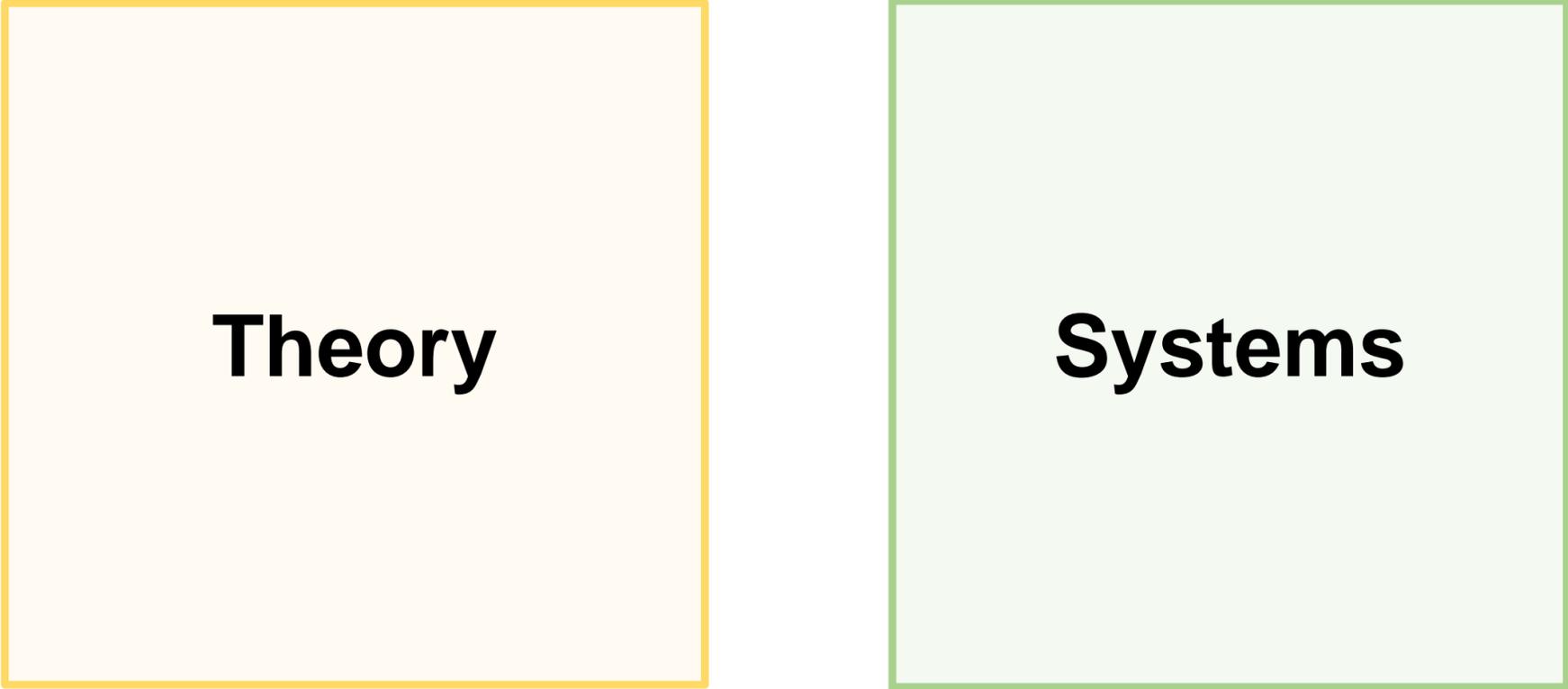# Amalur: Data Integration Meets Machine Learning

**Rihan Hai**, Christos Koutras, Andra Ionescu, Ziyu Li, Wenbo Sun, Jessie van Schijndel, Yan Kang, Asterios Katsifodimos

# Take-away

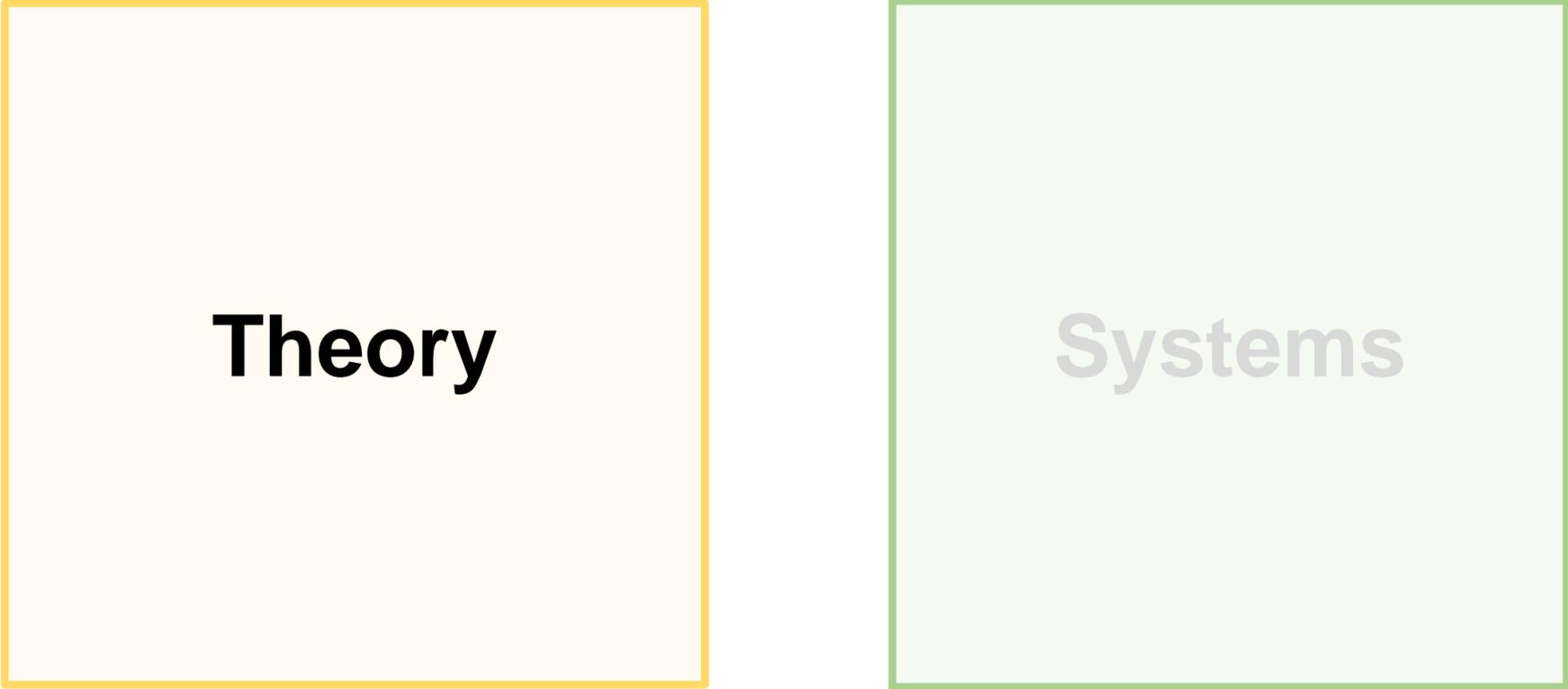**Q: Can we use data integration metadata to improve the effectiveness and efficiency of ML model training?**
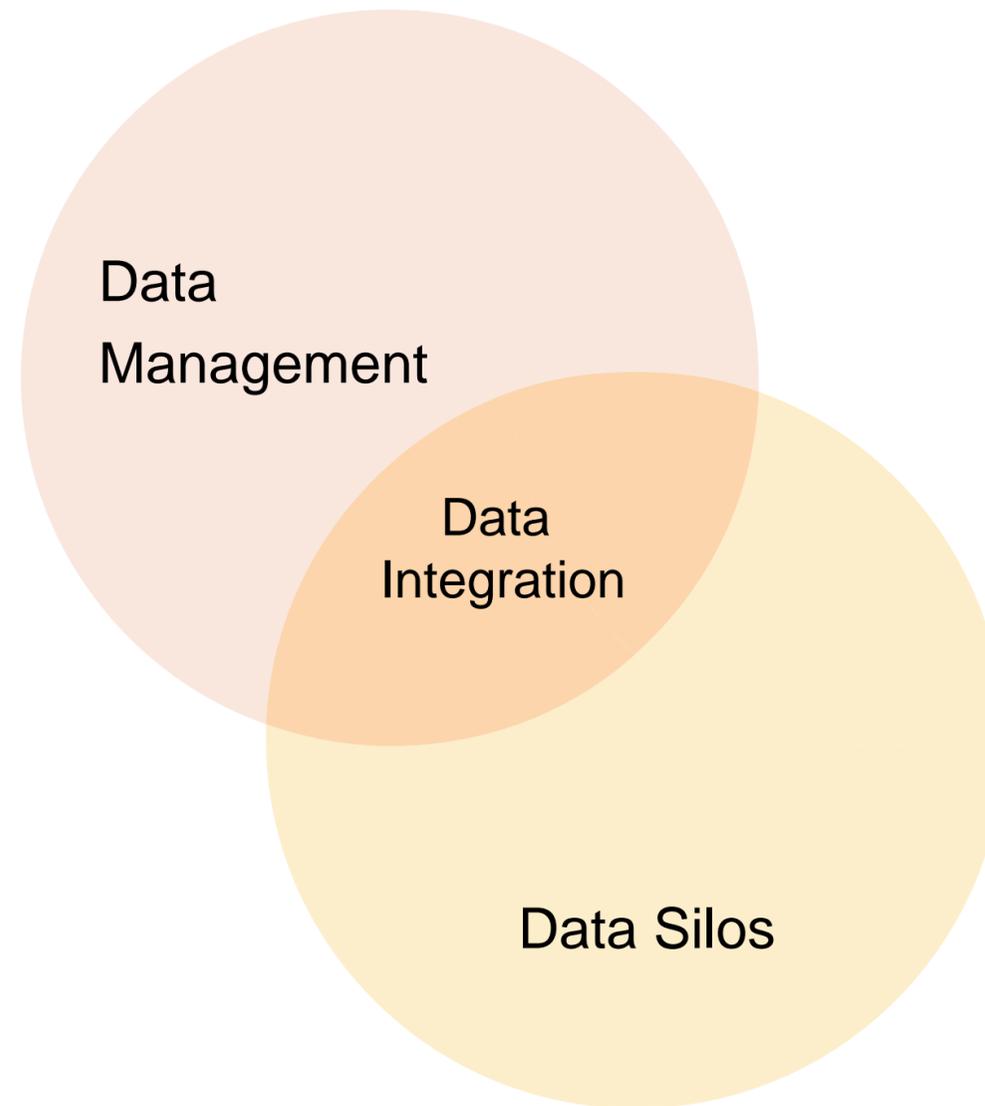
TUDelft

