# Advanced Data Management (CSCI 640/490)
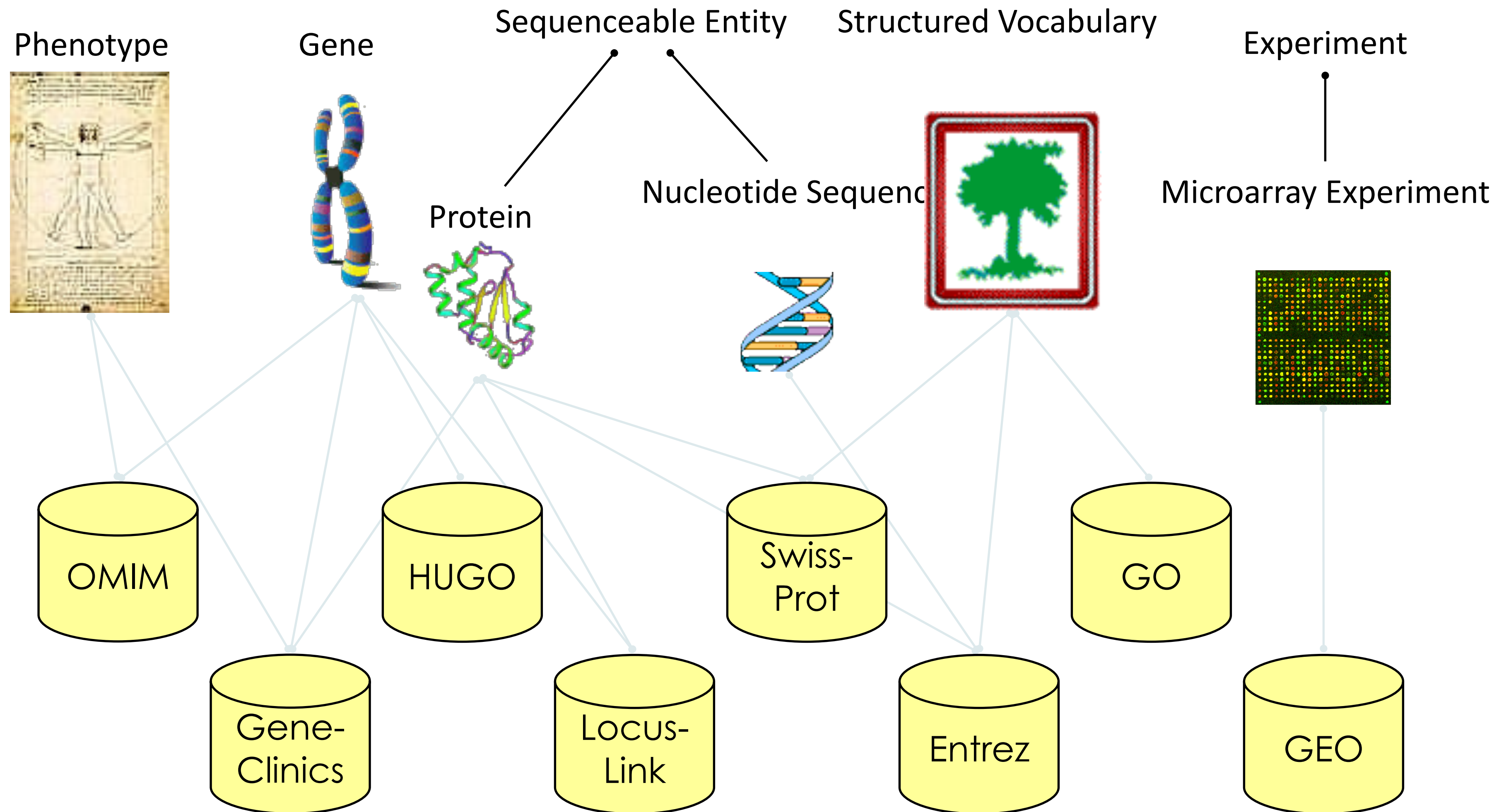
Data Fusion

Dr. David Koop

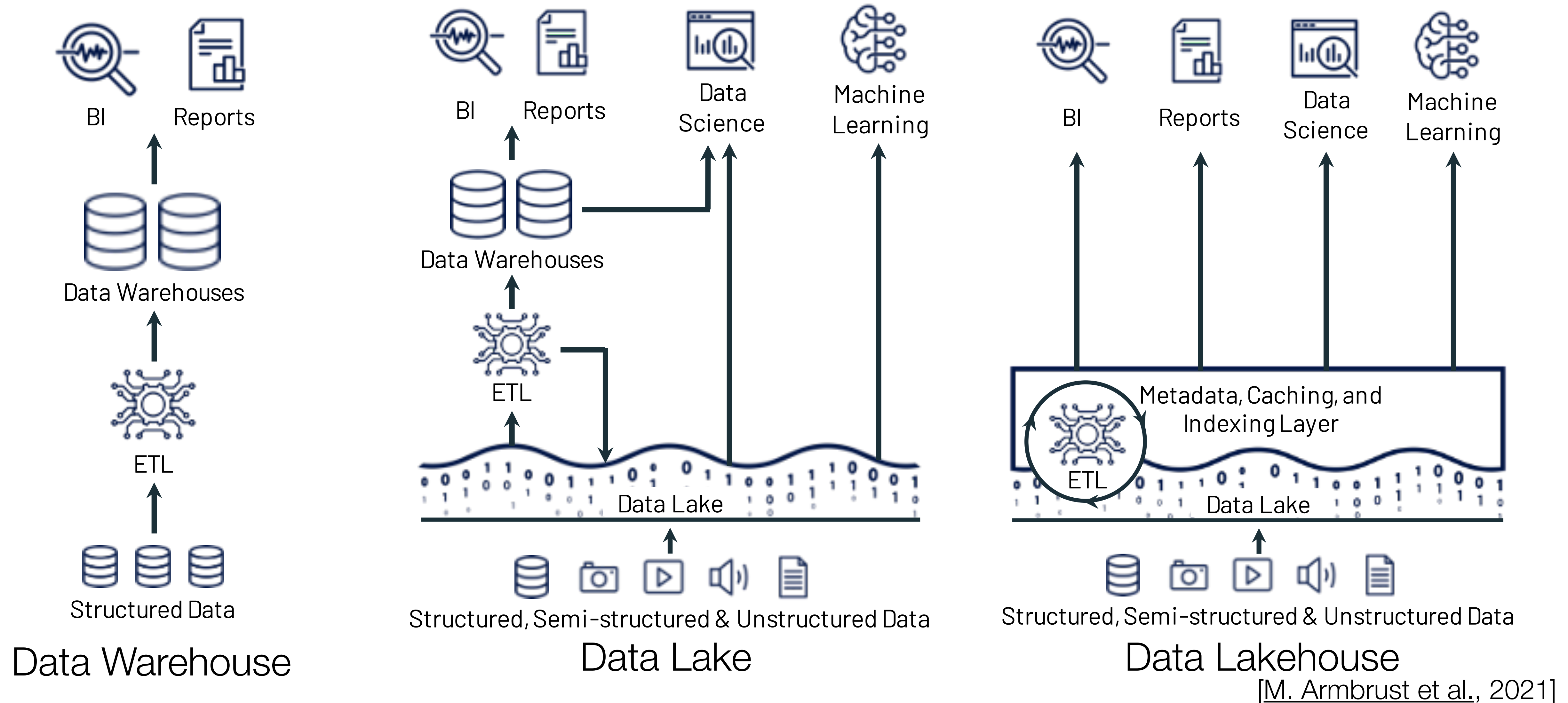Northern Illinois University

# Reading Quiz

# Data Integration: Combine Datasets with Different Data



[A. Doan et al., 2012]

