



IBM Research

Model Management and Schema Mappings: Theory and Practice (Part II)

Howard Ho

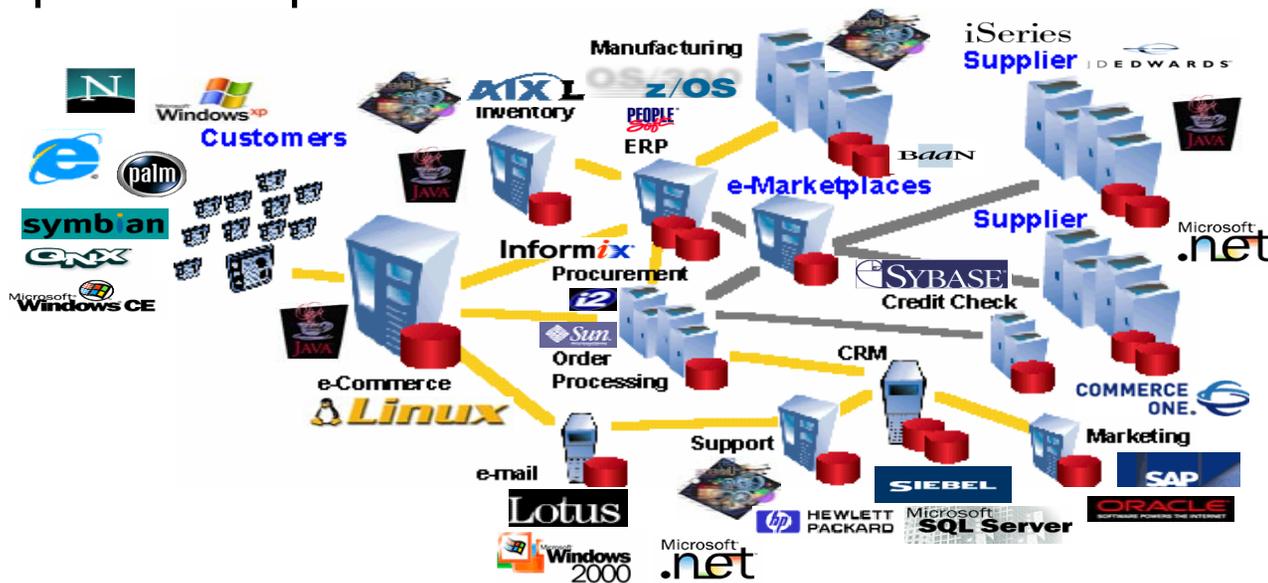
IBM Almaden Research Center

For VLDB07 Tutorial with Phil Bernstein, Microsoft Research



The Integration Challenge

- Complex and heterogeneous environments
 - Many different types of systems
 - Many inter-related applications
- Escalating needs
 - Variety, velocity, volume
- People are expensive