# Scope

Theory

Systems

# Scope

Theory

Systems

TUDelft

# Data integration -- in a nutshell

# Data integration -- System architecture



**Mediated Schema or Warehouse**

*Query reformulation / Query over materialized data*

*Source descriptions/ Transforms*

Wrapper/ Extractor

Wrapper/ Extractor

Wrapper/ Extractor

Wrapper/ Extractor

$RDBMS_1$

$HTML_1$

$RDBMS_2$

$XML_1$

Basic architecture of a general-purpose data integration system
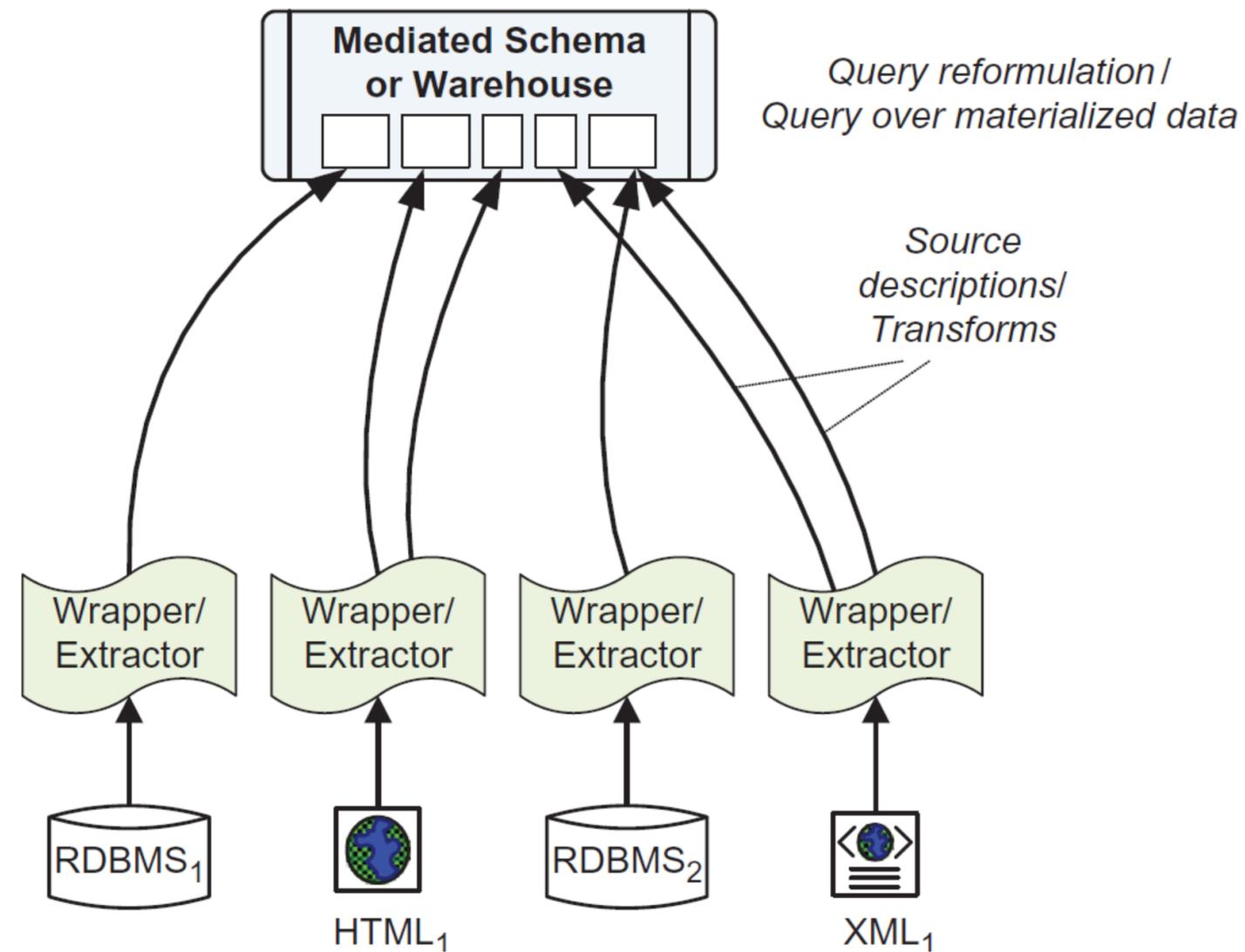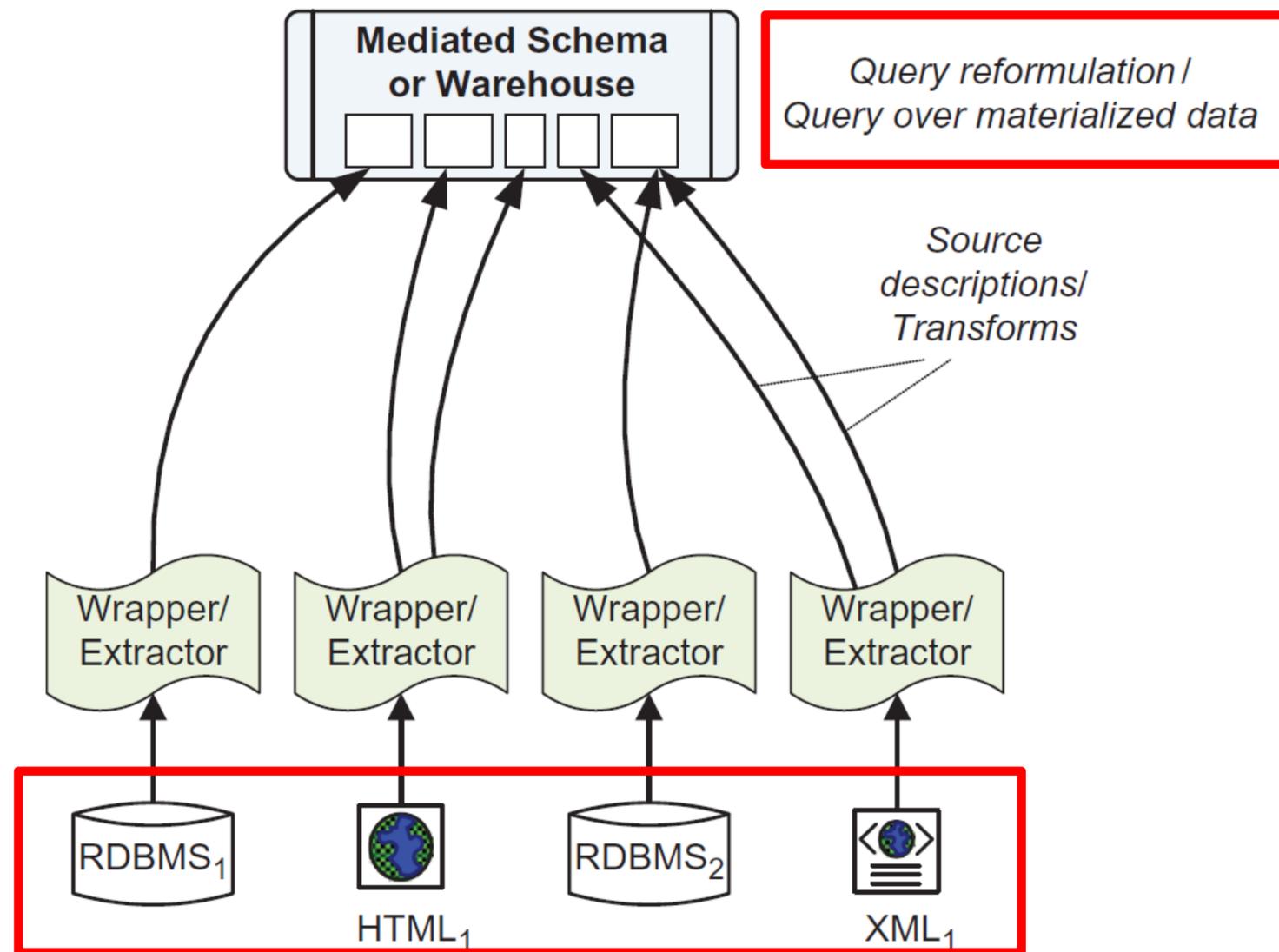
# Data integration -- System architecture