# Storage for Data Analysis
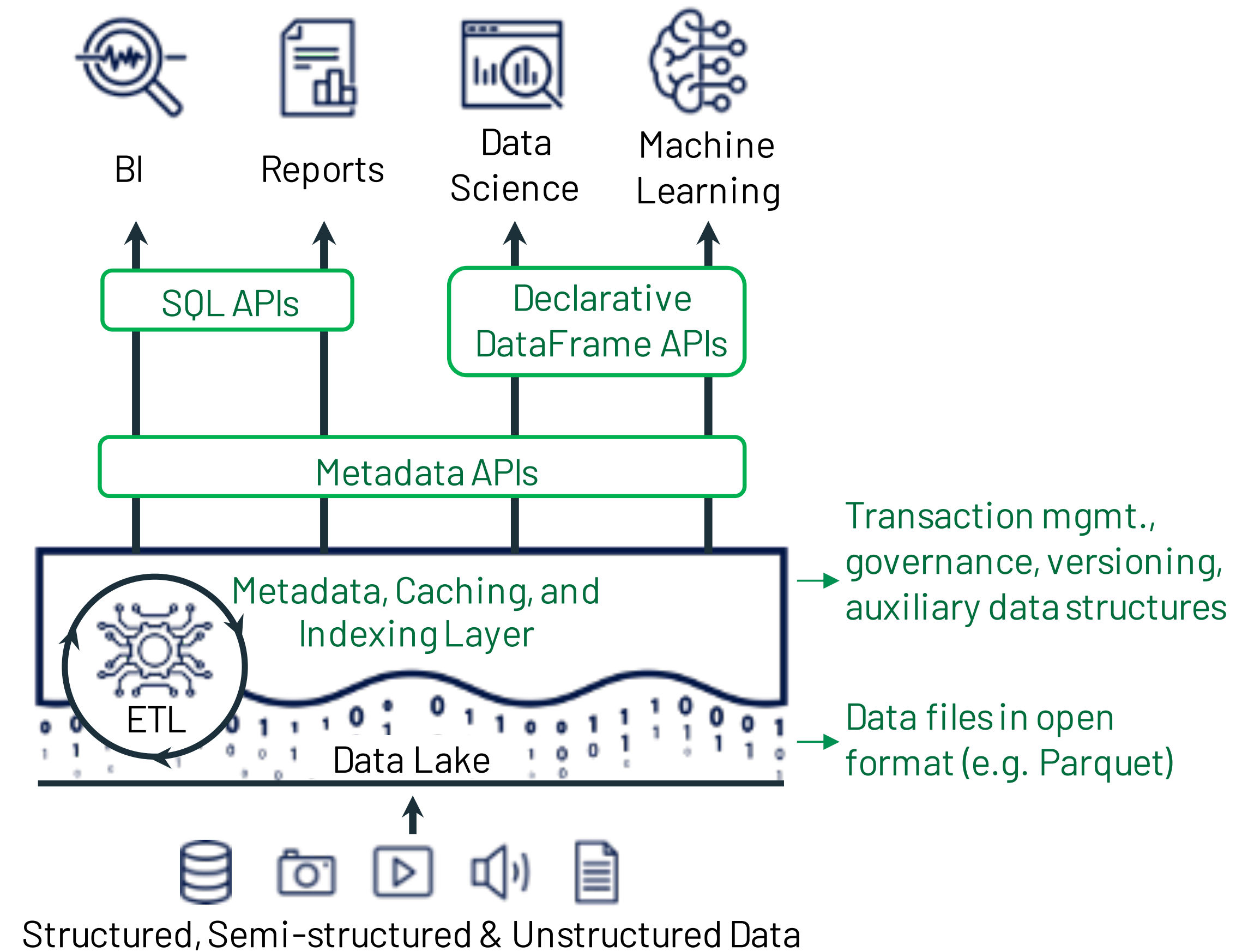


Data Warehouse

Data Lake
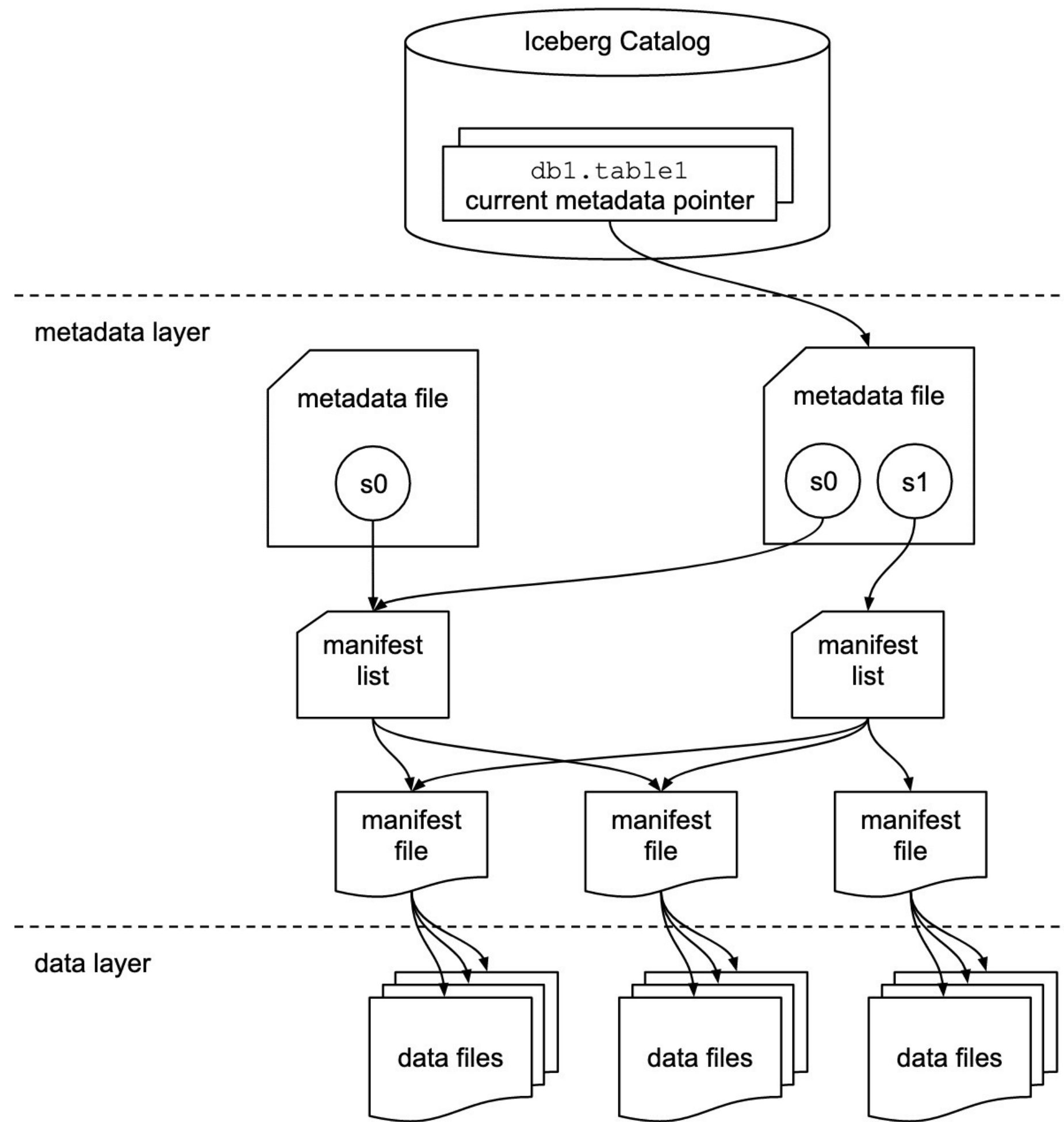
Data Lakehouse

[M. Armbrust et al., 2021]

# Data Lakehouse

- Solutions:

  - Reliable data management on data lakes: add metadata APIs

  - Support for machine learning and data science: allow use of declarative dataframe APIs

  - SQL performance: allow use of SQL APIs



[M. Armbrust et al., 2021]

