

Data Integration in Data Lakes

Rihan Hai

Web Information Systems

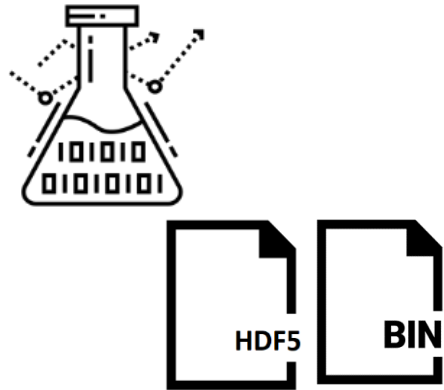
27.06.2022

Data Lake survey

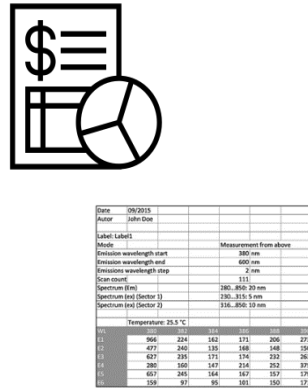


Problems of Big Data Management

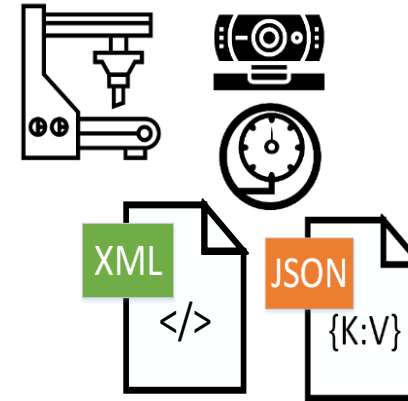
Scientific Experiments



Business Transactions



Industrial Production



Mobile Applications



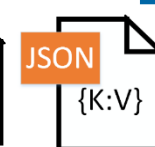
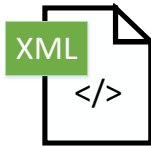
Data Lake



What Is a Data Lake?

A *data lake* is a flexible, scalable **data storage and management system**, where **raw data** from **heterogeneous** sources can be ingested and stored in their **original** format, and later queried in an **on-the-fly** manner.

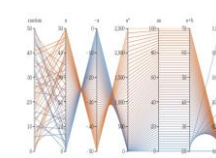
Large scale of
heterogeneous data



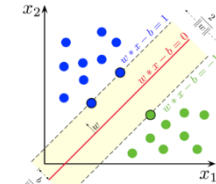
Data Lake



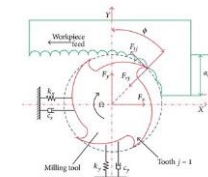
Insight



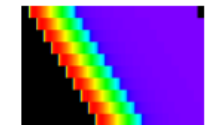
Visual Analytics



Data Analytics



Modeling



Simulation

Our contribution in data lake solution landscape

Table 1: Classification of data lake solutions based on functions

Tier	Functions	Systems
Ingestion	Metadata extraction	GEMMS [112]
		DATAMARAN [49]
		Skluma [129]
	Metadata modeling	GEMMS [61], [112]
		HANDLE [40]
		Data vault [54], [101]
		Diamantini et al. [31], [32], [33]
		Aurum [45]
	Sawadogo et al. [122]	
Maintenance	Dataset organization	GOODS [64], [65]
		DS-Prox [4], [5], [6]
		KAYAK [85], [86]
		Nargesian et al. [98]
		Ronin [105]
		Juneau [143]
	Related dataset discovery	Aurum [45]
		Brackenbury et.al. [14]
		JOSIE [145]
		D^3L [13]
		Juneau [70], [142], [143]
		PEXESO [37]
		RNLIM [116]
		DLN [11]
	Data integration	Constance [58], [59], [60], [62]
	Metadata enrichment	CoreDB [8], [9]
		D^4 [104]
		DomainNet [79]
		Constance [61]
		GOODS [64], [65]
	Data cleaning	CLAMS [44]
		Constance [61]
		Song et al. [130]
	Schema evolution	Klettke et al. [75]
	Data provenance	IBM tool [134]
		Suriarachchi et al. [132]
GOODS [64], [65]		
CoreDB [8], [9]		
Juneau [70], [142], [143]		
Exploration	Query-driven data discovery	JOSIE [145]
		D^3L [13]
		Juneau [70], [142], [143]
		Aurum [45]
	Heterogeneous data querying	Constance [58], [62]
		CoreDB [8], [9]
		Ontario [41], [74]
		Squerall [89]

[15]

Metadata Extraction

- Structures: DATAMARAN [Gao et al., 2018]
- Content and context: Skluma [Skluzacek et al., 2018]
- Schema and metadata properties: GEMMS [Quix et al., 2016b]**

Metadata Modeling

- Generic metadata model [Quix et al., 2016b, Hai et al., 2019]**
- Data Vault [Nogueira et al., 2018; Giebler et al., 2019]
- Network-based metadata model [Diamantini et al., 2018a,b]
- Enterprise knowledge graph [Fernandez et al., 2018]

Data Integration

Mapping

Query re

Schema

Dataset Organi

Discover Relat

Query-driven Discovery

- Explore related datasets with keyword queries/primitive-based query language [Nargesian et al., 2018; Brackenbury et al., 2018; Fernandez et al., 2018; Zhu et al., 2019]