Basic architecture of a general-purpose data integration system
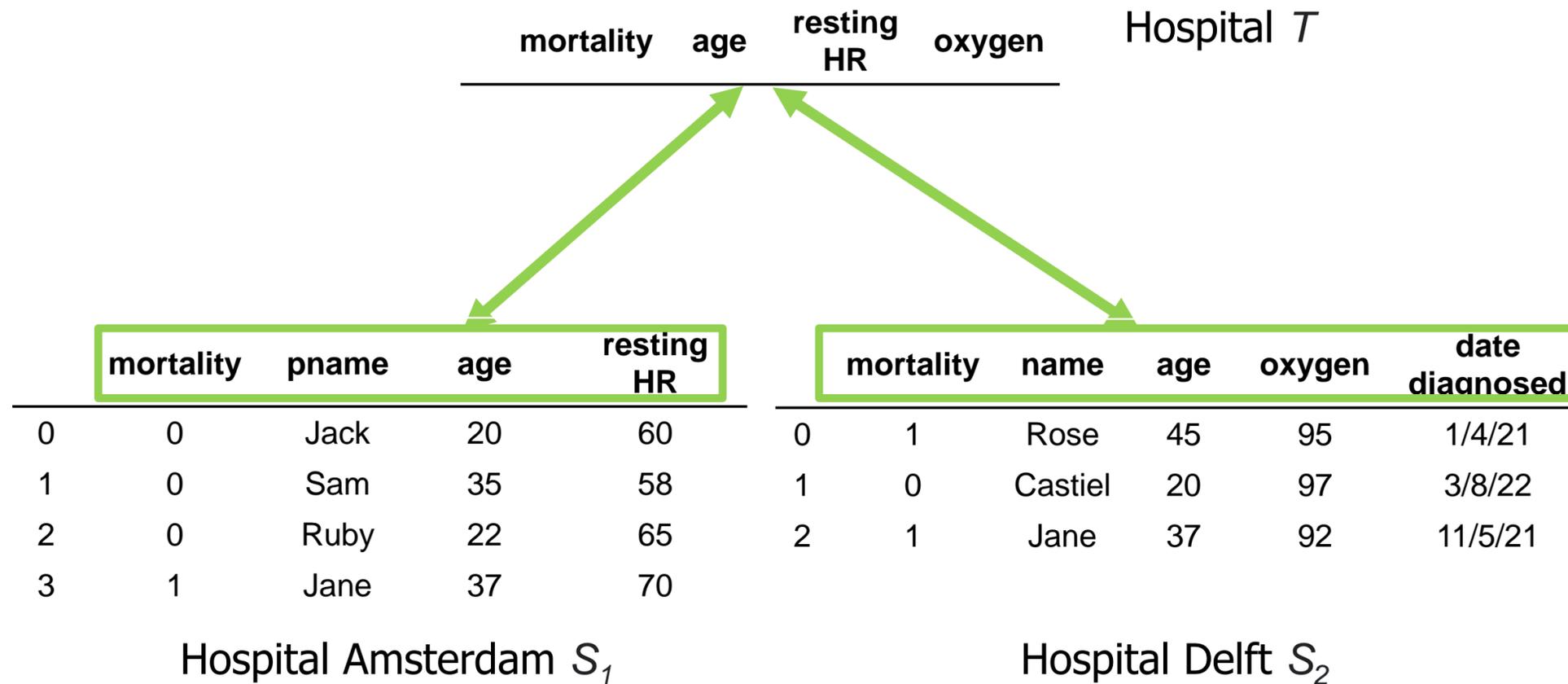
*Principles of Data Integration* Alon Halevy, AnHai Doan, and Zachary

# Data integration -- Schema matching

| | mortality | pname | age | resting HR |
|---|---|---|---|---|
| 0 | 0 | Jack | 20 | 60 |
| 1 | 0 | Sam | 35 | 58 |
| 2 | 0 | Ruby | 22 | 65 |
| 3 | 1 | Jane | 37 | 70 |

Hospital Amsterdam $S_1$

| | mortality | name | age | oxygen | date diagnosed |
|---|---|---|---|---|---|
| 0 | 1 | Rose | 45 | 95 | 1/4/21 |
| 1 | 0 | Castiel | 20 | 97 | 3/8/22 |
| 2 | 1 | Jane | 37 | 92 | 11/5/21 |

Hospital Delft $S_2$

# Data integration -- Schema merging

Target table schema
(a.k.a. global schema, mediated schema)

| mortality | age | resting HR | oxygen |
|-----------|-----|------------|--------|

Hospital $T$

| | mortality | pname | age | resting HR |
|---|-----------|-------|-----|------------|
| 0 | 0 | Jack | 20 | 60 |
| 1 | 0 | Sam | 35 | 58 |
| 2 | 0 | Ruby | 22 | 65 |
| 3 | 1 | Jane | 37 | 70 |

Hospital Amsterdam $S_1$

| | mortality | name | age | oxygen | date diagnosed |
|---|-----------|------|-----|--------|----------------|
| 0 | 1 | Rose | 45 | 95 | 1/4/21 |
| 1 | 0 | Castiel | 20 | 97 | 3/8/22 |
| 2 | 1 | Jane | 37 | 92 | 11/5/21 |

Hospital Delft $S_2$

# Data integration -- Schema mapping

Target table schema

| mortality | age | resting HR | oxygen |
|-----------|-----|------------|--------|

Hospital $T$



**Hospital Amsterdam $S_1$**

| | mortality | pname | age | resting HR |
|---|-----------|-------|-----|------------|
| 0 | 0 | Jack | 20 | 60 |
| 1 | 0 | Sam | 35 | 58 |
| 2 | 0 | Ruby | 22 | 65 |
| 3 | 1 | Jane | 37 | 70 |

**Hospital Delft $S_2$**

| | mortality | name | age | oxygen | date diagnosed |
|---|-----------|------|-----|--------|----------------|
| 0 | 1 | Rose | 45 | 95 | 1/4/21 |
| 1 | 0 | Castiel | 20 | 97 | 3/8/22 |
| 2 | 1 | Jane | 37 | 92 | 11/5/21 |

# Data integration -- Entity resolution

| | mortality | pname | age | resting HR |
|---|---|---|---|---|
| 0 | 0 | Jack | 20 | 60 |
| 1 | 0 | Sam | 35 | 58 |
| 2 | 0 | Ruby | 22 | 65 |
| 3 | 1 | Jane | 37 | 70 |

| | mortality | name | age | oxygen | date diagnosed |
|---|---|---|---|---|---|
| 0 | 1 | Rose | 45 | 95 | 1/4/21 |
| 1 | 0 | Castiel | 20 | 97 | 3/8/22 |
| 2 | 1 | Jane | 37 | 92 | 11/5/21 |