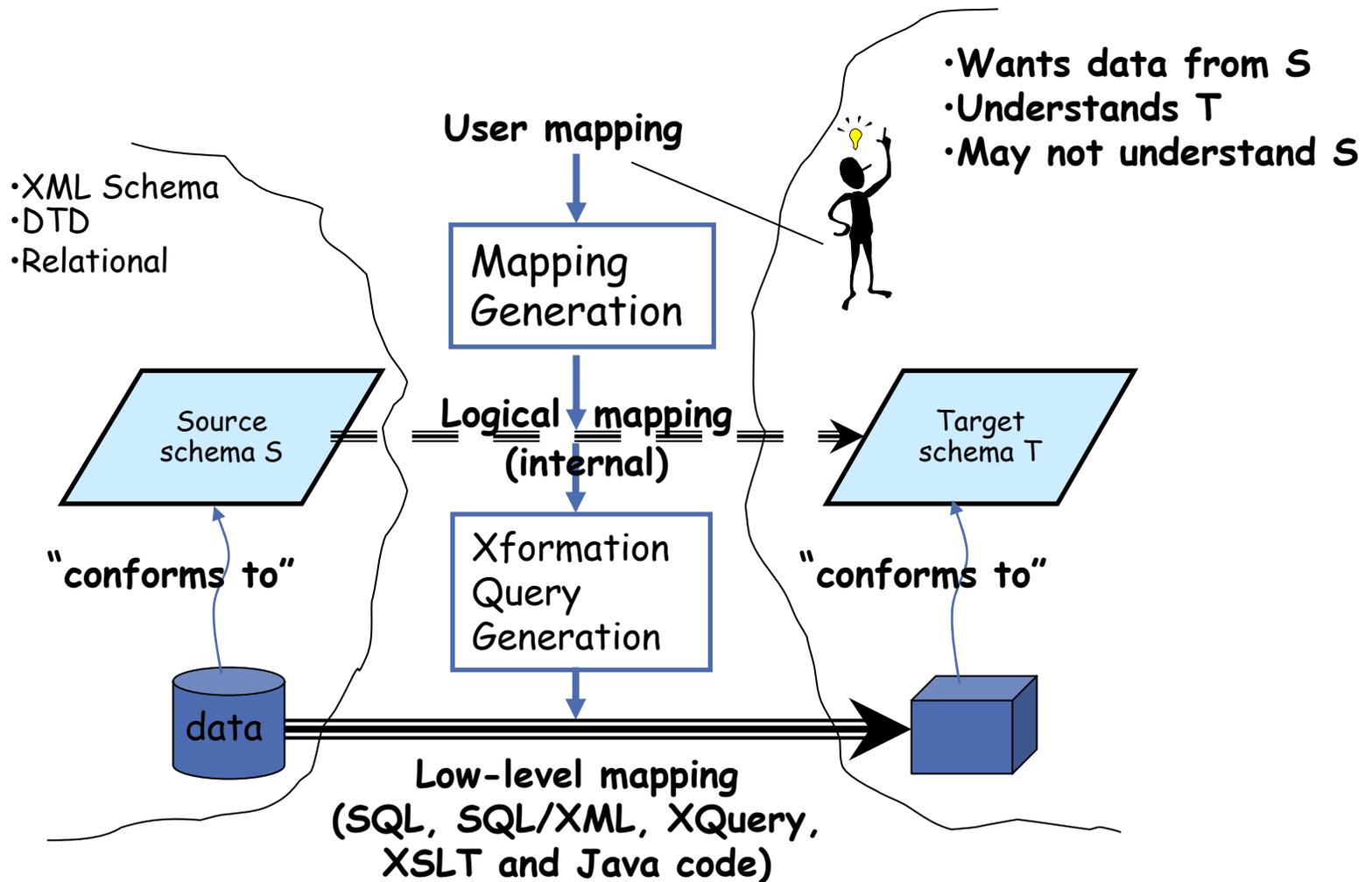


Outline

- Clio: Basic Features of a Schema Mapping System
 - Schema Matching
 - Schema Mapping
 - Query Generation
 - IBM Rational Data Architect Product
- MAUI: Advanced Features of a Schema Mapping System
 - Nested Mapping Model
 - Mapping-Based XML Transformation Engine
 - Schema Integration
 - Schema Evolution
- Clio2010: Mapping-Based Authoring of Data Flows
 - **ETL**: “Mapping ↔ ETL” Conversion
 - **Web-Service Composition**: Mapping for web-service data sources
 - **Mashups**: “Mapping → Mashup” Generation
 - **Reuse**: Mapping Polymorphism
- Conclusions and Future Directions



Clio: A Schema Mapping System

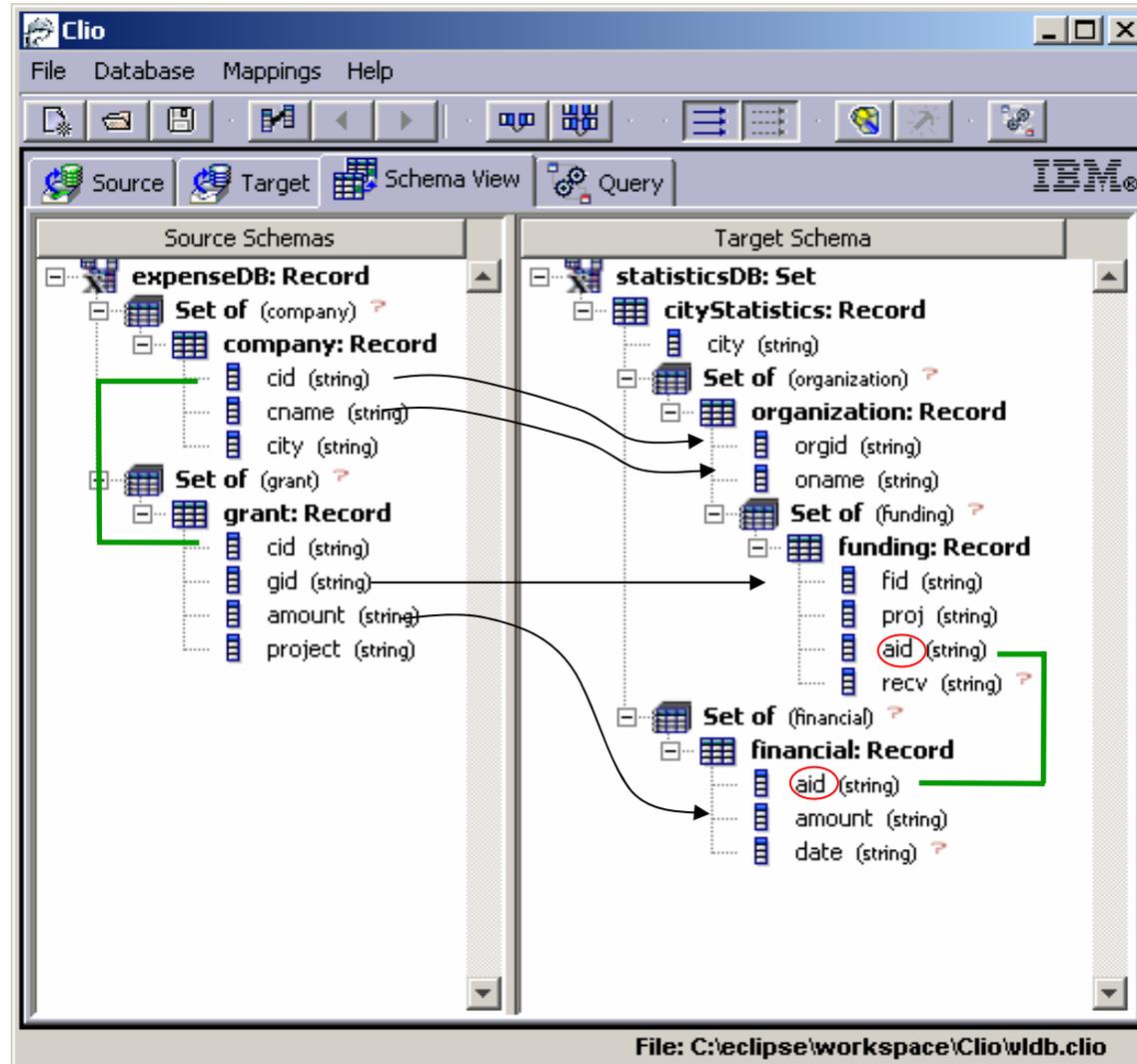


Logical mappings can be used for both **target materialization** or **query rewriting**



Major Features (and Challenges)

- Schemas can be arbitrarily different
 - Human friendly
 - Automatic discovery
- Element correspondences
 - Nested Relational Model
 - Nested Constraints
- Support Nested Structures
 - Nested Relational Model
 - Nested Constraints
- Produce Correct Grouping
- Preserve data meaning
 - Discover associations
 - Use constraints & schema
- Create New Target Values
- and ...