Query Heterogeneous Data with a Unified Interface

- Support multiple query languages** [Beheshti et al., 2017;2018; Hai et al., 2016,2018b]

richment

and named entities [Beheshti et al., 2017; 2018]
ricing for descriptive metadata [Halevy et al., 2016a;b]

functional dependencies [Hai et al., 2019b]

lity

ata objects violating functional dependencies [Hai et al.,

DF triples violating conditional denial constraints [Farid et al., 2016]

Our contribution in data lake solution landscape

■ Metadata Extraction

- Structures: DATAMARAN [Gao et al., 2018]
- Content and context: Skluma [Skluzacek et al., 2018]
- Schema and metadata properties: GEMMS [Quix et al., 2016b]**

■ Metadata Modeling

- Generic metadata model [Quix et al., 2016b, Hai et al., 2019]**
- Data Vault [Nogueira et al., 2018; Giebler et al., 2019]
- Network-based metadata model [Diamantini et al., 2018a,b]
- Enterprise knowledge graph [Fernandez et al., 2018]

■ Hadoop and File-based Storage Systems

- Hadoop Distributed File System [Stein et al., 2014; Boci et al., 2015]
- Azure data lake store [Ramakrishnan et al., 2017]

■ Single Data Store

- RDBMS [Zhu et al., 2019]
- NoSQL store [Walker et al., 2015]

■ Polystore Systems

- Google Dataset Search (GOODS) [Halevy et al., 2016b]
- CoreDB [Beheshti et al., 2017]
- Constance [Hai et al., 2016]**

■ Query-driven Discovery

- Explore related datasets with keyword queries/primitive-based query language [Nargesian et al., 2018; Brackenbury et al., 2018; Fernandez et al., 2018; Zhu et al., 2019]

■ Query Heterogeneous Data with a Unified Interface

- Support multiple query languages** [Beheshti et al., 2017;2018; Hai et al., 2016,2018b]

■ Data Integration

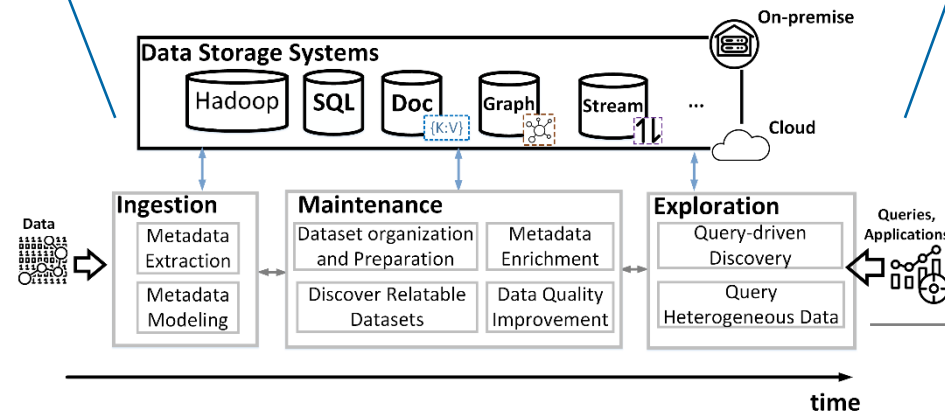
- Mapping generation: [Hai et al., 2018b, 2019a]**
- Query rewriting in Polystore based data lakes [Hai et al., 2018a]**
- Schema matching [Alserafi et al., 2020]

■ Dataset Organization and Preparation

■ Discover Relatable Datasets

■ Metadata Enrichment

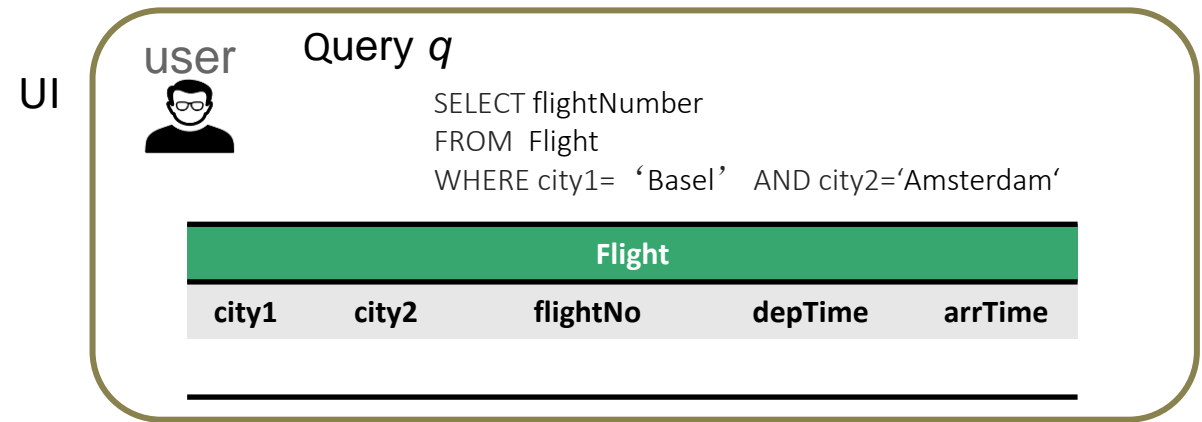
- Keywords and named entities [Beheshti et al., 2017; 2018]
- Crowd-sourcing for descriptive metadata [Halevy et al., 2016a;b]
- Relaxed functional dependencies [Hai et al., 2019b]**
- Data Quality**
- Detect data objects violating functional dependencies [Hai et al., 2019b]**
- Examine RDF triples violating conditional denial constraints [Farid et al., 2016]