Hospital Amsterdam $S_1$              Hospital Delft $S_2$

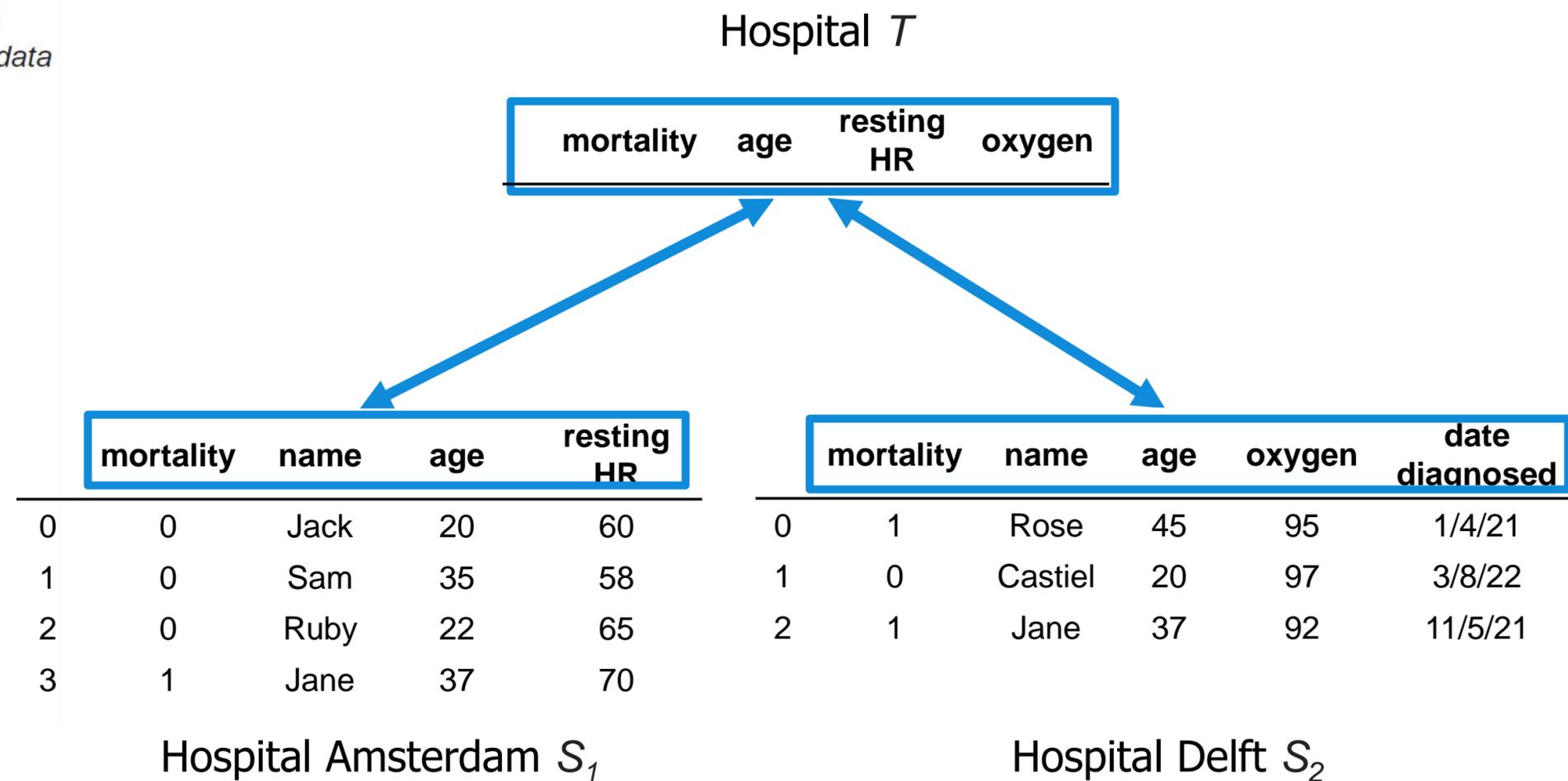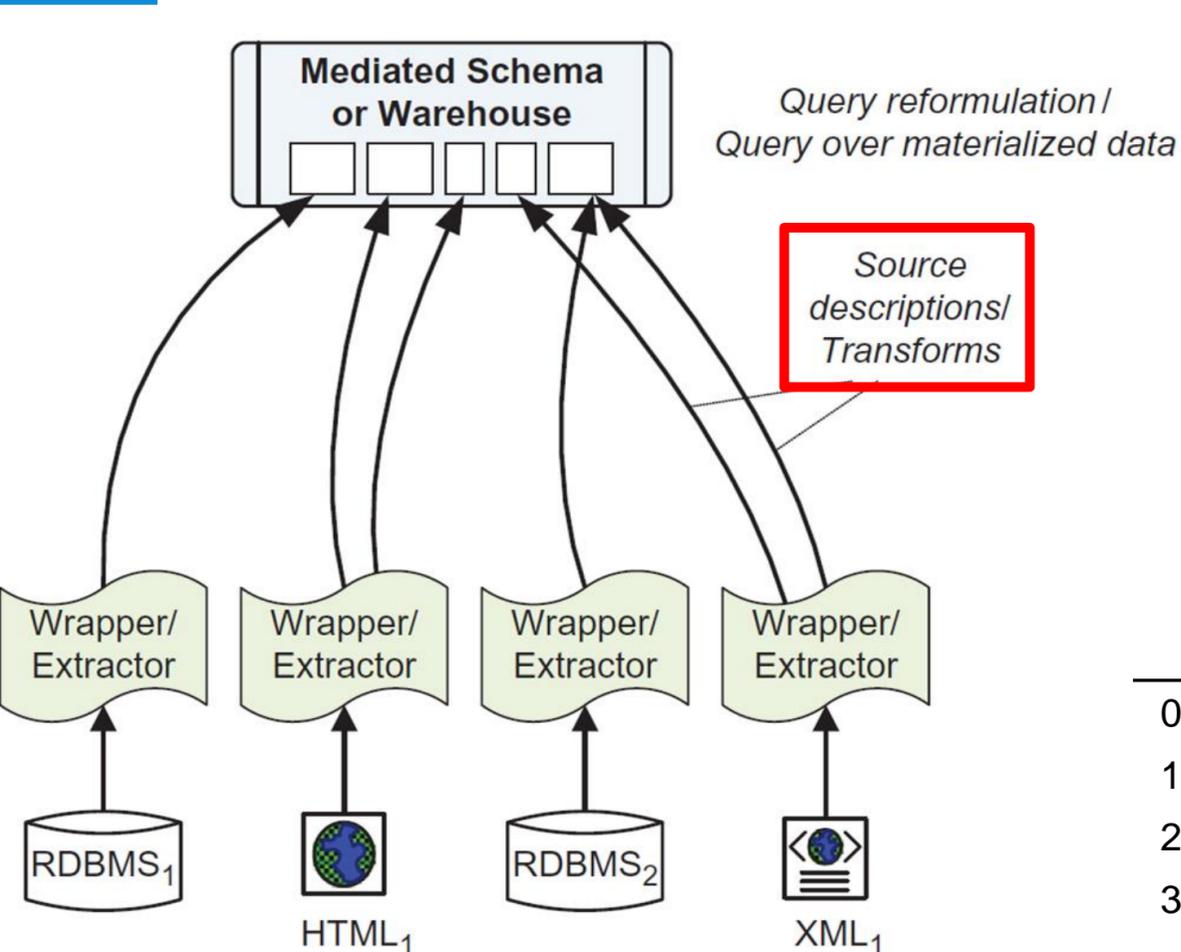# Data integration -- Query reformulation



Basic architecture of a general-purpose data integration system
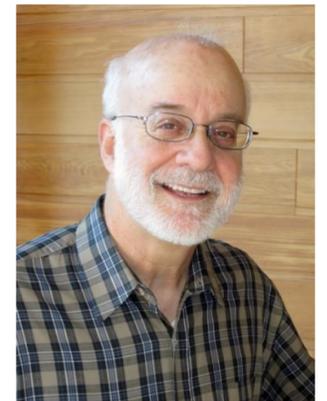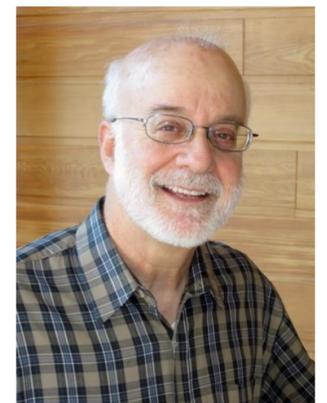
# Data integration -- Query reformulation



Basic architecture of a general-purpose data integration system

**Query q**

```
SELECT count(*)
FROM T
WHERE oxygen<94
```

**Hospital $T$**

| mortality | age | resting HR | oxygen |
|-----------|-----|------------|--------|

| | mortality | name | age | resting HR |
|---|-----------|------|-----|------------|
| 0 | 0 | Jack | 20 | 60 |
| 1 | 0 | Sam | 35 | 58 |
| 2 | 0 | Ruby | 22 | 65 |
| 3 | 1 | Jane | 37 | 70 |

Hospital Amsterdam $S_1$

| | mortality | name | age | oxygen | date diagnosed |
|---|-----------|------|-----|--------|----------------|
| 0 | 1 | Rose | 45 | 95 | 1/4/21 |
| 1 | 0 | Castiel | 20 | 97 | 3/8/22 |
| 2 | 1 | Jane | 37 | 92 | 11/5/21 |

Hospital Delft $S_2$

# Data integration -- Query reformulation



**Mediated Schema or Warehouse**

*Query reformulation / Query over materialized data*

*Source descriptions / Transforms*

Wrapper/Extractor   Wrapper/Extractor   Wrapper/Extractor   Wrapper/Extractor

$RDBMS_1$   $HTML_1$   $RDBMS_2$   $XML_1$

**Query q**

SELECT count(*)
FROM T
WHERE oxygen<94

**Hospital** $T$

| mortality | age | resting HR | oxygen |
|-----------|-----|------------|--------|

**New query $q_1$**

| | mortality | name | age | resting HR |
|---|-----------|------|-----|------------|
| 0 | 0 | Jack | 20 | 60 |
| 1 | 0 | Sam | 35 | 58 |
| 2 | 0 | Ruby | 22 | 65 |
| 3 | 1 | Jane | 37 | 70 |

**Hospital Amsterdam** $S_1$

**New query $q_2$**

| | mortality | name | age | oxygen | date diagnosed |
|---|-----------|------|-----|--------|----------------|
| 0 | 1 | Rose | 45 | 95 | 1/4/21 |
| 1 | 0 | Castiel | 20 | 97 | 3/8/22 |
| 2 | 1 | Jane | 37 | 92 | 11/5/21 |

**Hospital Delft** $S_2$

Basic architecture of a general-purpose data integration system

# Data integration -- Schema mapping



Basic architecture of a general-purpose data integration system

Hospital *T*

| mortality | age | resting HR | oxygen |
|-----------|-----|------------|--------|

Hospital Amsterdam $S_1$

| | mortality | name | age | resting HR |
|---|-----------|------|-----|------------|
| 0 | 0 | Jack | 20 | 60 |
| 1 | 0 | Sam | 35 | 58 |
| 2 | 0 | Ruby | 22 | 65 |
| 3 | 1 | Jane | 37 | 70 |

Hospital Delft $S_2$

| | mortality | name | age | oxygen | date diagnosed |
|---|-----------|------|-----|--------|----------------|
| 0 | 1 | Rose | 45 | 95 | 1/4/21 |
| 1 | 0 | Castiel | 20 | 97 | 3/8/22 |
| 2 | 1 | Jane | 37 | 92 | 11/5/21 |

*Principles of Data Integration* Alon Halevy, AnHai Doan, and Zachary

# Mapping Languages – First-order Logic


Ronald Fagin
IBM Fellow

# Mapping Languages – First-order Logic


Ronald Fagin
IBM Fellow

" ***Tgds*** are one of the two major types
of *database dependencies* "

TUDelft

# Mapping Languages − First-order Logic

**Source-to-target Tuple generating dependencies (s-t tgds) [Beeri and Vardi, 1984]**

**Definition 3.3.1** ([Fagin, 2009]). *Tuple-generating dependencies (tgds) are formulas of* the form

$$\forall x(\varphi(\mathbf{x}) \to \exists \mathbf{y} \ \psi(\mathbf{x}, \mathbf{y})), \text{ where}$$

(1) $\varphi(\mathbf{x})$ is a conjunction of atomic formulas, all with variables among the variables in $\mathbf{x}$.

(2) every variable in x appears in $\varphi(\mathbf{x})$ (but not necessarily in $\psi(\mathbf{x}, \mathbf{y})$).

(3) $\psi(\mathbf{x}, \mathbf{y})$ is a conjunction of atomic formulas, all with variables among the variables in $\mathbf{x}$ and $\mathbf{y}$.

Ronald Fagin
IBM Fellow

" ***Tgds*** are one of the two major types of *database dependencies* "

# The world of schema mapping is far beyond this...

- **More Definitions**

  - Constants & labeled nulls

  - Source and target instances

  - Homomorphism, universal solutions, canonical solution, core, certain answers

  - GAV, LAV and GLAV

  - Closed World Assumption vs. Open World Assumption

  - Term (plain term), function symbol

  - Skolemization, normalization, implication and Logical Equivalence problems

  - ...