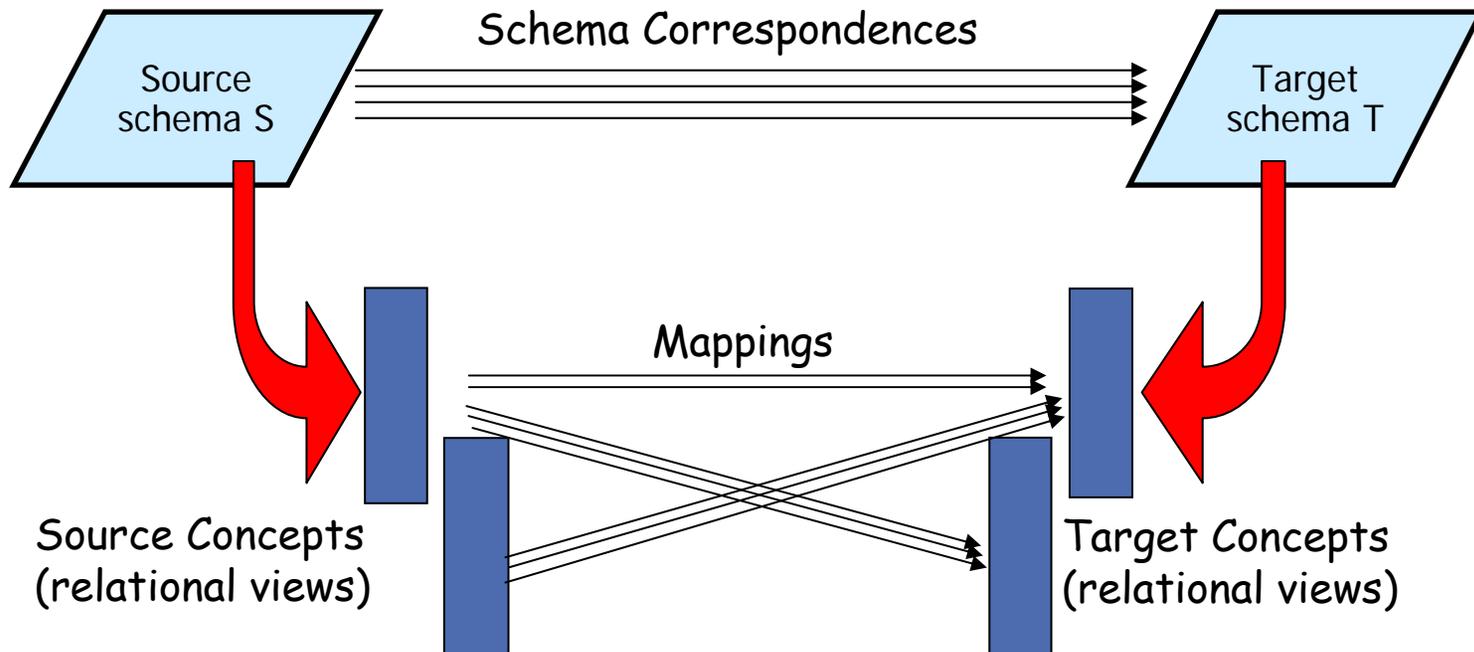


Generate Transformation Queries (XQuery)

```
<?xml version="1.0" encoding="UTF-8"?>
<statisticsDB>
  <cityStatistics>
    <city/>,
    distinct (
      FOR
        $x0 IN $doc/expenseDB/grant,
        $x1 IN $doc/expenseDB/company
      WHERE
        $x1/cid/text() = $x0/cid/text()
      RETURN
        <organization>
          <orgid> $x0/cid/text() </orgid>,
          <oname> $x1/cname/text() </oname>,
          distinct (
            FOR
              $x0L1 IN $doc/expenseDB/grant,
              $x1L1 IN $doc/expenseDB/company
            WHERE
              $x1L1/cid/text() = $x0L1/cid/text() AND
              $x1/cname/text() = $x1L1/cname/text() AND
              $x0/cid/text() = $x0L1/cid/text()
            RETURN
              <funding>
                <fid> "Sk35(", $x0L1/amount/text(), ", ", $x1L1/cname/text(), ", ", $x0L1/cid/text(), ")" </fid>,
                <proj> "Sk36(", $x0L1/amount/text(), ", ", $x1L1/cname/text(), ", ", $x0L1/cid/text(), ")" </proj>,
                <aid> "Sk32(", $x0L1/amount/text(), ", ", $x1L1/cname/text(), ", ", $x0L1/cid/text(), ")" </aid>
              </funding> )
          </organization> ),
    distinct (
      FOR
        $x0 IN $doc/expenseDB/grant,
        $x1 IN $doc/expenseDB/company
      WHERE
        $x1/cid/text() = $x0/cid/text()
      RETURN
        <financial>
          <aid> "Sk32(", $x0/amount/text(), ", ", $x1/cname/text(), ", ", $x0/cid/text(), ")" </aid>,
          <amount> $x0/amount/text() </amount>
        </financial> )
    </cityStatistics>
  </statisticsDB>
```