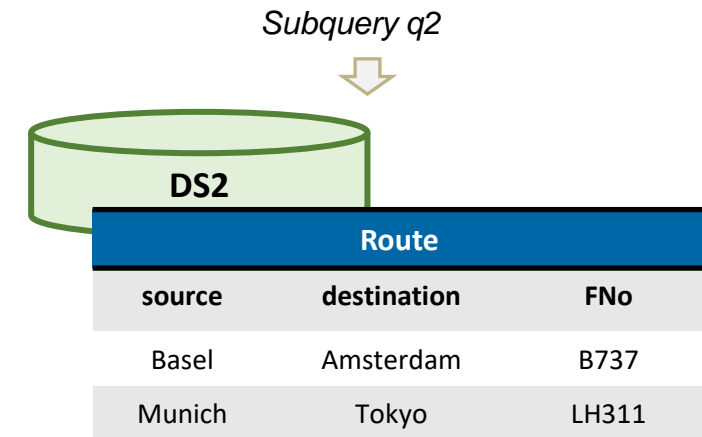
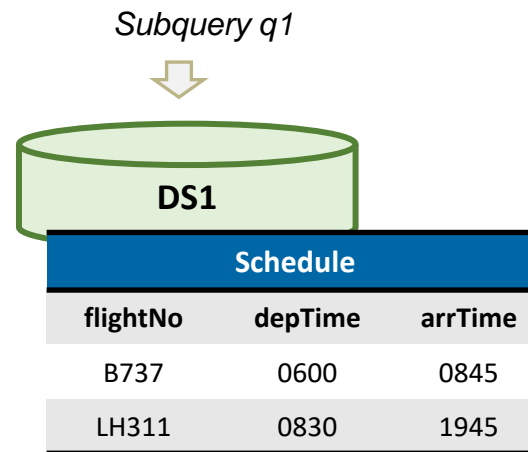
When is data integration?

Goal:

- Combine multiple heterogeneous data sources
- Provide a uniform query interface for heterogeneous data



How to rewrite q to $q1$ and $q2$?



Motivation | Data Integration and Metadata Management

- **Schema Matching**

Discover the correspondences among source schemas

- **Schema Merging**

Resolve several related source schemas, and build the integrated schema

- **Schema Mapping**

Capture semantic relationships between the source schemas and the integrated/target schema

- **Entity linkage**

Discover records across different data sources, which represent the same entity

- **Data cleaning**

Detect and correct corrupt or inaccurate records

- **Query Reformulation**

Translate user queries posed on the global/target schema to a set of subqueries based on source schemas

Schedule		
flightNo	depTime	arrTime
B737	0600	0845
LH311	0830	1945

Route		
source	destination	FNo
Basel	Amsterdam	B737
Munich	Tokyo	LH311

Data Integration key operations

- **Schema Matching**

Discover the correspondences among source schemas

- **Schema Merging**

Resolve several related source schemas, and build the integrated schema

- **Schema Mapping**

Capture semantic relationships between the source schemas and the integrated/target schema

- **Entity linkage**

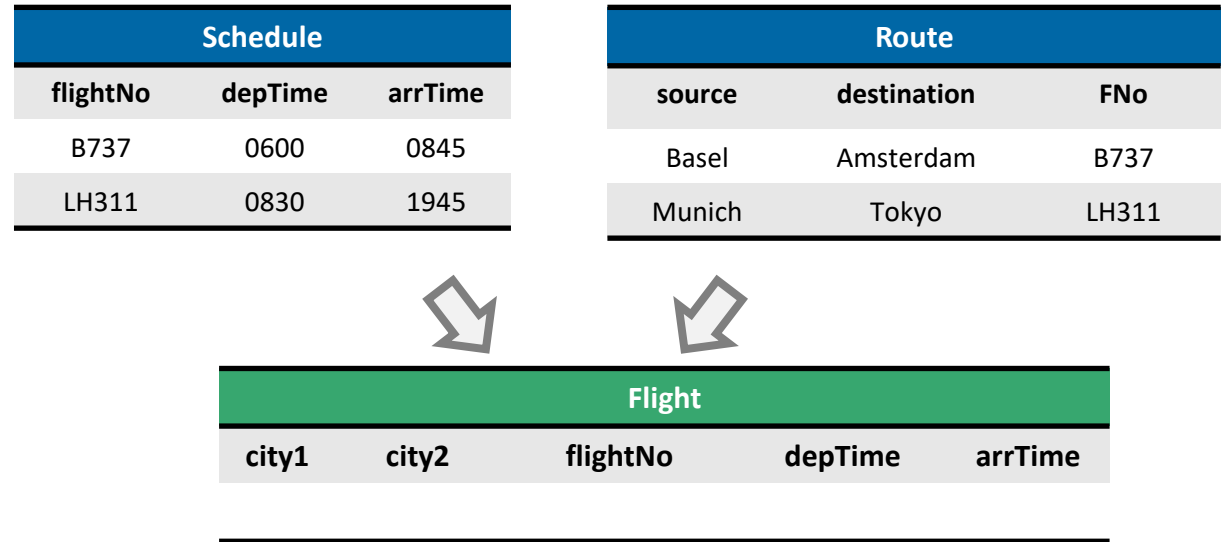
Discover records across different data sources, which represent the same entity

- **Data cleaning**

Detect and correct corrupt or inaccurate records

- **Query Reformulation**

Translate user queries posed on the global/target schema to a set of subqueries based on source schemas



Data Integration key operations

- **Schema Matching**

Discover the correspondences among source schemas

- **Schema Merging**

Resolve several related source schemas, and build the integrated schema

- **Schema Mapping**

Capture semantic relationships between the source schemas and the integrated/target schema