- **More People & Topics**

  - **Maurizio Lenzerini** (DI & Ontology related Rewriting)

  - **Alin Deutsch** (Query rewriting, EGDs)

  - **Lucian Popa** (Data Exchange)

  - **Phokion G. Kolaitis** (Nested Mappings)

  - **Balder Ten Cate, Wang-chiew Tan** (Data Examples)

  - **Renée Miller** (Data Exchange, Open Data Integration)

  - …

Schema mappings for nested data?



Expressive power



Table 3.1: Comparison of structural and reasoning properties

| Structural/Reasoning properties | Tgds | Nested tgds | Plain SO tgds | SO tgds |
|---|---|---|---|---|
| Closure under target homomorphisms | Yes | Yes | Yes | No |
| Admitting universal solutions | Yes | Yes | Yes | Yes |
| Allowing for conjunctive query rewriting | Yes | Yes | Yes | Yes |
| Decidable for implication problem | Yes | Yes | Unknown | No |
| Decidable for logical equivalence problem | Yes | Yes | Unknown | No |

# Adding Machine Learning



Data Management

Data Integration

Data Silos
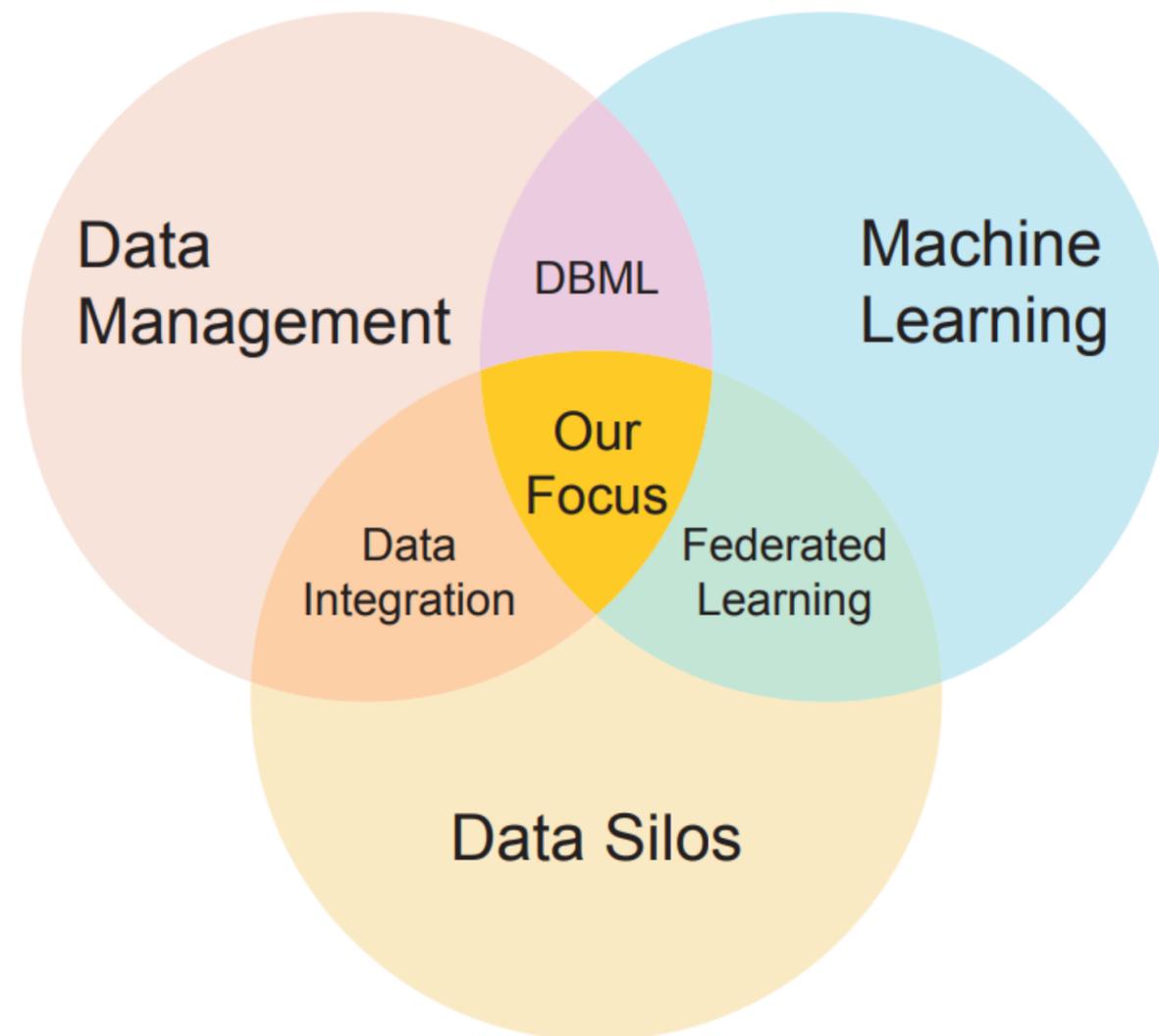
TUDelft

# Adding Machine Learning
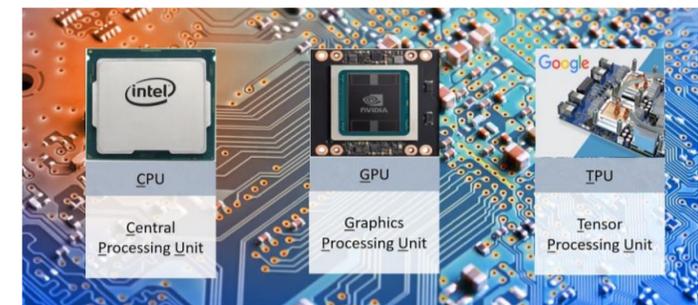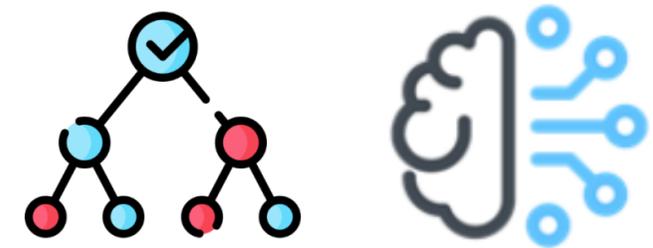
# Adding Machine Learning



ML models

# Adding Machine Learning



ML models



New hardware

# Represent dataset relationships in ML use cases

Step.1: Formal description of ML use cases with
database dependencies

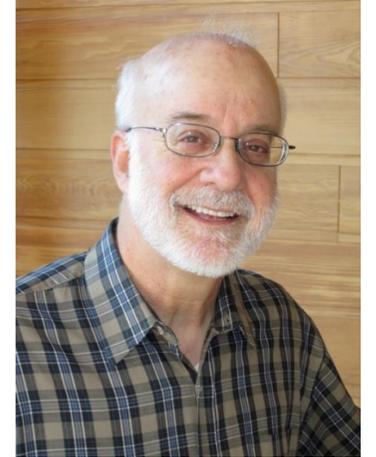*"**Tgds** are one of the two major types*
*of database dependencies"*

**Table 1: Four example data integration scenarios for feature augmentation and federated learning**

| No. | Dataset Relationship | Schema mappings | Example use cases |
|-----|----------------------|-----------------|-------------------|
| 1 | Full outer join | $m_1 : \forall m, n, a, hr, o, dd \ (S_1(m, n, a, hr) \land S_2(m, n, a, o, dd) \to T(m, a, hr, o))$ <br> $m_2 : \forall m, n, a, hr \ (S_1(m, n, a, hr) \to \exists o \ T(m, a, hr, o))$ <br> $m_3 : \forall m, n, a, o, dd \ (S_2(m, n, a, o, dd) \to \exists hr \ T(m, a, hr, o))$ | Feature augmentation, Federated learning, … |
| 2 | Inner join | $m_1 : \forall m, n, a, hr, o, dd \ (S_1(m, n, a, hr) \land S_2(m, n, a, o, dd) \to T(m, a, hr, o))$ | Feature augmentation, (Vertical) federated learning, … |
| 3 | Left join | $m_1 : \forall m, n, a, hr, o, dd \ (S_1(m, n, a, hr) \land S_2(n, a, o, dd) \to T(m, a, hr, o))$ <br> $m_2 : \forall m, n, a, hr \ (S_1(m, n, a, hr) \to \exists o \ T(m, a, hr, o))$ | Feature augmentation, (Vertical) federated learning, … |
| 4 | Union | $m_2 : \forall m, n, a, hr, o \ (S_1(m, n, a, hr, o) \to T(m, a, hr, o))$ <br> $m_3 : \forall m, n, a, hr, o, dd \ (S_2(m, n, a, hr, o, dd) \to T(m, a, hr, o))$ | Data sample augmentation, (Horizontal) federated learning, … |

# Transform dataset relationships in ML use cases

Step.1: Formal description of ML use cases with database dependencies

Step.2: Transform metadata into matrices

## Column matching  (schema mapping)

$$m_1: \forall m, n, a, hr, o, dd \; (S_1(m, n, a, hr) \wedge S_2(m, n, a, o, dd) \longrightarrow T(m, a, hr, o))$$

$$M_1 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \qquad M_2 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

(a) Mapping matrix

## Row matching (entity resolution)

| | | | | | | |
|---|---|---|---|---|---|---|
| $S_1$ | 3 | 1 | Jane | 37 | 70 | |
| $S_2$ | 2 | 1 | Jane | 37 | 92 | 11/5/21 |

$$I_1 \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \qquad I_2 \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