(Flat) Mapping Generation



- **Step 1.** Extraction of "concepts" (in each schema).
 - Concept = one category of data that can exist in the schema
- **Step 2.** Mapping generation
 - Enumerate all non-redundant maps between pairs of concepts
- [Popa, Velegakis, Miller, Hernandez, Fagin. VLDB02]
- [Fagin, Kolatios, Miller, Popa. ICDT 03]
- [Haas, Hernandez, Ho, Popa, Roth. SIGMOD 05]



(Flat) Mapping Example

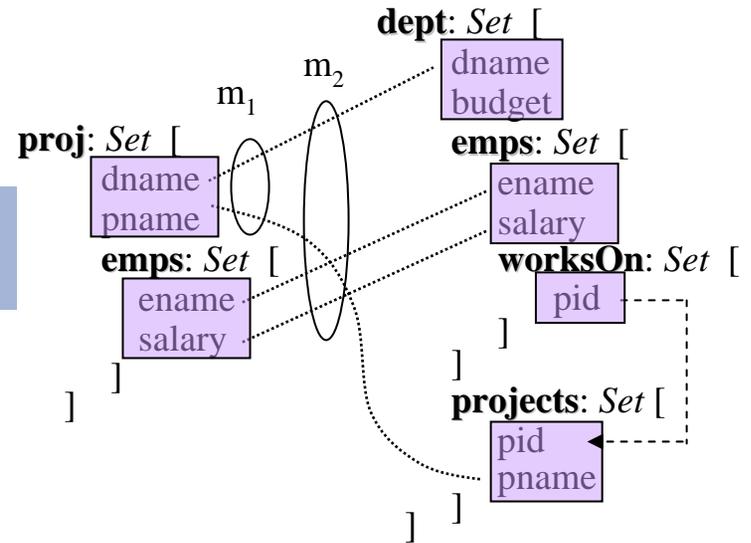
- m_1 maps **proj** to **dept-projects**

m_1 :
 $\forall(p_0 \text{ in } \text{proj})$
 $\exists(d \text{ in } \text{dept}) \exists(p \text{ in } d.\text{projects})$
 $p_0.\text{dname} = d.\text{dname}$
 $\wedge p_0.\text{pname} = p.\text{pname}$

The concept of
 “project of a
 department”

m_2 :
 $\forall(p_0 \text{ in } \text{proj}) \forall(e_0 \text{ in } p_0.\text{emps})$
 $\exists(d \text{ in } \text{dept}) \exists(p \text{ in } d.\text{projects})$
 $\exists(e \text{ in } d.\text{emps}) \exists(w \text{ in } e.\text{worksOn})$
 $w.\text{pid} = p.\text{pid}$
 $\wedge p_0.\text{dname} = d.\text{dname}$
 $\wedge p_0.\text{pname} = p.\text{pname}$
 $\wedge e_0.\text{ename} = e.\text{ename}$
 $\wedge e_0.\text{salary} = e.\text{salary}$

expression for
dept-emps-worksOn-projects



- m_2 maps **proj-emps** to **dept-emps-worksOn-projects**

The concept of
 “project of an
 employee of a
 department”

- Two 'basic' mappings (or *source-to-target tgds* or *GLAV formulas*)



IBM Rational Data Architect Product

The screenshot displays the IBM Rational Data Architect interface. The main workspace is divided into three panes: Source, Mappings, and Target. The Source pane shows a database schema for 'AMALGAM.dbm' with tables like ARTICLE, AUTHOR, and BOOK. The Target pane shows an XML schema with elements like article, author, and book. Yellow arrows indicate mappings between source columns and target elements. The Properties pane at the bottom shows a 'Discovered Mapping' for the source path /AMALGAM/S1/ARTICLE/VOL (INTEGER) to a target element.

Source	Location	Data type
/AMALGAM/S1/ARTICLE/VOL		INTEGER

Targets	Location	Data type
---------	----------	-----------



IBM Rational Data Architect Product

- Schema Matching, Schema Mapping and Query Generation Technologies from Clio
- Value Correspondences in the GUI
 - Blue Lines: Confirmed by the users
 - Gray Lines: Suggested by the schema matching algorithms
- Schema Matching
 - Five different algorithms: two **name-based** (including **thesaurus lookup**) and three **instance-based**
 - Users can choose
 - Any weighted combination of the 5 schema matching algorithms
 - Source or target
 - One element (element/attribute or column) or a group of elements (subtree or table)
 - The value k (for the top-k matches)
 - The system returns the top-k matches for each element
- Current Release
 - Source is relational (other IBM products support XML sources)
 - Target can be relational (generates SQL) or XML (generates SQL/XML)
 - Mapping is standardized within IBM, as an EMF in-memory object and as a serialized XML document



Outline

- Clio: Basic Features of a Schema Mapping System
- MAUI: Advanced Features of a Schema Mapping System
 - Nested Mapping Model
 - [Fuxman, Hernandez, Ho, Miller, Papotti, Popa. VLDB 06]
 - Mapping-Based XML Transformation Engine
 - Schema Integration
 - Schema Evolution
- Clio2010: Mapping-Based Authoring of Data Flows
- Conclusions and Future Directions



New Nested-Mapping Engine for Clio

- Existing Clio engine is based on a flat mapping model
 - Pros: easier to implement
 - Cons:
 - Fragmentation into many overlapping mappings
 - Inefficiency in execution
 - Redundancy in the output data
 - No user-defined grouping semantics
- New Clio engine is based on a nested-mapping model
 - Cons: more challenging to design and implement
 - Pros: overcomes the above problems in the flat mapping model