# Apache Iceberg



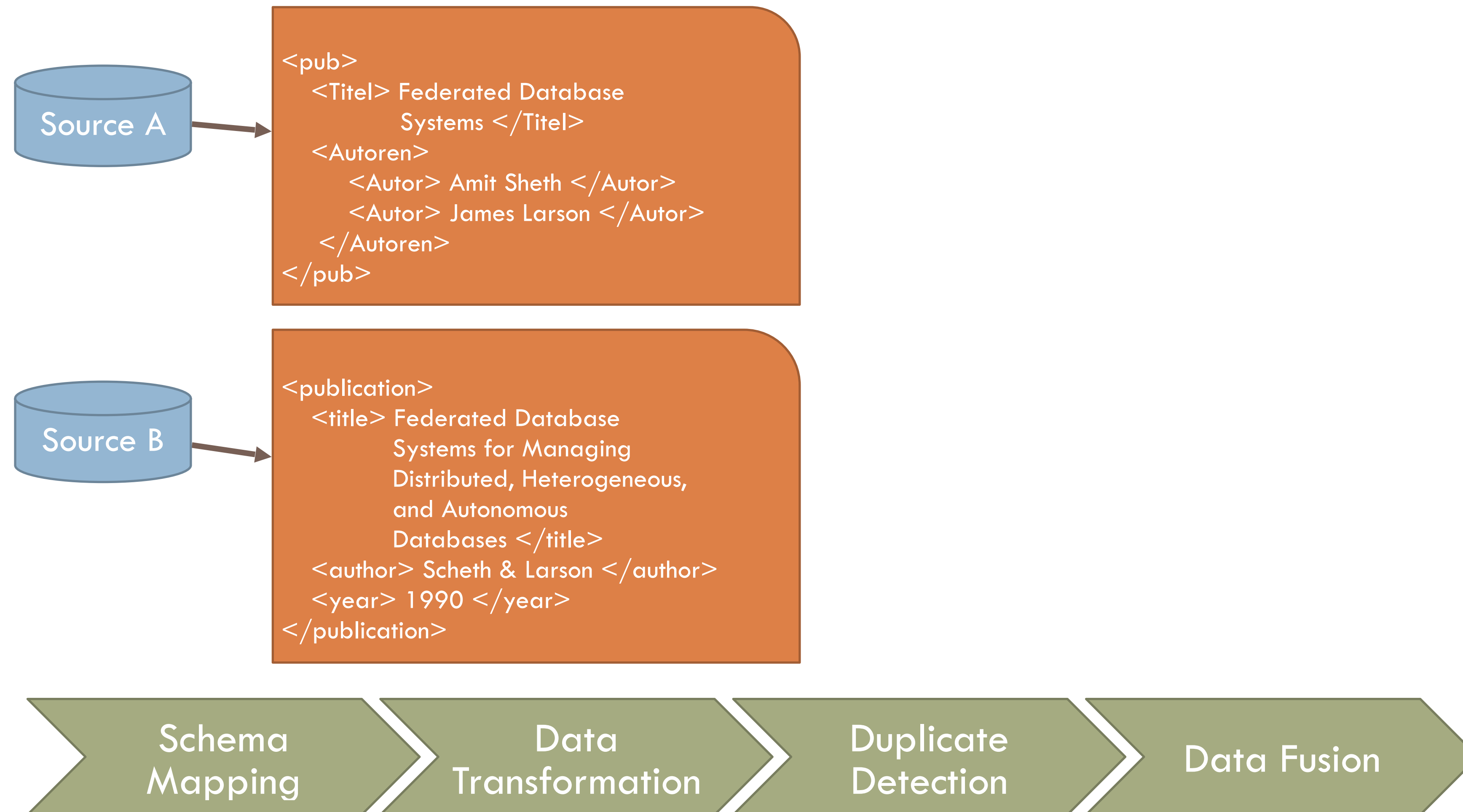- Data Files: store the actual data (parquet, avro, orc)
- Manifest Files: track a group of data files (and delete files); have metadata for filtering (min/max)
- Manifest Lists: which manifest files make up a table at a given point in time (snapshot)
- Metadata file: keeps track of table creates or data add/delete
- Catalog: Tracks the tables and pointer to the most recently created metadata file