- **Entity linkage**

Discover records across different data sources, which represent the same entity

- **Data cleaning**

Detect and correct corrupt or inaccurate records

- **Query Reformulation**

Translate user queries posed on the global/target schema to a set of subqueries based on source schemas

Schedule		
flightNo	depTime	arrTime
B737	0600	0845
LH311	0830	1945

Route		
source	destination	FNo
Basel	Amsterdam	B737
Munich	Tokyo	LH311



Flight				
city1	city2	flightNo	depTime	arrTime

- Transformation: executable mapping constraint
 - Constructs target instances from source instances
 - E.g., SQL query, XSLT, C# program
- Mapping formalism: Tuple-generating dependency (TGD)

$$\forall c_1, c_2, d \ (Route(c_1, c_2, d) \rightarrow \exists a_1, a_2 \ (Flight(c_1, c_2, d, a_1, a_2)))$$

Data Integration key operations

- **Schema Matching**

Discover the correspondences among source schemas

- **Schema Merging**

Resolve several related source schemas, and build the integrated schema

- **Schema Mapping**

Capture semantic relationships between the source schemas and the integrated/target schema

- **Entity linkage**

Discover records across different data sources, which represent the same entity

- **Data cleaning**

Detect and correct corrupt or inaccurate records

- **Query Reformulation**

Translate user queries posed on the global/target schema to a set of subqueries based on source schemas

Schedule				Route		
flightNo	depTime	arrTime		source	destination	FNo
B737	0600	0845	↔	Basel	Amsterdam	B737
LH311	0830	1945		Munich	Tokyo	LH311

Data Integration key operations

- **Schema Matching**

Discover the correspondences among source schemas

- **Schema Merging**

Resolve several related source schemas, and build the integrated schema

- **Schema Mapping**

Capture semantic relationships between the source schemas and the integrated/target schema

- **Entity linkage**

Discover records across different data sources, which represent the same entity

- **Data cleaning**

Detect and correct corrupt or inaccurate records

- **Query Reformulation**

Translate user queries posed on the global/target schema to a set of subqueries based on source schemas

Schedule		
flightNo	depTime	arrTime
B737	0600	0845
LH311	0830	1945



Route		
source	destination	FNo
Basel	Amsterdam	B737?
Munich	Tokyo	LH311

Data Integration key operations

- **Schema Matching**

Discover the correspondences among source schemas

- **Schema Merging**

Resolve several related source schemas, and build the integrated schema

- **Schema Mapping**

Capture semantic relationships between the source schemas and the integrated/target schema

- **Entity linkage**

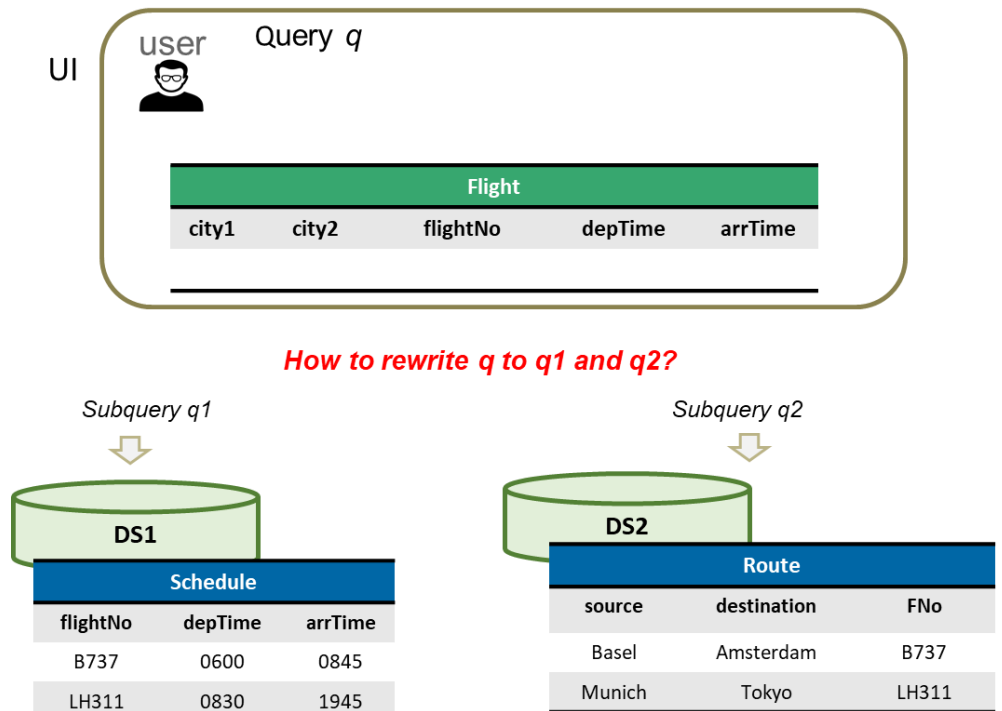
Discover records across different data sources, which represent the same entity