(Flat) Mapping Example

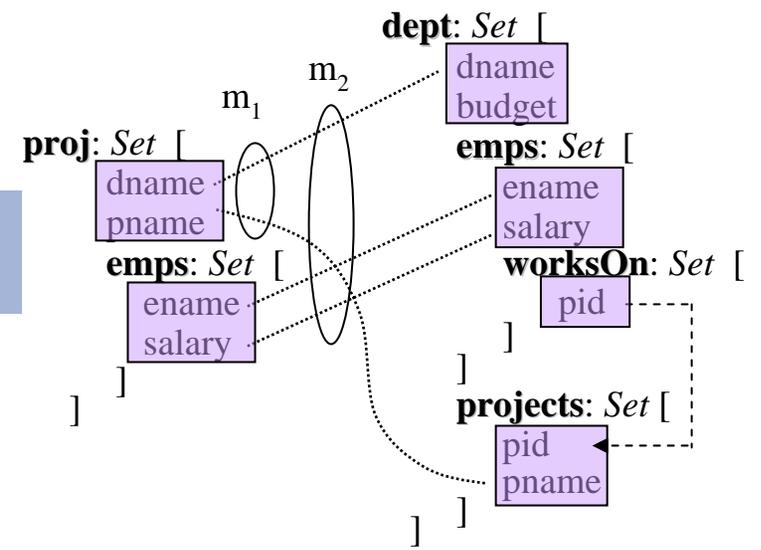
- m_1 maps **proj** to **dept-projects**

$$\begin{aligned}
 m_1: & \\
 & \forall (p_0 \text{ in } \text{proj}) \\
 & \quad \exists (d \text{ in } \text{dept}) \exists (p \text{ in } d.\text{projects}) \\
 & \quad \quad p_0.\text{dname} = d.\text{dname} \\
 & \quad \quad \wedge p_0.\text{pname} = p.\text{pname}
 \end{aligned}$$

The concept of
“project of a
department”

$$\begin{aligned}
 m_2: & \\
 & \forall (p_0 \text{ in } \text{proj}) \forall (e_0 \text{ in } p_0.\text{emps}) \\
 & \quad \exists (d \text{ in } \text{dept}) \exists (p \text{ in } d.\text{projects}) \\
 & \quad \exists (e \text{ in } d.\text{emps}) \exists (w \text{ in } e.\text{worksOn}) \\
 & \quad \quad w.\text{pid} = p.\text{pid} \\
 & \quad \quad \wedge p_0.\text{dname} = d.\text{dname} \\
 & \quad \quad \wedge p_0.\text{pname} = p.\text{pname} \\
 & \quad \quad \wedge e_0.\text{ename} = e.\text{ename} \\
 & \quad \quad \wedge e_0.\text{salary} = e.\text{salary}
 \end{aligned}$$

expression for
dept-emps-worksOn-projects



- m_2 maps **proj-emps** to **dept-emps-worksOn-projects**

The concept of
“project of an
employee of a
department”

- Two 'basic' mappings (or *source-to-target tgds* or *GLAV formulas*)



Correlating Mapping Formulas

$m_1: \forall(p_0 \text{ in proj})$
 $\exists(d \text{ in dept}) \exists(p \text{ in d.projects})$
 $p_0.dname = d.dname \wedge p_0.pname = p.pname$

$m_2: \forall(p_0 \text{ in proj}) \forall(e_0 \text{ in } p_0.emps)$
 $\exists(d \text{ in dept}) \exists(p \text{ in d.projects}) \exists(e \text{ in d.emps}) \exists(w \text{ in e.worksOn})$
 $w.pid = p.pid$
 $\wedge p_0.dname = d.dname \wedge p_0.pname = p.pname$
 $\wedge e_0.ename = e.ename \wedge e_0.salary = e.salary$

proj tuples mapped only once

Submapping, correlated to the parent mapping

For every proj tuple, we map all employees, as a **group**.
 (Source grouping is preserved)

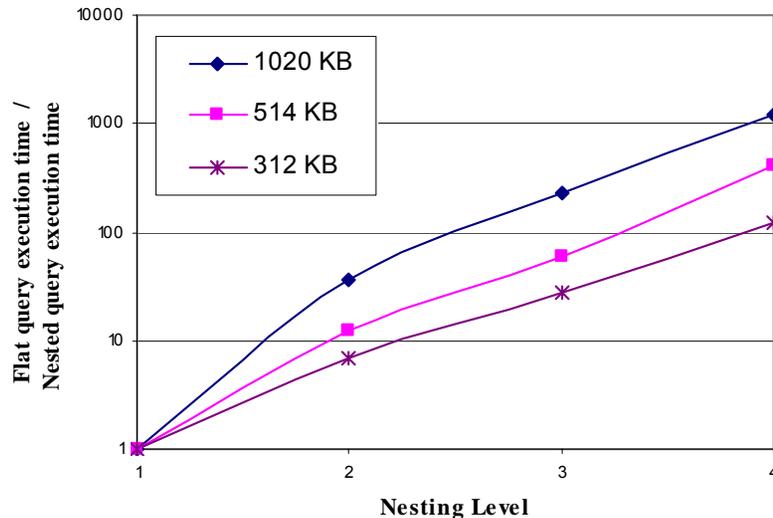
Replace with

$n: \forall(p_0 \text{ in proj})$
 $\exists(d \text{ in dept}) \exists(p \text{ in d.projects})$
 $p_0.dname = d.dname \wedge p_0.pname = p.pname$
 $\wedge [\forall(e_0 \text{ in } p_0.emps)$
 $\exists(e \text{ in d.emps}) \exists(w \text{ in e.worksOn})$
 $w.pid = p.pid$
 $\wedge e_0.ename = e.ename \wedge e_0.salary = e.salary$
 $]$