[A. Merced, 2022]

# Information Integration

Source A

```
<pub>
   <Titel> Federated Database
          Systems </Titel>
   <Autoren>
       <Autor> Amit Sheth </Autor>
       <Autor> James Larson </Autor>
   </Autoren>
</pub>
```

Source B

```
<publication>
   <title> Federated Database
          Systems for Managing
          Distributed, Heterogeneous,
          and Autonomous
          Databases </title>
   <author> Scheth & Larson </author>
   <year> 1990 </year>
</publication>
```

Schema Mapping → Data Transformation → Duplicate Detection → Data Fusion

[L. Dong and F. Naumann, 2009]

# Information Integration

Source A

Source B

```
<pub>
  <Titel> Federated Database
       Systems </Titel>
  <Autoren>
    <Autor> Amit Sheth </Autor>
    <Autor> James Larson </Autor>
  </Autoren>
</pub>
```

```
<pub>
  <title> </title>
  <Autoren>
    <author> </author>
    <author> </author>
  </Autoren>
  <year> </year>
</pub>
```

```
<publication>
  <title> Federated Database
       Systems for Managing
       Distributed, Heterogeneous,
       and Autonomous
       Databases </title>
  <author> Sheth & Larson </author>
  <year> 1990 </year>
</publication>
```

Schema Integration

Schema Mapping

| Schema Mapping | Data Transformation | Duplicate Detection | Data Fusion |

[L. Dong and F. Naumann, 2009]

# Information Integration

Transformation queries or views

Source A

```
<pub>
    <Titel> Federated Database
            Systems </Titel>
    <Autoren>
        <Autor> Amit Sheth </Autor>
        <Autor> James Larson </Autor>
    </Autoren>
</pub>
```

XQuery

```
<pub>
    <title> Federated Database
            Systems </title>
    <Autoren>
        <author> Amit Sheth </author>
        <author> James Larson </author>
    </Autoren>
</pub>
<pub>
    <title> Federated Database Systems for
            Managing Distributed,
            Heterogeneous, and Autonomous
            Databases </title>
    <Autoren>
        <author> Scheth & Larson </author>
    </Autoren>
    <year> 1990 </year>
</pub>
```

Source B

```
<publication>
    <title> Federated Database
            Systems for Managing
            Distributed, Heterogeneous,
            and Autonomous
            Databases </title>
    <author> Scheth & Larson </author>
    <year> 1990 </year>
</publication>
```

XQuery

Schema Mapping → Data Transformation → Duplicate Detection → Data Fusion

[L. Dong and F. Naumann, 2009]

# Information Integration



Source A

```
<pub>
   <Titel> Federated Database
          Systems </Titel>
   <Autoren>
      <Autor> Amit Sheth </Autor>
      <Autor> James Larson </Autor>
   </Autoren>
</pub>
```

Source B

```
<publication>
   <title> Federated Database
          Systems for Managing
          Distributed, Heterogeneous,
          and Autonomous
          Databases </title>
   <author> Scheth & Larson </author>
   <year> 1990 </year>
</publication>
```

```
<pub>
   <title> Federated Database
          Systems </title>
   <Autoren>
      <author> Amit Sheth </author>
      <author> James Larson </author>
   </Autoren>
</pub>
<pub>
   <title> Federated Database Systems for
          Managing Distributed,
          Heterogeneous, and Autonomous
          Databases </title>
   <Autoren>
      <author> Scheth & Larson </author>
   </Autoren>
 <year> 1990 </year>
</pub>
```

Schema Mapping → Data Transformation → Duplicate Detection → Data Fusion

[L. Dong and F. Naumann, 2009]

# "Duplicate Detection" has many Duplicates

[L. Dong and F. Naumann, 2009]