(b) Indicator matrix

TUDelft

# Scope

Theory

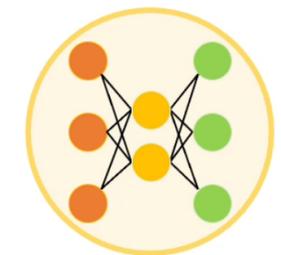**Systems**

# Amalur: Model Lake



Source: https://www.behance.net/gallery/38395965/Mother-Earth

## Web Information Systems (WIS) group



**Data Lake**



arXiv > cs > arXiv:2106.09592

**Computer Science > Databases**

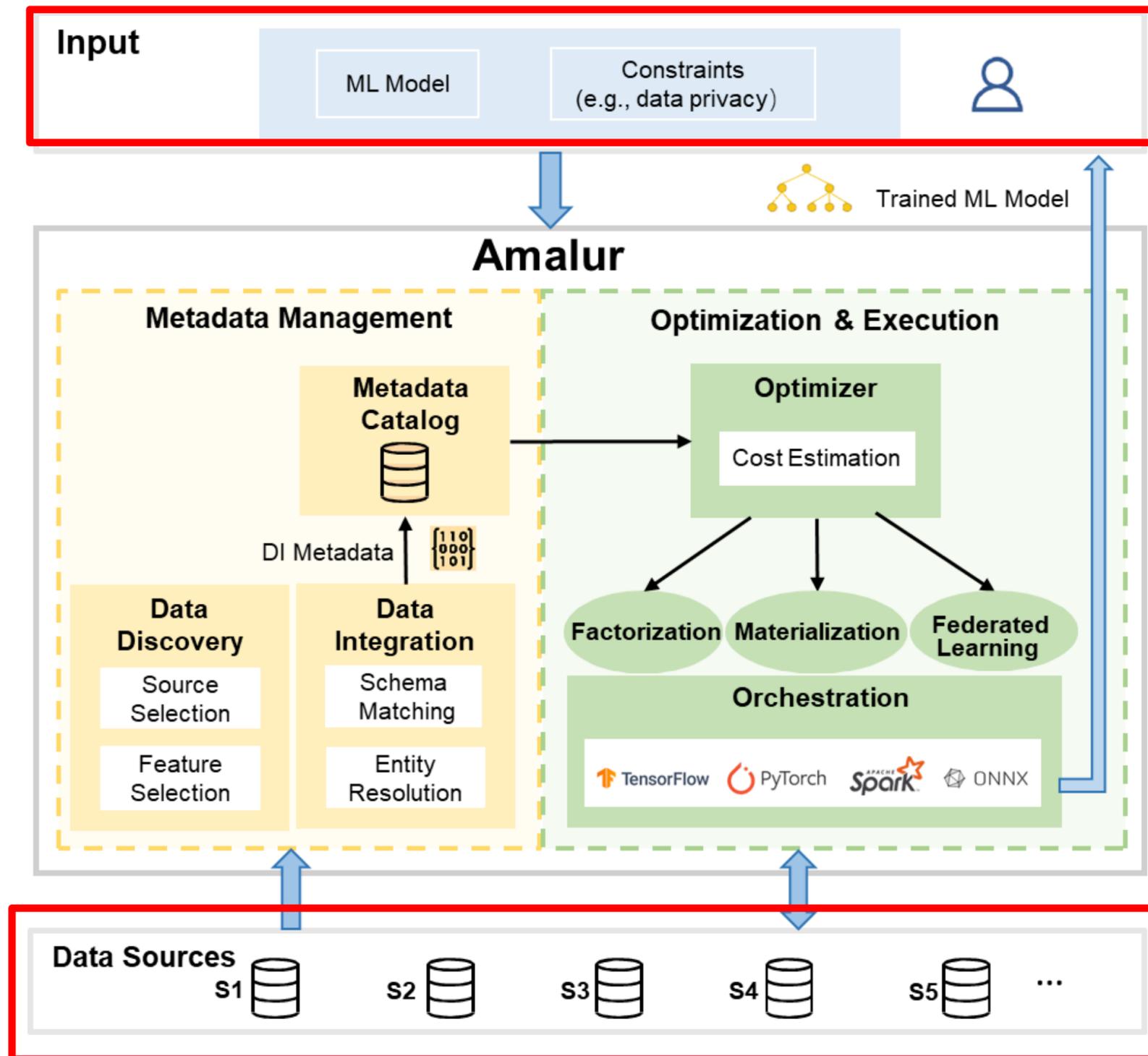[Submitted on 17 Jun 2021 (v1), last revised 17 Feb 2023 (this version, v2)]

**Data Lakes: A Survey of Functions and Systems**

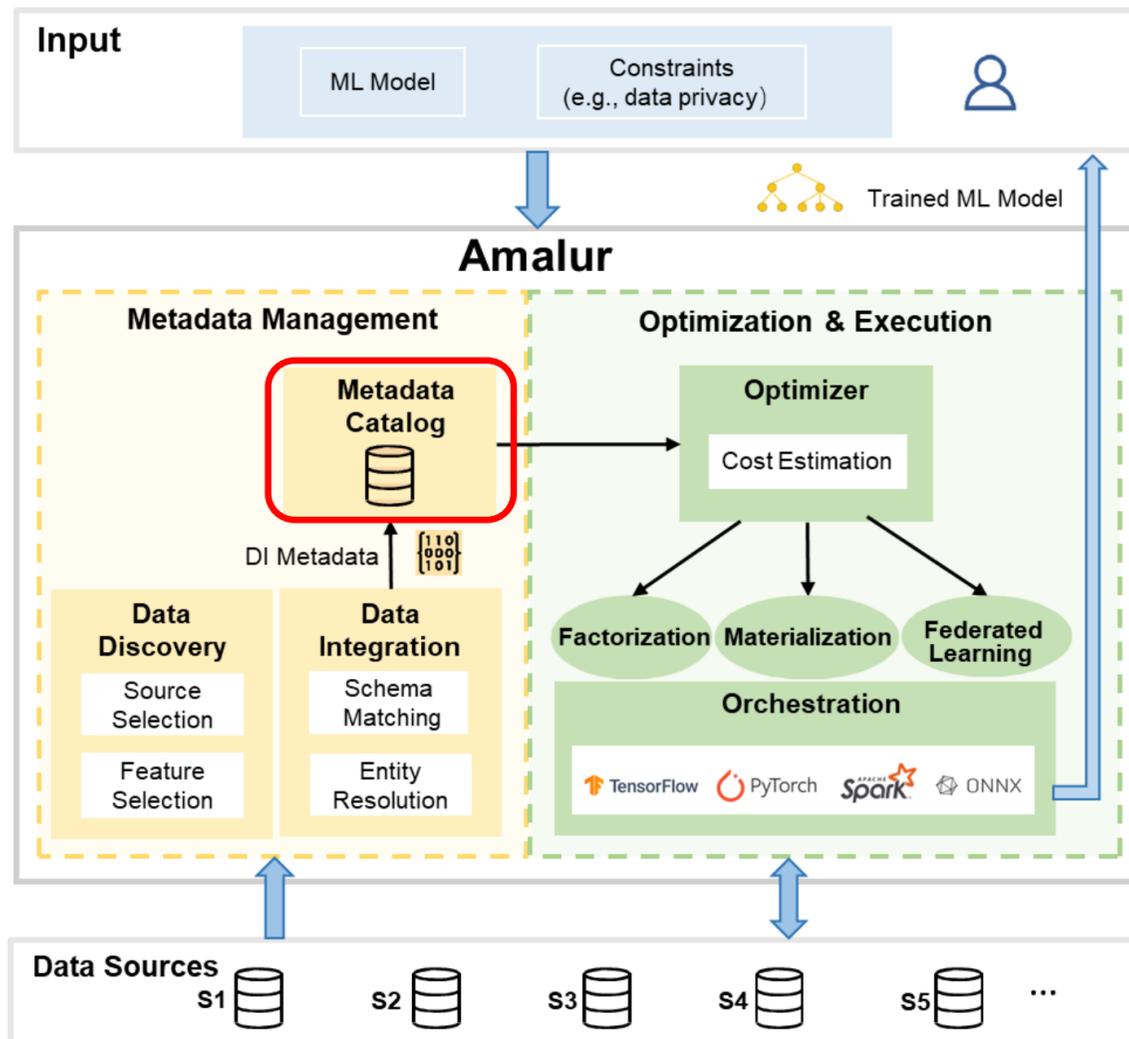Rihan Hai, Christos Koutras, Christoph Quix, Matthias Jarke

**Model Zoo**



**Macaroni: Crawling & Enriching Metadata from Public Model Zoos**

Ziyu Li, Henk Kant, Rihan Hai, Asterios Katsifodimos, and Alessandro Bozzon

Delft University of Technology

TUDelft

# Amalur: Model Lake

# Amalur: Workflow

# Amalur: Workflow



Column matching (schema mapping)