- **Data cleaning**

Detect and correct corrupt or inaccurate records

- **Query Reformulation**

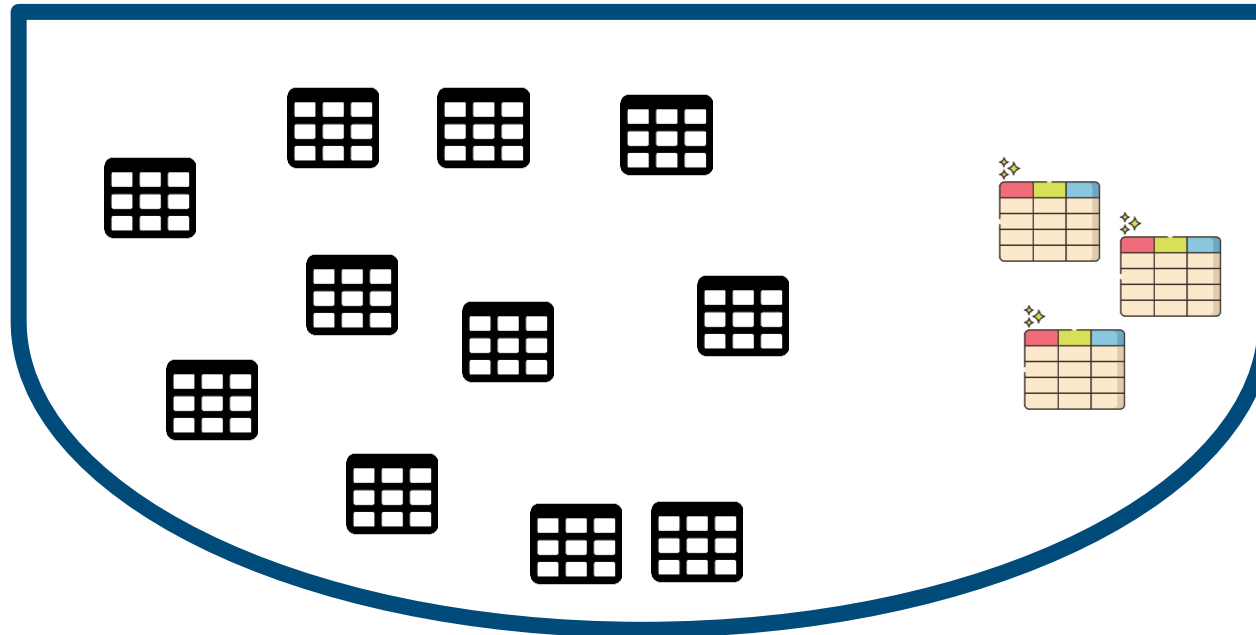
Translate user queries posed on the global/target schema to a set of subqueries based on source schemas



Challenges: Data integration in data lakes

- Open data: massive datasets, not all needed
- Various data: heterogeneous data models (e.g., relational, semi-structured)
- Insufficient metadata: structureless datasets, missing semantics

Data Lake



Challenges: Data integration in data lakes

- Open data: massive datasets, not all needed
- Various data: heterogeneous data models (e.g., relational, semi-structured)
- Insufficient metadata: structureless datasets, missing semantics

Table 3: Comparison of related dataset discovery approaches in data lakes

Systems	Relatedness criteria	Similarity metrics	Applied technique
<i>Aurum</i> [45]	Instance value overlap Attribute name PK-FK candidate	Jaccard similarity (MinHash) Cosine similarity (TF-IDF)	Hypergraph
<i>Brackenbury et.al.</i> [14]	Instance value overlap Attribute name Semantics Descriptive metadata	Jaccard similarity (MinHash)	-
<i>JOSIE</i> [145]	Instance value overlap	Intersection size of sets	Inverted Index
<i>D³L</i> [13]	Instance value overlap Attribute name Semantics Data value representation pattern (Numerical) data distribution	Jaccard similarity (MinHash) Cosine similarity (Random projections)	5-dim Euclidean space
<i>Juneau</i> [70], [142], [143]	Instance value overlap Domain overlap Attribute name Key constraint New attributes rate New instance rate Variable dependency Descriptive metadata Null Values	Jaccard similarity	Workflow graph Variable dependency graph
<i>PEXESO</i> [37]	(Textual) instance values	Any similarity function in a metric space	High-dimensional vectors Hierarchical grids Inverted Index
<i>RNLIM</i> [116]	Table name Attribute name Attribute data type Attribute value domain	-	BERT [30]
<i>DLN</i> [11]	Attribute name Instance values	Jaccard similarity Cosine similarity	Classification models

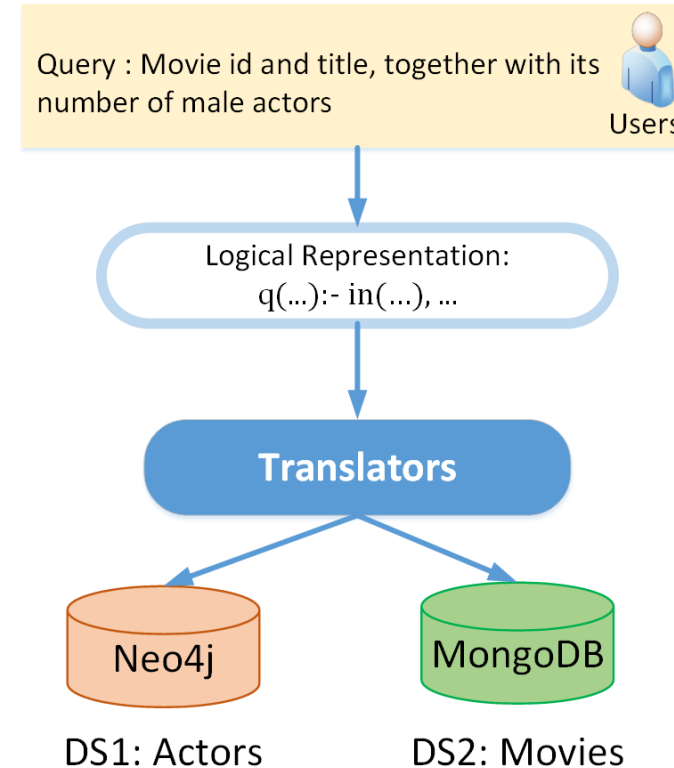
Challenges: Data integration in data lakes