# "Duplicate Detection" has many Duplicates
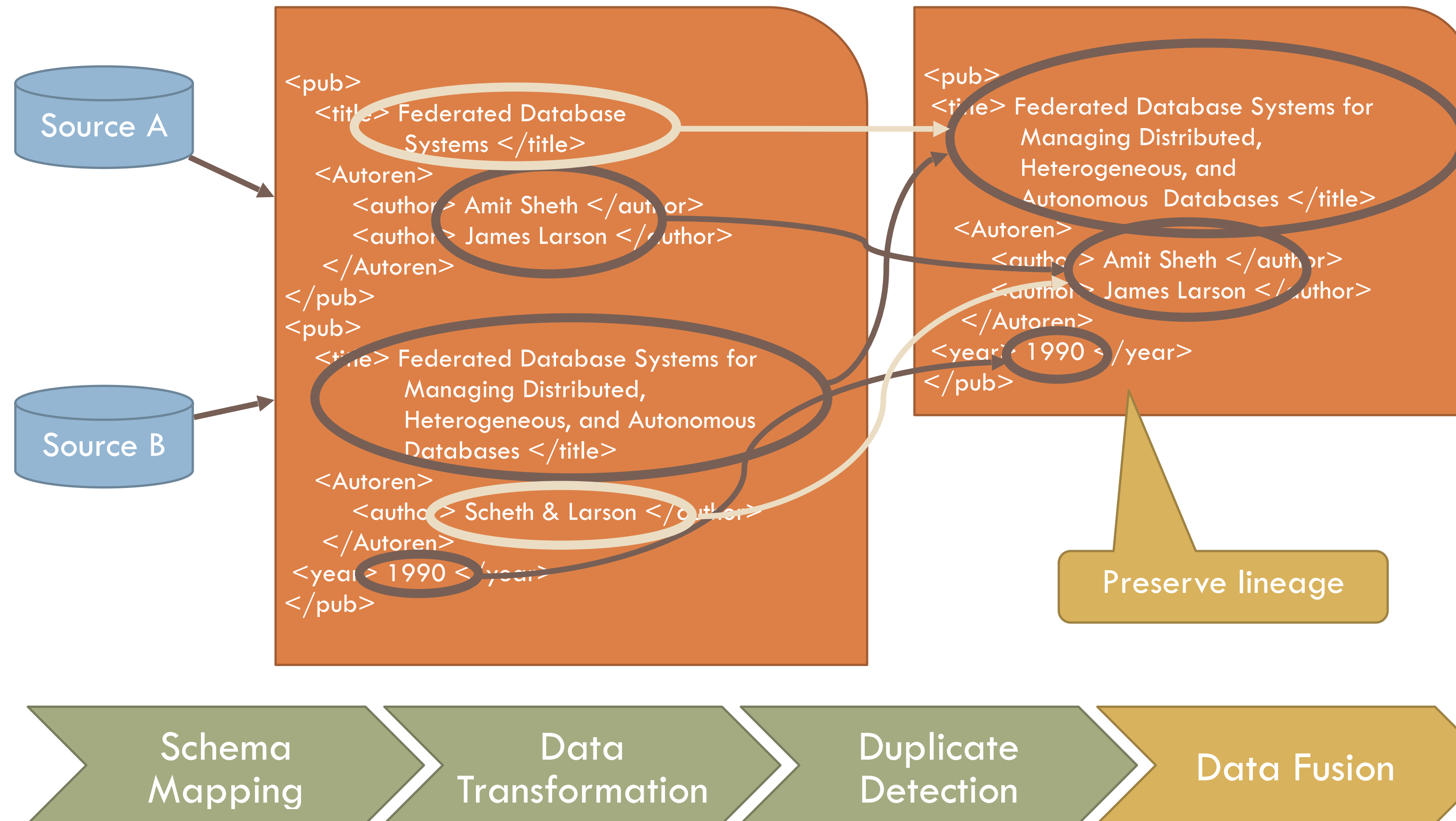
Household matching

Doubles

Duplicate detection

Mixed and split citation problem

Record linkage

Object identification

Match

Deduplication

Fuzzy match

Object consolidation

Entity resolution

Entity clustering

Approximate match

Identity uncertainty

Reference reconciliation

Merge/purge

Hardening soft databases

Reference matching

Householding

[L. Dong and F. Naumann, 2009]

# Record Linkage Process



[P. Christen , 2019]

# Record Linkage Techniques

- Deterministic matching

  - Rule-based matching (complex to build and maintain)

- Probabilistic record linkage [Fellegi and Sunter, 1969]

  - Use available attributes for linking (often personal information, like names, addresses, dates of birth, etc.)

  - Calculate match weights for attributes

- "Computer science" approaches

  - Based on machine learning, data mining, database, or information retrieval techniques

  - Supervised classification: Requires training data (true matches)

  - Unsupervised: Clustering, collective, and graph based

[P. Christen , 2019]

# Information Integration

[L. Dong and F. Naumann, 2009]

# Assignment 3

- Ask a Manager Salary Data

- Use Polars & OpenRefine
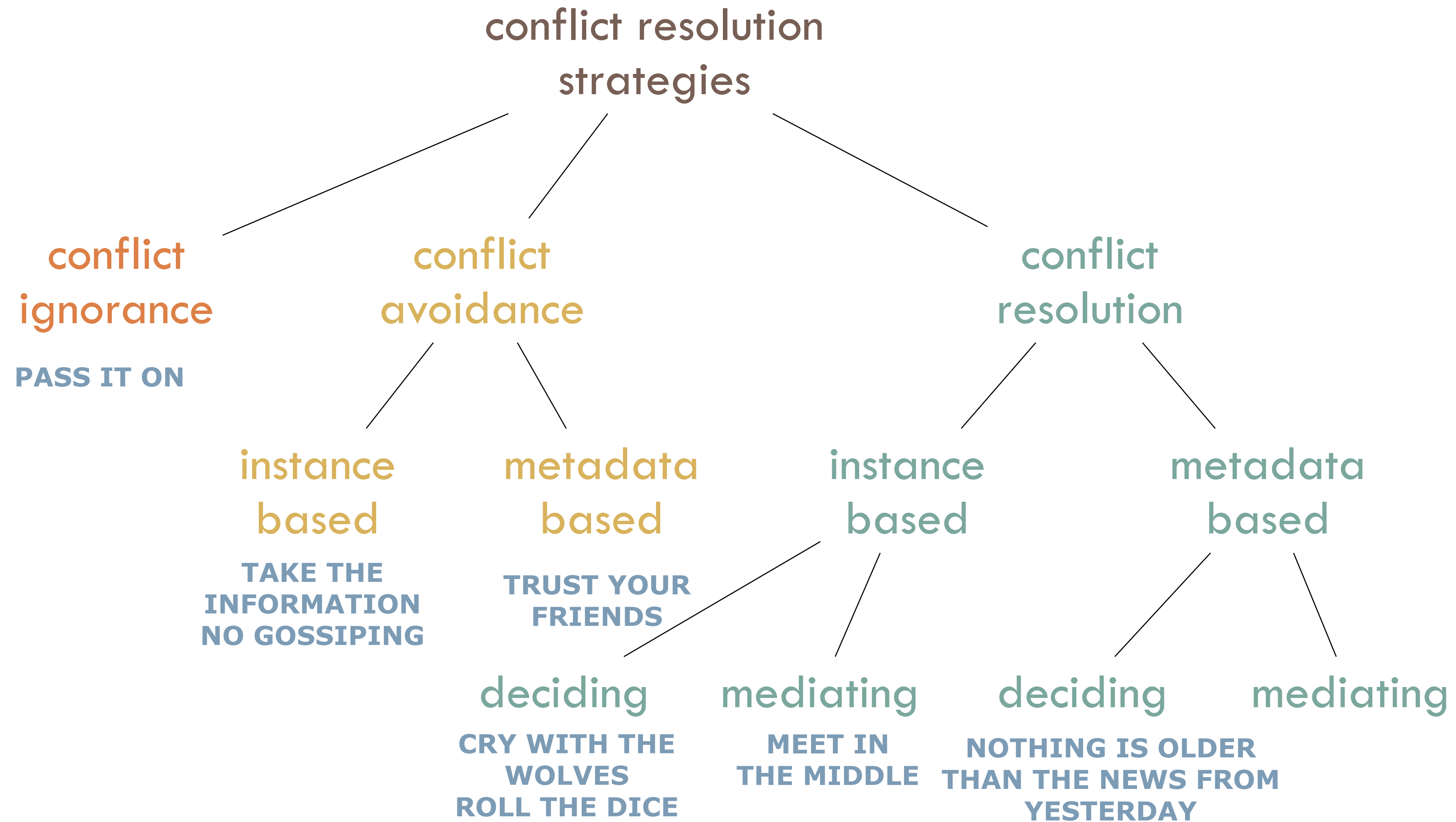
- Moved deadline to next Tuesday, October. 21