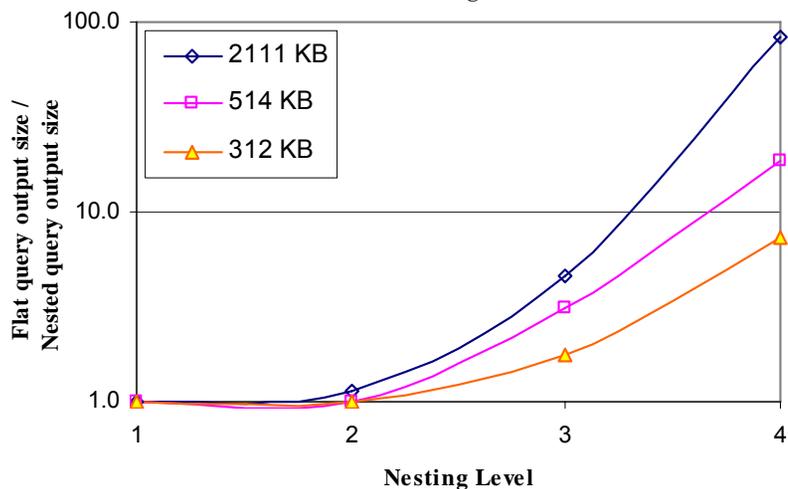
This is a **nested mapping**



Performance



Execution time for flat:
22mins
Execution time for nested: **1.1s**



- Size of generated data (flat) – including duplicates: **45MB**
- Size of generated data (nested): **552KB**

The nested mapping generates much more efficient execution and less redundant data

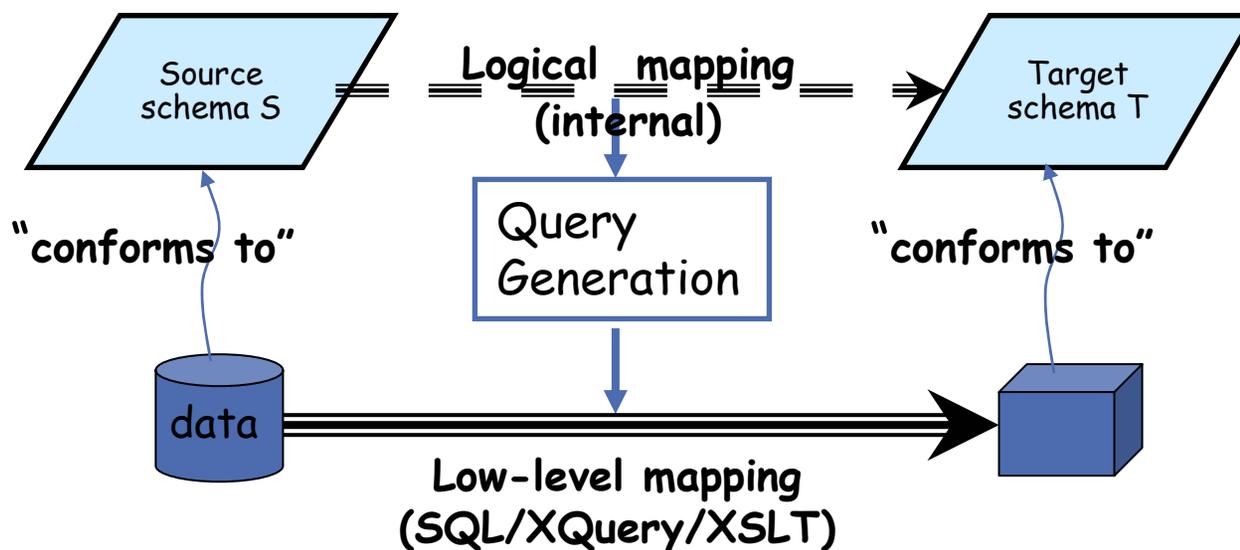


Outline

- Clio: Basic Features of a Schema Mapping System
- MAUI: Advanced Features of a Schema Mapping System
 - Nested Mapping Model
 - Mapping-Based XML Transformation Engine
 - [Jiang, Ho, Popa, Han. WWW 07]
 - Schema Integration
 - Schema Evolution
- Clio2010: Mapping-Based Authoring of Data Flows
- Conclusions and Future Directions