- Open data: massive datasets, not all needed
- Various data: heterogeneous data models (e.g., relational, semi-structured)
- Insufficient metadata: structureless datasets, missing semantics

Query Rewriting in Data Lakes [Hai et al., ADBIS 2018]

Query rewriting for heterogeneous data

- User query might need data from multiple sources
 - Different source schemas and formats
 - Diverse query languages



Query Processing over Diverse Polystore [Hai et al., ADBIS 2018]

■ Reformulate and process subqueries

- ❶ Rewrite an input query on integrated schema to logical rules in predicates
- ❷❸ Generate subqueries based on source schemas and stores
- ❹ Execute subqueries on each data store
- ❺❻ Retrieve the query results, and create the final integrated results

■ Optimize query execution

- ❸ Push down selection predicates to the data sources
- ❺ Reduce the amount of data that has to be returned and loaded into Spark

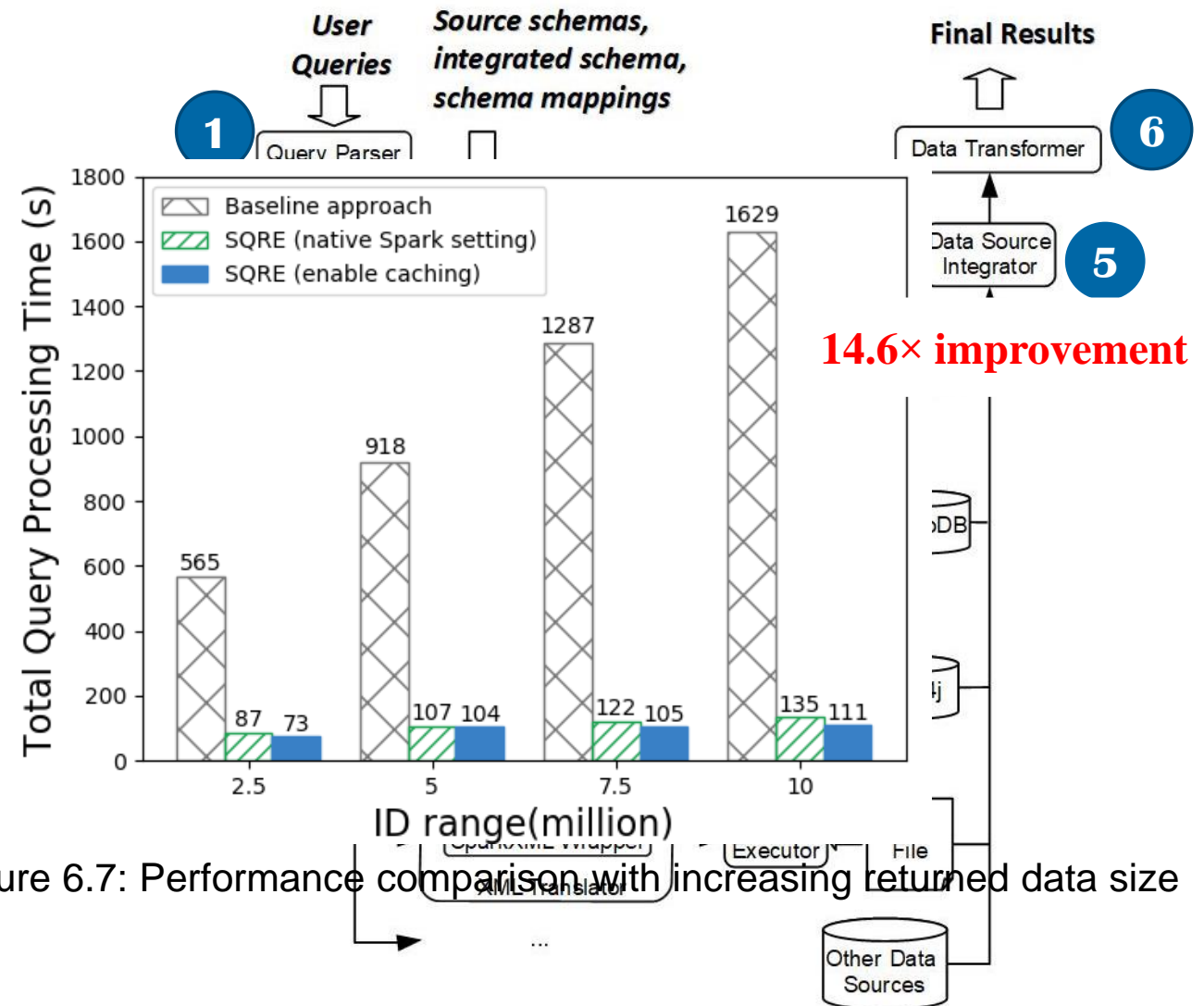


Figure 6.7: Performance comparison with increasing returned data size

Challenges: Data integration in data lakes

- Open data: massive datasets, not all needed
- Various data: heterogeneous data models (e.g., relational, semi-structured)
- Insufficient metadata: structureless datasets, missing semantics

Table 2: Comparison of DAG-based dataset organization approaches in Sec. 6.1.3

System	<i>KAYAK [85], [86]</i> (<i>pipeline</i>)	<i>KAYAK [85], [86]</i> (<i>task dependency</i>)	<i>Nargesian et al. [98]</i>	<i>Juneau [143]</i> (<i>variable dependency</i>)
Function	Represent the primitives of a data preparation pipeline	Enforce correct execution sequence of tasks while parallelization	Semantic navigation	Measure table relatedness w.r.t. notebook workflow
Node	Primitives	Atomic tasks for data preparation operations	Sets of attributes	Notebook variables
Edge	Sequential execution order of two primitives	Sequential execution order of two tasks	Containment relationships	Notebook functions (as edge labels)
Edge direction	From the previous primitive to the subsequent primitive	From the previous task to the subsequent task	From the superset to the subset	From the input variable of the function to the output variable