$m_1: \forall m,n,a,hr,o,dd \ (S_1(m,n,a,hr) \wedge S_2(m,n,a,o,dd) \rightarrow T(m,a,hr,o))$

$$M_1 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \qquad M_2 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$
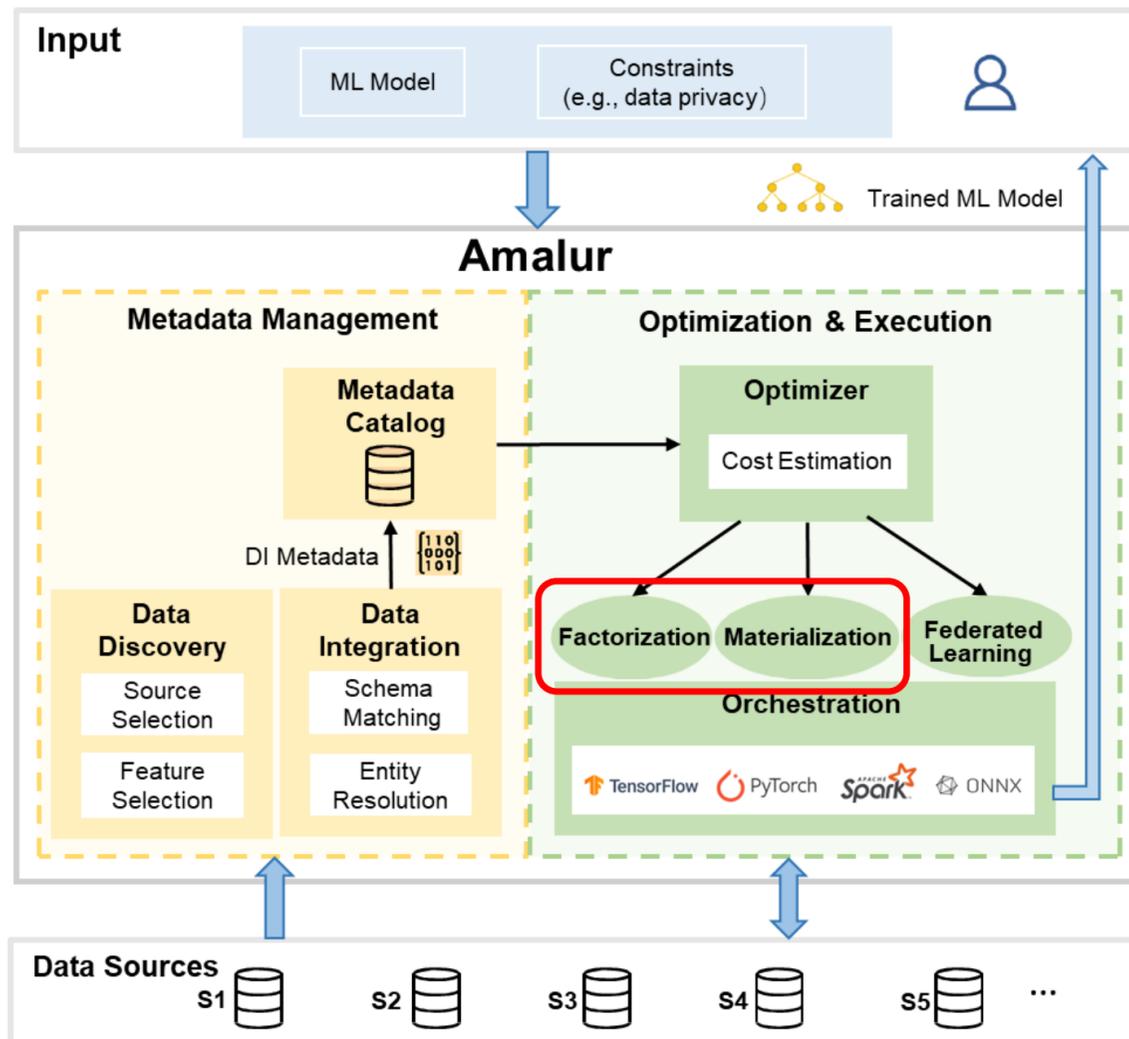
(a) Mapping matrix

Row matching (entity resolution)

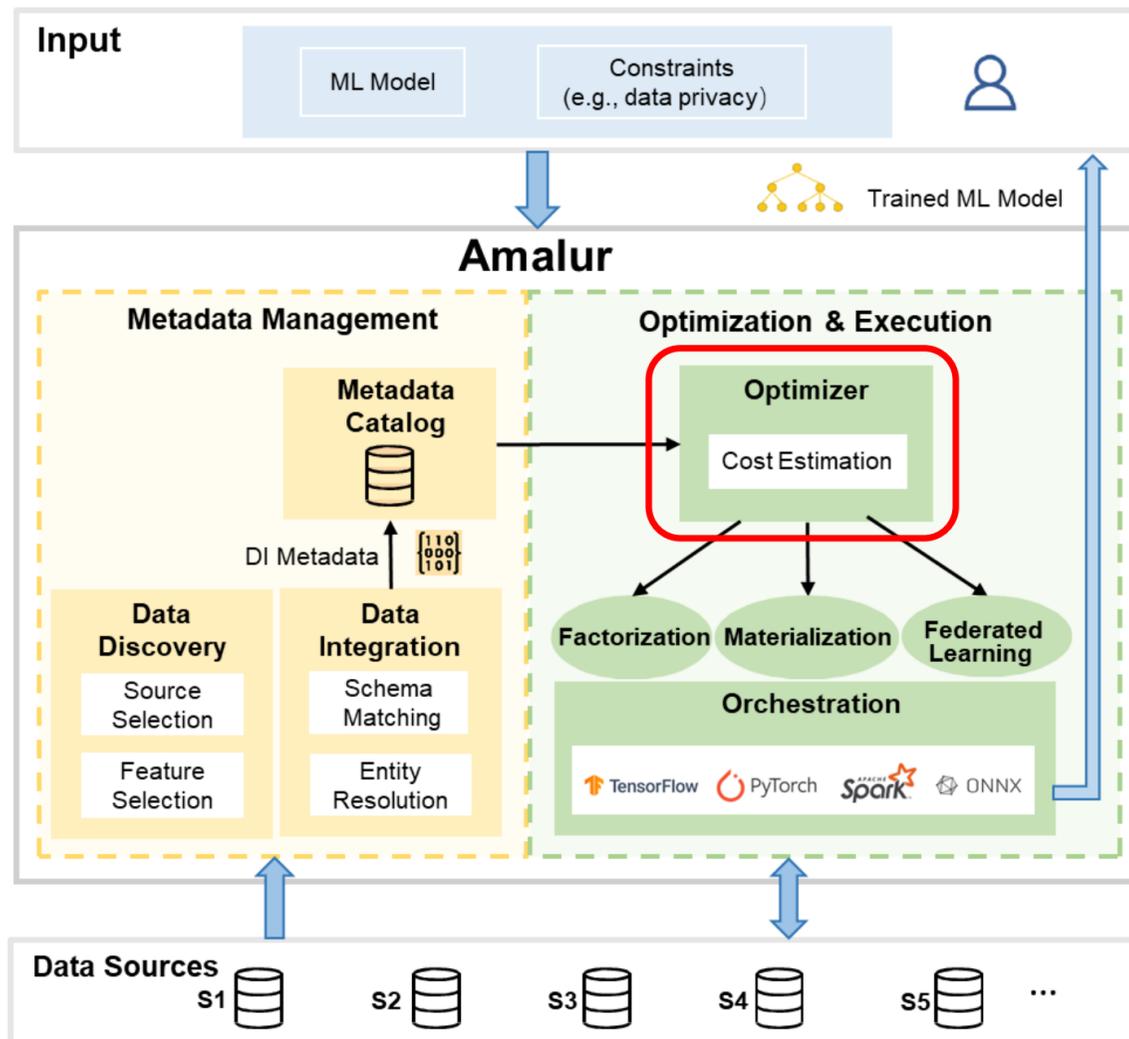| | | | | | | |
|---|---|---|---|---|---|---|
| $S_1$ | 3 | 1 | Jane | 37 | 70 | |
| $S_2$ | 2 | 1 | Jane | 37 | 92 | 11/5/21 |

$$I_1 \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \qquad I_2 \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

(b) Indicator matrix

Factorize: $TX \rightarrow T_1 X + (T_2 \circ R_2)X$

Materialize: $T \rightarrow I_1 D_1 M_1^T + (I_2 D_2 M_2^T \circ R_2)$

# Amalur: Workflow



Input — ML Model — Constraints (e.g., data privacy)

Trained ML Model

**Amalur**

Metadata Management
- Metadata Catalog
- DI Metadata
- Data Discovery
  - Source Selection
  - Feature Selection
- Data Integration
  - Schema Matching
  - Entity Resolution

Optimization & Execution
- Optimizer
  - Cost Estimation
- Factorization
- Materialization
- Federated Learning
- Orchestration
  - TensorFlow, PyTorch, Spark, ONNX

Data Sources: S1, S2, S3, S4, S5 …

**Column matching (schema mapping)**

$m_1: \forall m, n, a, hr, o, dd \ (S_1(m, n, a, hr) \wedge S_2(m, n, a, o, dd) \rightarrow T(m, a, hr, o))$

$M_1 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \quad M_2 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

(a) Mapping matrix

**Row matching (entity resolution)**

| | | | | | | |
|---|---|---|---|---|---|---|
| $S_1$ | 3 | 1 | Jane | 37 | 70 | |
| $S_2$ | 2 | 1 | Jane | 37 | 92 | 11/5/21 |