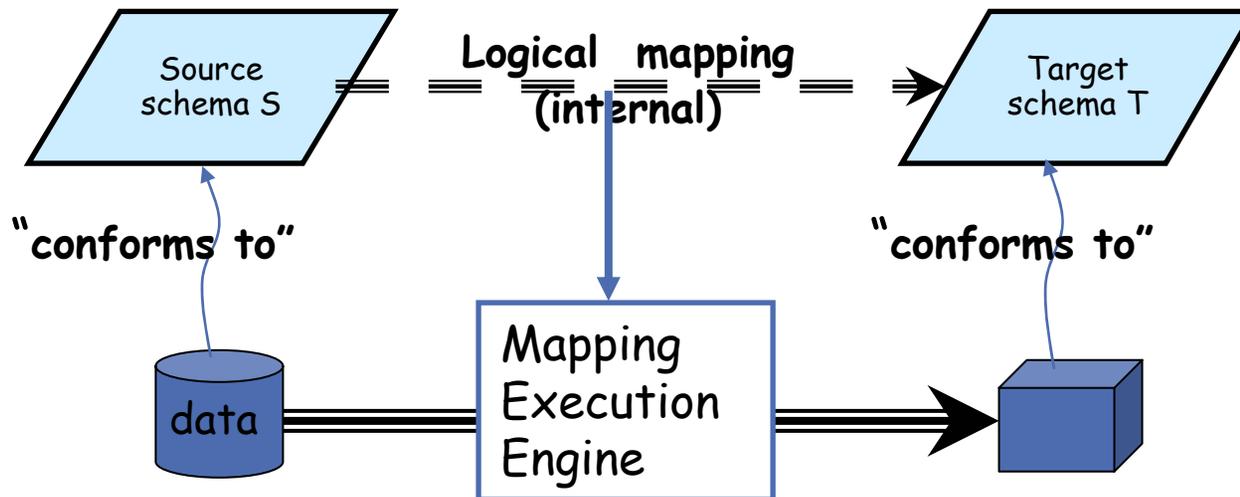
Performance in Executing Mappings



- Mappings are translated into general purpose query languages
 - E.g., SQL, XQuery, XSLT
- Query optimization issues are left to the runtime engine of each of these languages to decide
 - i.e., Q.O. decisions are not encoded in the queries.
- Idea: Execute mappings directly in our own runtime.



Mapping Execution Engine in Java



- Motivation
 - Grouping (on multiple levels) of transformed values.
 - XQuery & XSLT have no specific constructs (hence algorithms) for such grouping tasks
- Scalability and efficiency
 - Controllable memory resource usage
 - Speed of transformation almost linear to input sizes
- Mapping-based
 - Model high-level mapping semantics using IBM Mapping Specification Language (MSL) standard