Future work: data lake for AI

Data Lake (2015-2020)



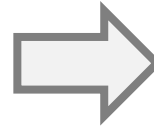
RWTHAACHEN
UNIVERSITY



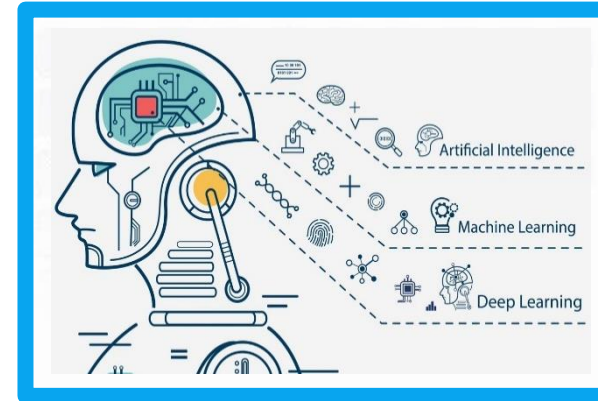
Christoph Quix



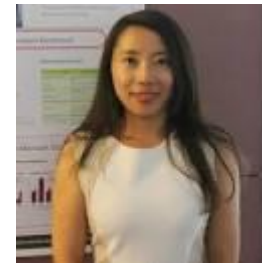
Matthias Jarke



Intelligent, scalable Data Lake (2021-)



TU Delft Delft
University of
Technology



Rihan Hai



Asterios Katsifodimos

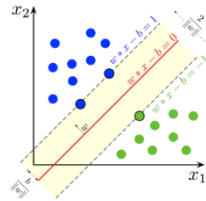
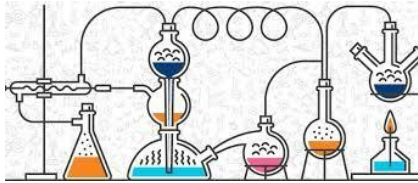
Amalur: a data science platform

Use cases:

Healthcare



Chemistry

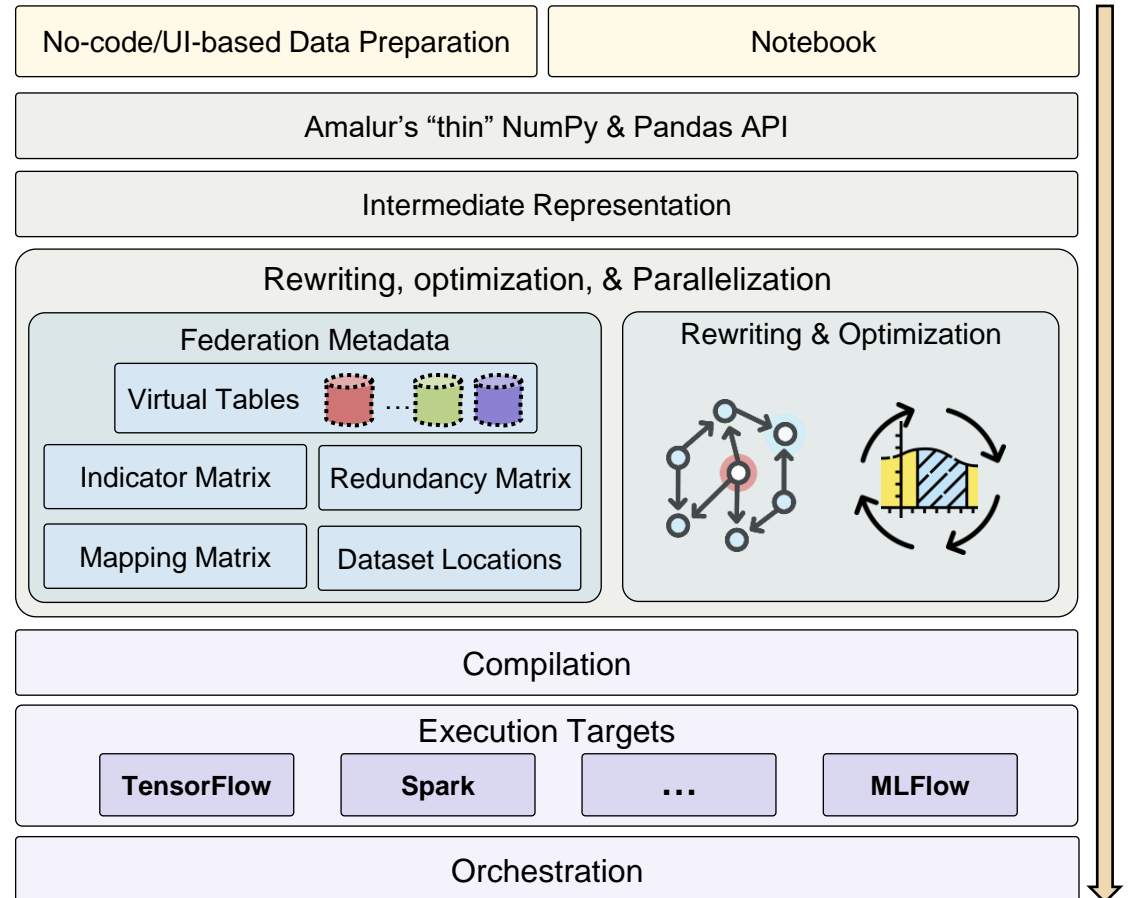


Data Analytics

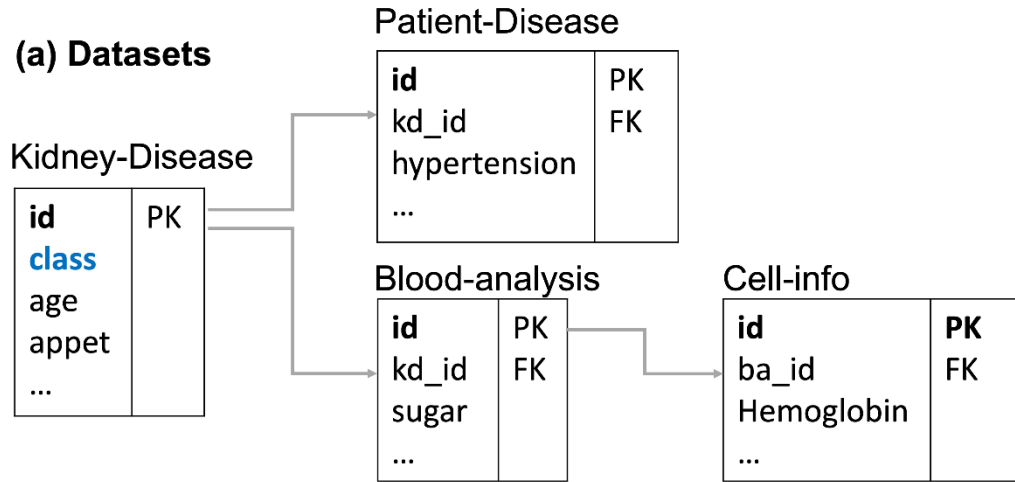


Insight

Amalur



Early results so far



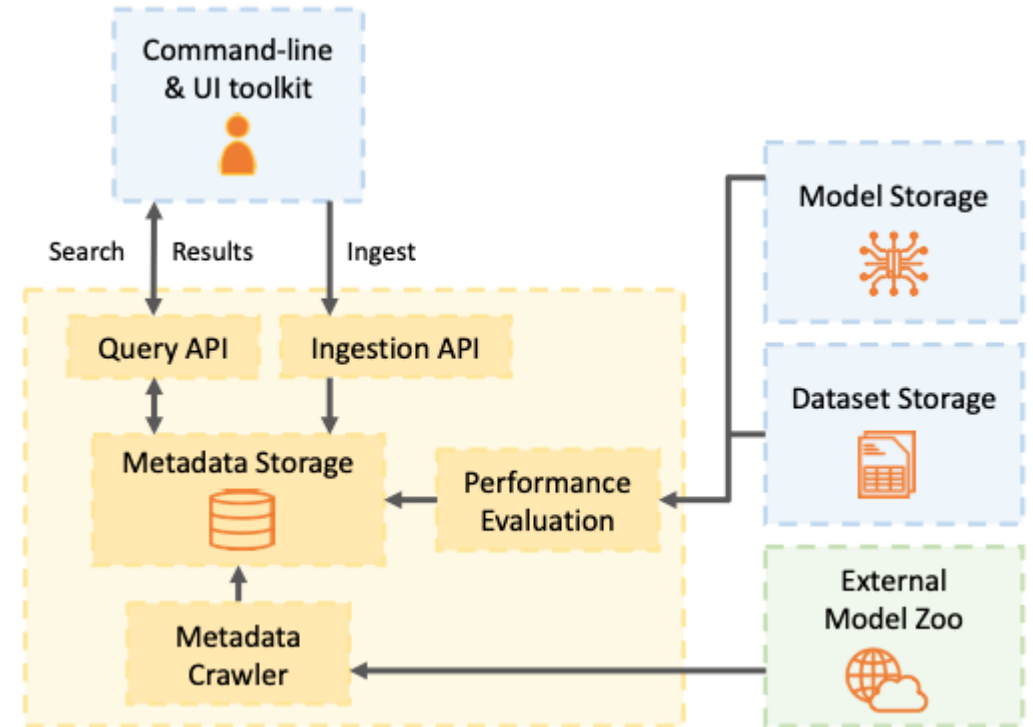
(b) Ranked list of join path

p_1 : Kidney-Disease \rightarrow Blood-analysis \rightarrow Cell-info

p_2 : Kidney-Disease \rightarrow Blood-analysis

p_3 : Kidney-Disease \rightarrow Patient-Disease

(c) DT_base	(d) DT_JoinAll	(e) DT_ p_2	(f) DT_ p_1
acc = 71.25%	acc = 96.25%	acc = 80.25%	acc = 95.75%
Tree_dep= 4	Tree_dep= 10	Tree_dep= 8	Tree_dep= 4



Andra Ionescu, et al. "Join Path Based Data Augmentation for Decision Trees." ICDE Workshops. 2022.

Ziyu Li, et al. "Metadata Representations for Queryable ML Model Zoos." ICML Workshops. 2022.

Conclusion

- Data integration is a difficult problem, even more difficult in data lakes
 - Data structured with various schemas
 - Heterogeneous data models
 - Diverse query languages and systems
- New challenges, new technologies
 - Related dataset discovery
 - Dataset organization
 - Metadata extraction/dataset profiling
 - Entity resolution using ML/DL
 - Data cleaning for ML
 - Factorized Databases
 - Federated learning
 - ...