# Data Fusion

# Data Fusion

- Problem: Given a duplicate, create a single object representation while resolving conflicting data values.

- Difficulties:

  - Null values: Subsumption and complementation

  - Contradictions in data values

  - Uncertainty & truth: Discover the true value and model uncertainty in this process

  - Metadata: Preferences, recency, correctness

  - Lineage: Keep original values and their origin

  - Implementation in DBMS: SQL, extended SQL, UDFs, etc.

# Conflict Resolution Strategies



conflict resolution strategies

- conflict ignorance
  - PASS IT ON
- conflict avoidance
  - instance based
    - TAKE THE INFORMATION
    - NO GOSSIPING
  - metadata based
    - TRUST YOUR FRIENDS
- conflict resolution
  - instance based
    - deciding
      - CRY WITH THE WOLVES
      - ROLL THE DICE
    - mediating
      - MEET IN THE MIDDLE
  - metadata based
    - deciding
      - NOTHING IS OLDER THAN THE NEWS FROM YESTERDAY
    - mediating

[L. Dong and F. Naumann, 2009]

# Integrating Conflicting Data: The Role of Source Dependence

X. L. Dong, L. Berti-Equille, and D. Srivastava

Northern Illinois University

# Discussion

- What is the paper's main contribution?

- Do you buy the argument? Any issues with the experiments?

- Can you think of any scenarios where the proposed technique will fail?

- Questions?

# Example Problem

[X L Dong et al., 2009]

# Example Problem

| | S1 | S2 | S3 |
|---|---|---|---|
| Stonebraker | MIT | Berkeley | MIT |
| Dewitt | MSR | MSR | UWisc |
| Bernstein | MSR | MSR | MSR |
| Carey | UCI | AT&T | BEA |
| Halevy | Google | Google | UW |

[X L Dong et al., 2009]

# Naive Voting Works

| | S1 | S2 | S3 |
|---|---|---|---|
| Stonebraker | MIT | Berkeley | MIT |
| Dewitt | MSR | MSR | UWisc |
| Bernstein | MSR | MSR | MSR |
| Carey | UCI | AT&T | BEA |
| Halevy | Google | Google | UW |

[X L Dong et al., 2009]

# Naive Voting Only Works if Data Sources are Independent

[X L Dong et al., 2009]

Northern Illinois University

# Naive Voting Only Works if Data Sources are Independent

|  | S1 | S2 | S3 | S4 | S5 |
|---|---|---|---|---|---|
| Stonebraker | MIT | Berkeley | MIT | MIT | MS |
| Dewitt | MSR | MSR | UWisc | UWisc | UWisc |
| Bernstein | MSR | MSR | MSR | MSR | MSR |
| Carey | UCI | AT&T | BEA | BEA | BEA |
| Halevy | Google | Google | UW | UW | UW |

[X L Dong et al., 2009]

# S4 and S5 copy from S3

| | S1 | S2 | S3 | S4 | S5 |
|---|---|---|---|---|---|
| Stonebraker | MIT | Berkeley | MIT | MIT | MS |
| Dewitt | MSR | MSR | UWisc | UWisc | UWisc |
| Bernstein | MSR | MSR | MSR | MSR | MSR |
| Carey | UCI | AT&T | BEA | BEA | BEA |
| Halevy | Google | Google | UW | UW | UW |

[X L Dong et al., 2009]

# S4 and S5 copy from S3

| | S1 | S2 | S3 | S4 | S5 |
|---|---|---|---|---|---|
| Stonebraker | MIT | Berkeley | MIT | MIT | MS |
| Dewitt | MSR | MSR | UWisc | UWisc | UWisc |
| Bernstein | MSR | MSR | MSR | MSR | MSR |
| Carey | UCI | AT&T | BEA | BEA | BEA |
| Halevy | Google | Google | UW | UW | UW |

[X L Dong et al., 2009]

# Challenges in Dependence Discovery

| | S1 | S2 | S3 | S4 | S5 |
|---|---|---|---|---|---|
| Stonebraker | MIT | Berkeley | MIT | MIT | MS |
| Dewitt | MSR | MSR | UWisc | UWisc | UWisc |
| Bernstein | MSR | MSR | MSR | MSR | MSR |
| Carey | UCI | AT&T | BEA | BEA | BEA |
| Halevy | Google | Google | UW | UW | UW |

