disadvantage
Real time processing only

Which functions must be implemented by the application developer when using a map-reduce framework?

Map

Shuffle

Reduce

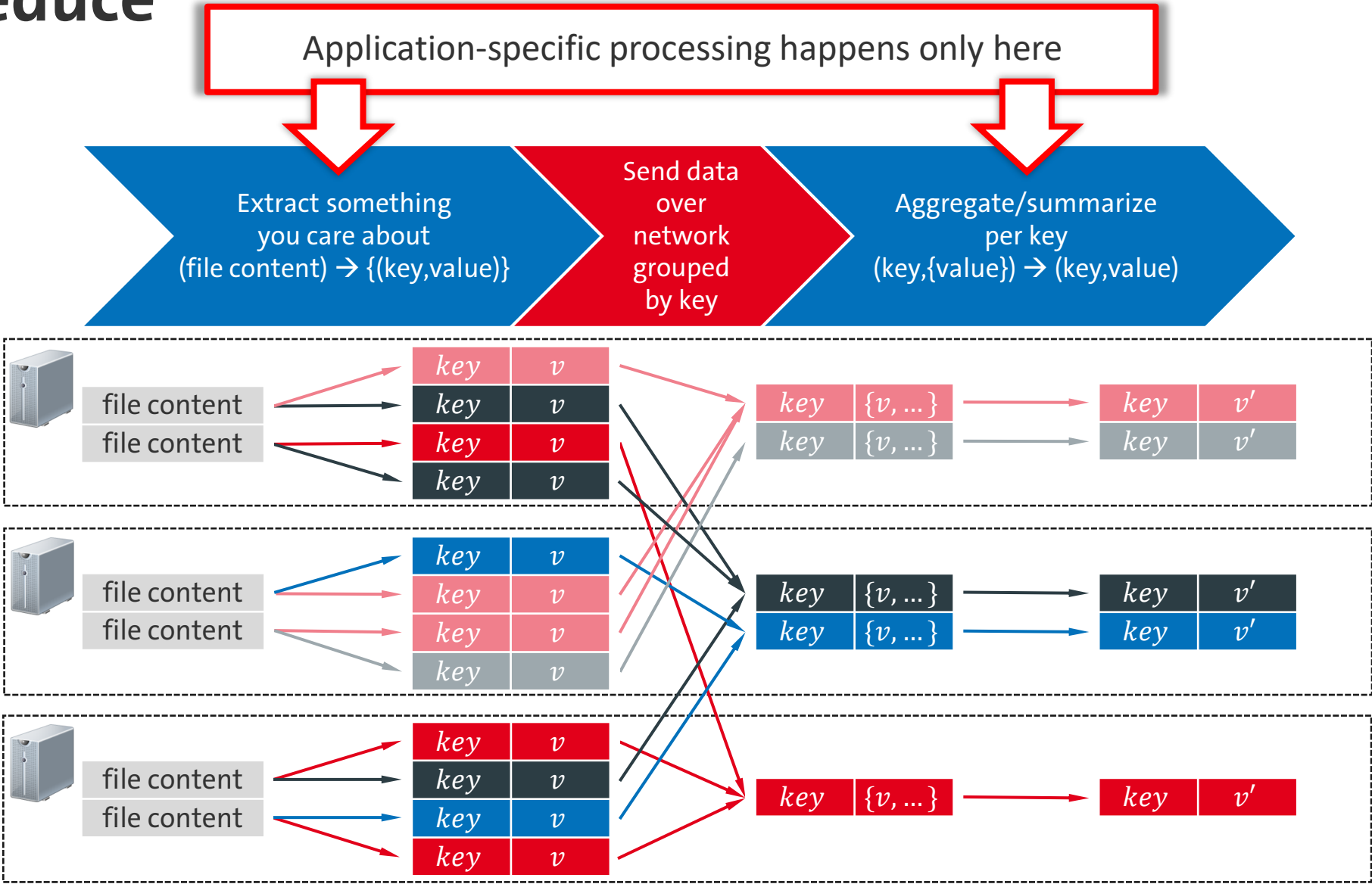
Which functions must be implemented by the application developer when using a map-reduce framework?

Map

Shuffle

Reduce

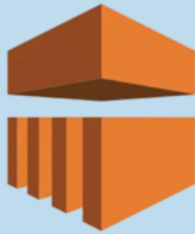
Map Reduce



Name 3 frameworks for batch processing

Framework Overview

Batch



amazon
EMR

Spark



Flink

Stream



APACHE
STORM™

samza



kafka