$I_1 \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad I_2 \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$
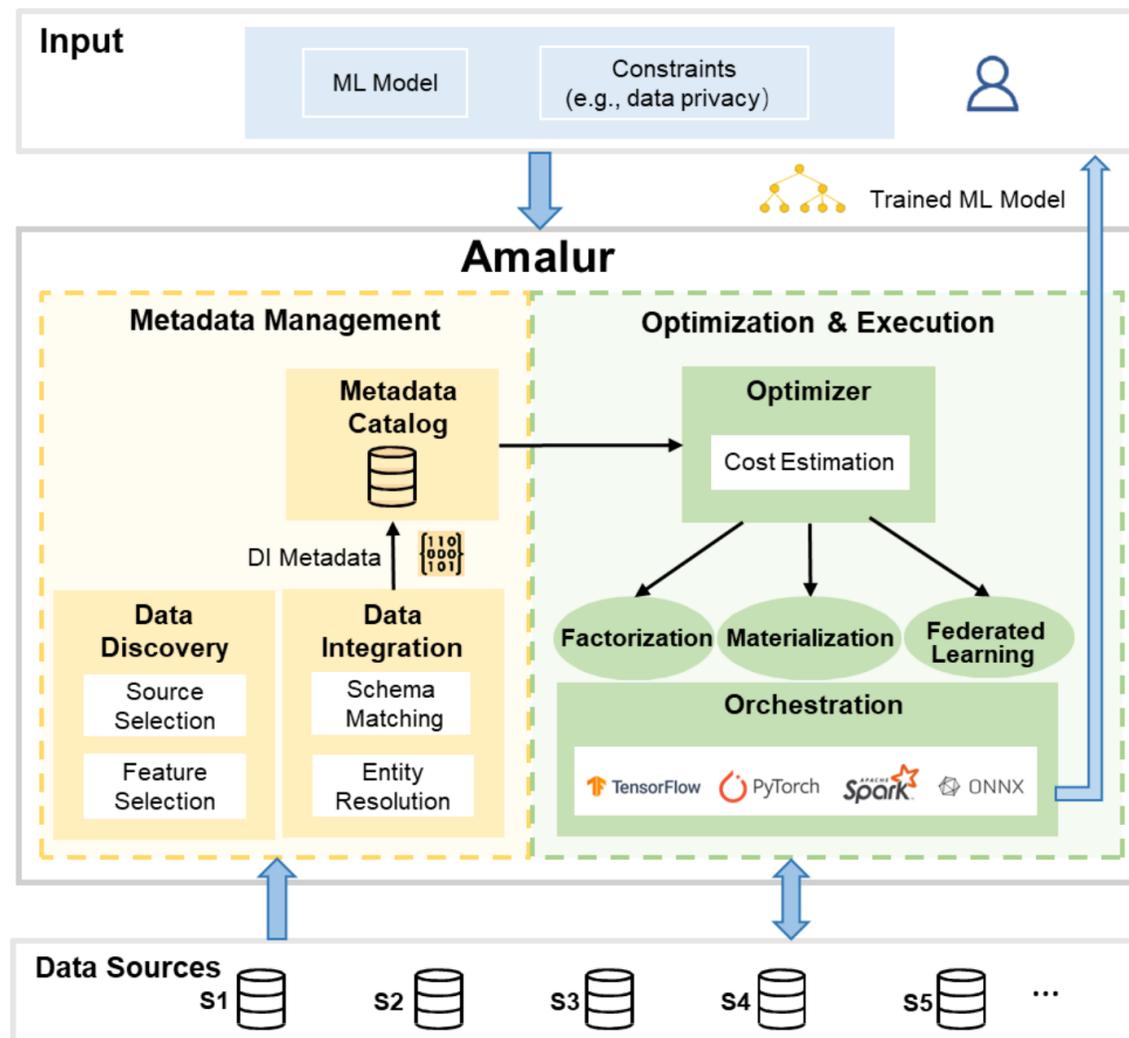
(b) Indicator matrix

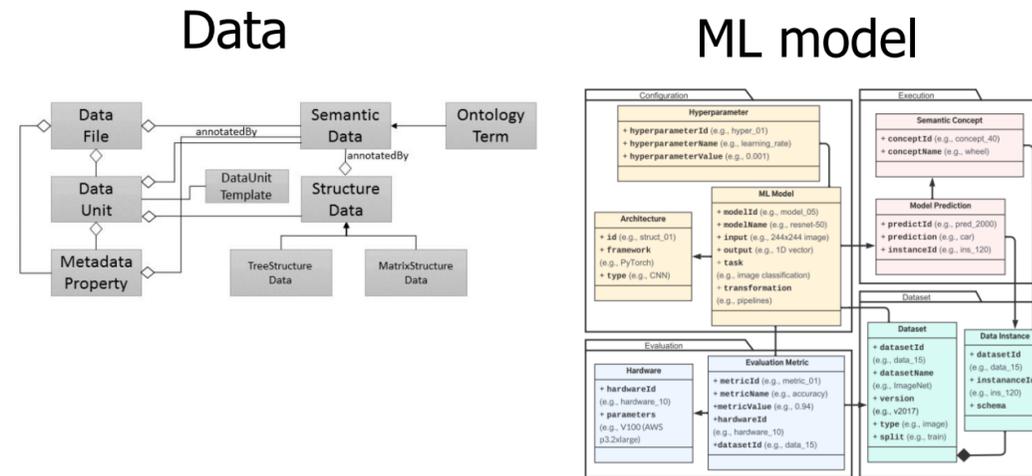Factorize: $TX \rightarrow T_1 X + (T_2 \circ R_2)X$

Materialize: $T \rightarrow I_1 D_1 M_1^T + (I_2 D_2 M_2^T \circ R_2)$

# The whole picture

## Metadata Representations

Data      ML model



Conceptual
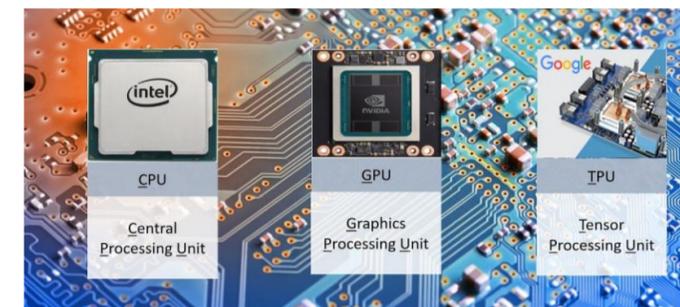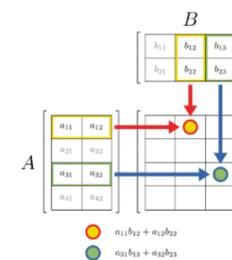
## Model Lake Amalur



## FO logic

$$m_1 : \; \forall s, n, se, a, p, d \; (S_1(s, n, se, a) \wedge S_2(s, n, se, p, d) \rightarrow T(s, se, a, p))$$

Logical

## Matrices

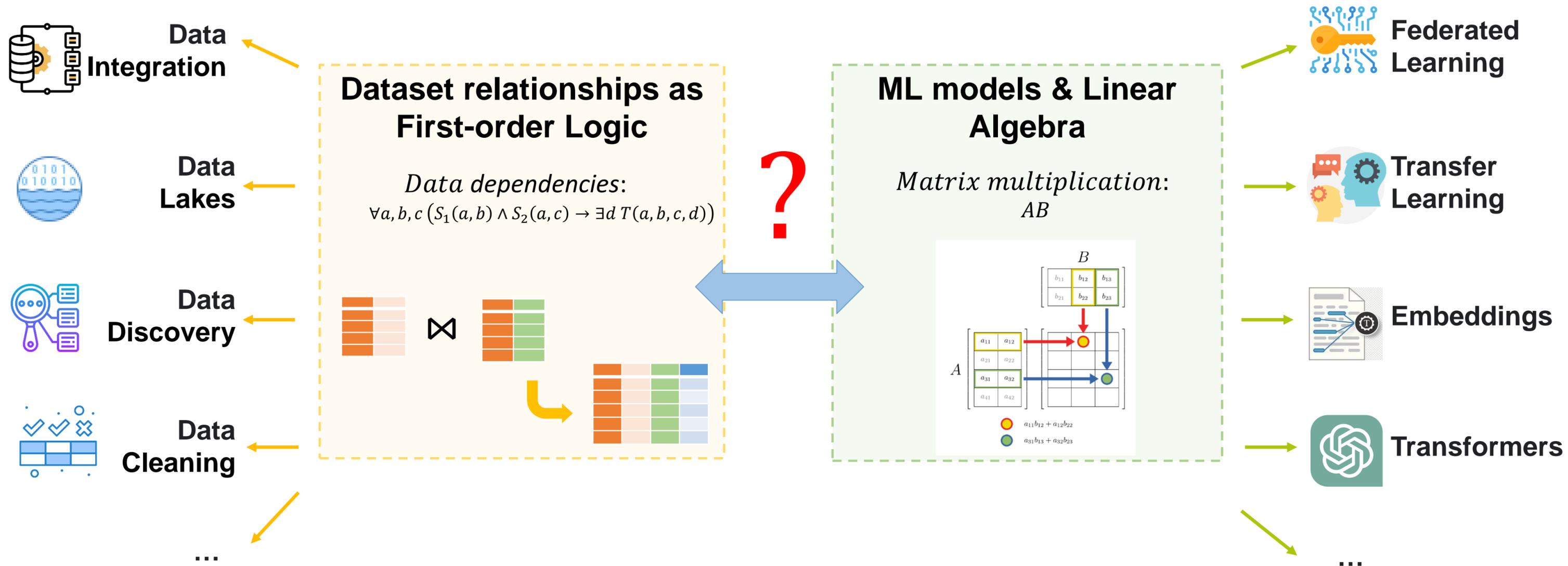$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

Physical

# Vision   Data Integration meets Machine Learning

*Q: Can we use data integration metadata to improve the effectiveness and efficiency of ML model training?*

Data Integration

Data Lakes

Data Discovery

Data Cleaning

...

**Dataset relationships as First-order Logic**

*Data dependencies*:
$$\forall a, b, c \left( S_1(a,b) \wedge S_2(a,c) \rightarrow \exists d\, T(a,b,c,d) \right)$$

**?**

**ML models & Linear Algebra**

*Matrix multiplication*:
*AB*

Federated Learning

Transfer Learning

Embeddings

Transformers

...

# We're hiring
## -- Work, live, love at Delft

- **DBML, data lakes, and federated learning**

  - PhD students
  - postdoctoral researcher

  Contact: R.Hai@tudelft.nl

- **Scalable & Consistent Function-as-a-Service Systems**
  - 2 PhD students
  - 1 postdoctoral researcher

  Contact: a.katsifodimos@tudelft.nl