[X L Dong et al., 2009]

# Challenges in Dependence Discovery

| | S1 | S2 | S3 | S4 | S5 |
|---|---|---|---|---|---|
| Stonebraker | MIT | Berkeley | MIT | MIT | MS |
| Dewitt | MSR | MSR | UWisc | UWisc | UWisc |
| Bernstein | MSR | MSR | MSR | MSR | MSR |
| Carey | UCI | AT&T | BEA | BEA | BEA |
| Halevy | Google | Google | UW | UW | UW |

[X L Dong et al., 2009]

# Challenges in Dependence Discovery

2. With only a snapshot it is hard to decide which source is a copier.

|  | S1 | S2 | S3 | S4 | S5 |
|---|---|---|---|---|---|
| Stonebraker | MIT | Berkeley | MIT | MIT | MS |
| Dewitt | MSR | MSR | UWisc | UWisc | UWisc |
| Bernstein | MSR | MSR | MSR | MSR | MSR |
| Carey | UCI | AT&T | BEA | BEA | BEA |
| Halevy | Google | Google | UW | UW | UW |

[X L Dong et al., 2009]

# Challenges in Dependence Discovery

1. Sharing common data does not in itself imply copying.

2. With only a snapshot it is hard to decide which source is a copier.

| | S1 | S2 | S3 | S4 | S5 |
|---|---|---|---|---|---|
| Stonebraker | MIT | Berkeley | MIT | MIT | MS |
| Dewitt | MSR | MSR | UWisc | UWisc | UWisc |
| Bernstein | MSR | MSR | MSR | MSR | MSR |
| Carey | UCI | AT&T | BEA | BEA | BEA |
| Halevy | Google | Google | UW | UW | UW |

3. A copier can also provide or verify some data by itself, so it is inappropriate to ignore all of its data.

[X L Dong et al., 2009]

# Source Dependence

- Source dependence: two sources S and T deriving the same part of data directly or transitively from a common source (can be one of S or T).

  - Independent source

  - Copier

    - copying part (or all) of data from other sources

    - may verify or revise some of the copied values

    - may add additional values

- Assumptions

  - Independent values

  - Independent copying

  - No loop copying

[X L Dong et al., 2009]

Northern Illinois University

# Core Case

- Conditions
  - Same source accuracy
  - Uniform false-value distribution
  - Categorical value
- Proposition: W. independent "good" sources, Naïve voting selects values with highest probability to be true.
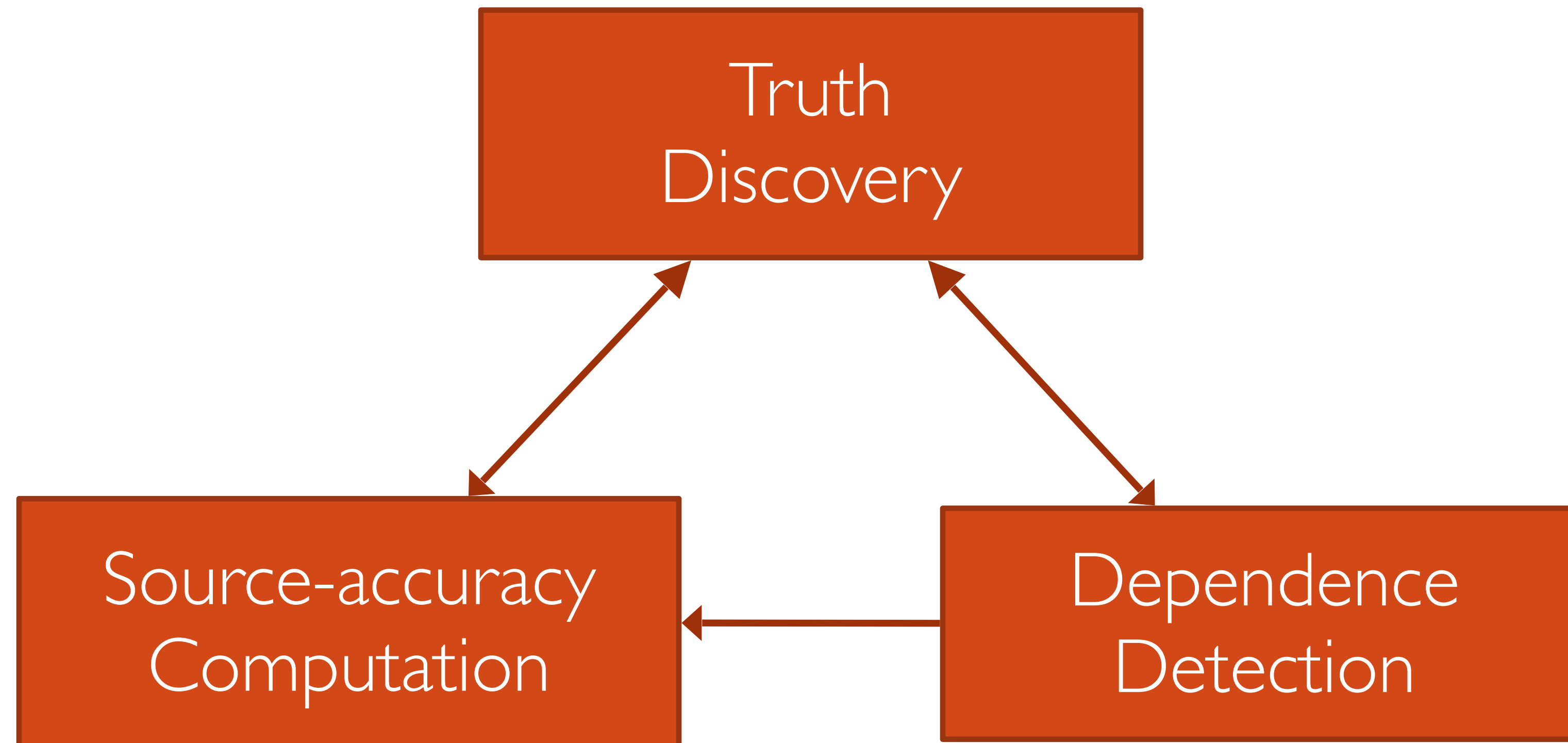
# Ideas

- If two sources share a lot of false values, they are more likely to be dependent.

- S1 is more likely to copy from S2, if the accuracy of the common data is highly different from the accuracy of S1.

[X L Dong et al., 2009]
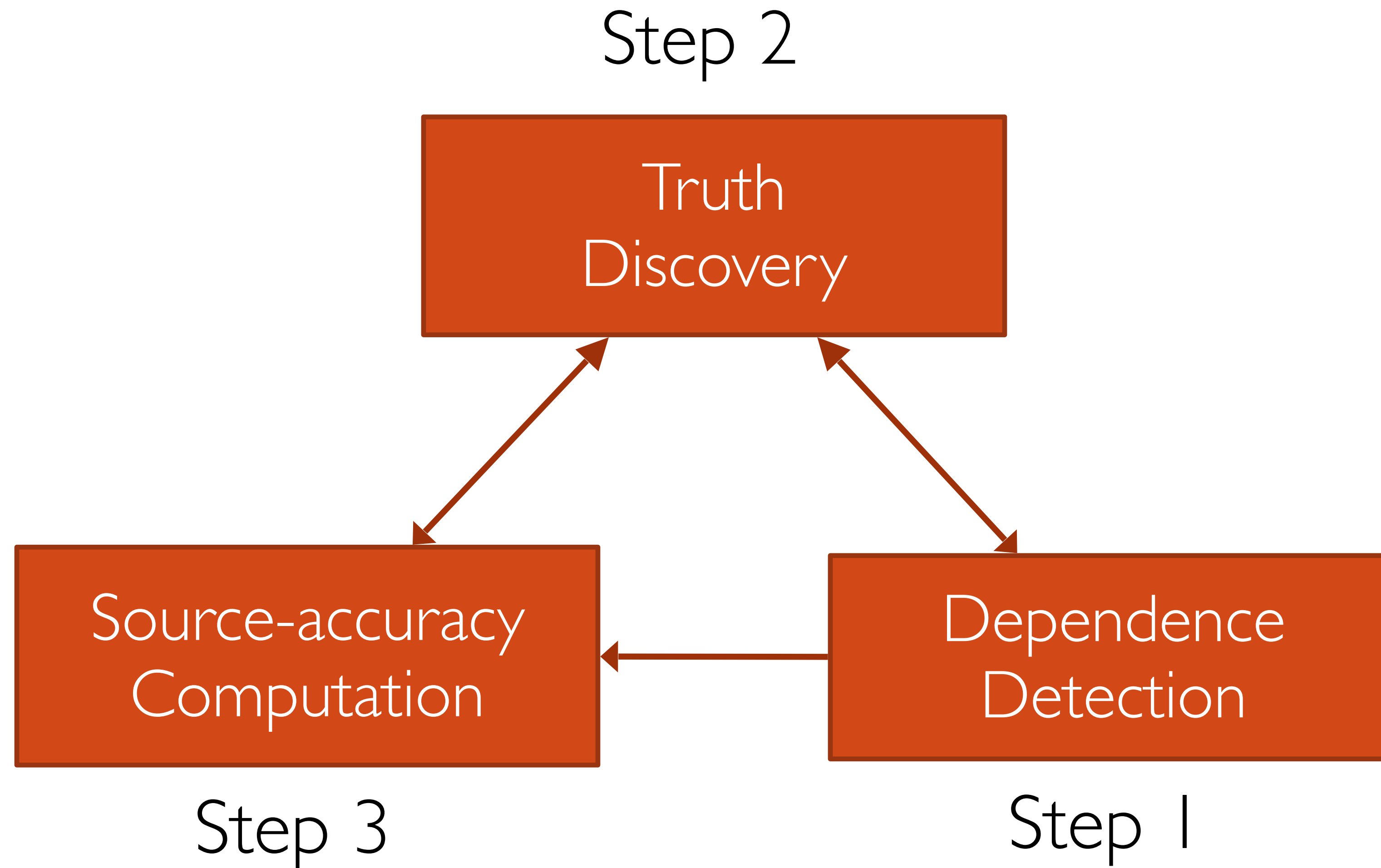
# Combining Accuracy and Dependence



Truth Discovery → Dependence Detection → Source-accuracy Computation → Truth Discovery

[X L Dong et al., 2009]

# Combining Accuracy and Dependence

Step 2

Truth
Discovery

Source-accuracy
Computation

Dependence
Detection

Step 3

Step 1

[X L Dong et al., 2009]

# The Motivating Example

|  | S1 | S2 | S3 | S4 | S5 |
|---|---|---|---|---|---|
| Stonebraker | MIT | Berkeley | MIT | MIT | MS |
| Dewitt | MSR | MSR | UWisc | UWisc | UWisc |
| Bernstein | MSR | MSR | MSR | MSR | MSR |
| Carey | UCI | AT&T | BEA | BEA | BEA |
| Halevy | Google | Google | UW | UW | UW |



[X L Dong et al., 2009]

# The Motivating Example

| Accuracy | S1 | S2 | S3 | S4 | S5 |
|---|---|---|---|---|---|
| *Round 1* | .52 | .42 | .53 | .53 | .53 |
| *Round 2* | .63 | .46 | .55 | .55 | .55 |
| *Round 3* | .71 | .52 | .53 | .53 | .37 |
| *Round 4* | .79 | .57 | .48 | .48 | .31 |
| … | … | … | … | … | … |
| *Round 11* | .97 | .61 | .40 | .40 | .21 |

| Value Confidence | Carey | | | Halevy | |
|---|---|---|---|---|---|
| | **UCI** | AT&T | BEA | **Google** | UW |
| *Round 1* | 1.61 | 1.61 | 2.0 | 2.1 | 2.0 |
| *Round 2* | 1.68 | 1.3 | 2.12 | 2.74 | 2.12 |
| *Round 3* | 2.12 | 1.47 | 2.24 | 3.59 | 2.24 |
| *Round 4* | 2.51 | 1.68 | 2.14 | 4.01 | 2.14 |
| … | … | … | … | … | … |
| *Round 11* | 4.73 | 2.08 | 1.47 | 6.67 | 1.47 |

[X L Dong et al., 2009]