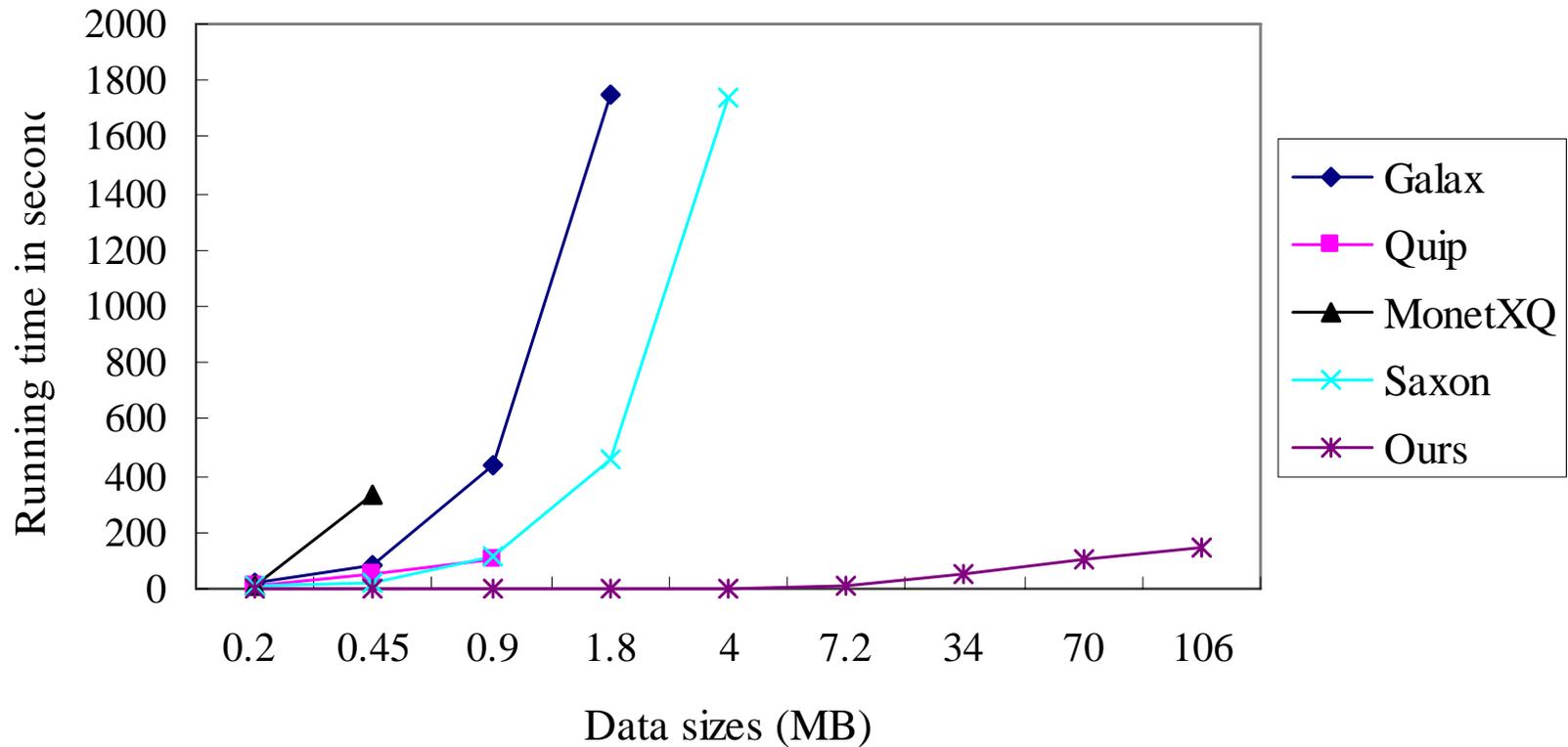


Three-Phase Algorithm for Executing Mappings

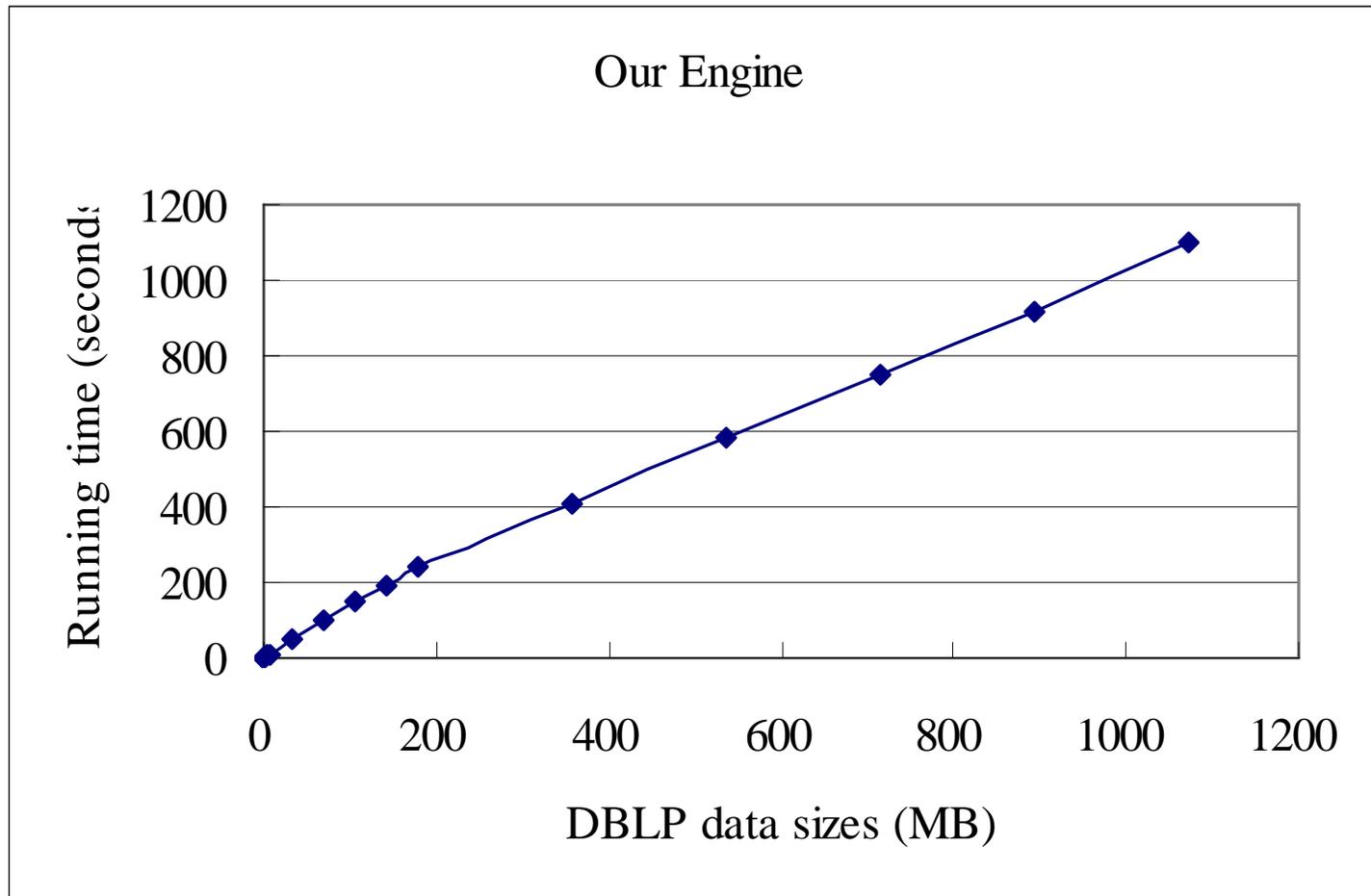
- Phase 1 (Extract): tuple extraction
 - Streaming holistic twig join
 - Streaming and stored/indexed XML
 - SQL queries through ODBC
 - Relational source
- Phase 2 (Transform): generate an XML tree from each tuple
- Phase 3 (Merge): data merging
 - A dynamic, scalable merge algorithm
 - Hash-based vs. sort-based algorithms



XML-to-XML Comparative Results



XML-to-XML Scalability Results



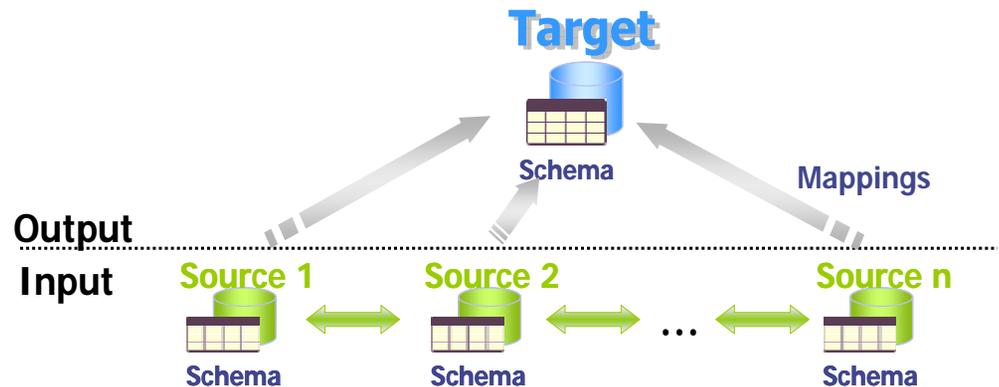
Outline

- Clio: Basic Features of a Schema Mapping System
- MAUI: Advanced Features of a Schema Mapping System
 - Nested Mapping Model
 - Mapping-Based XML Transformation Engine
 - Schema Integration
 - [Chiticariu, Hernandez, Popa, Kolaitis. VLDB 07 demo]
 - Schema Evolution
- Clio2010: Mapping-Based Authoring of Data Flows
- Conclusions and Future Directions

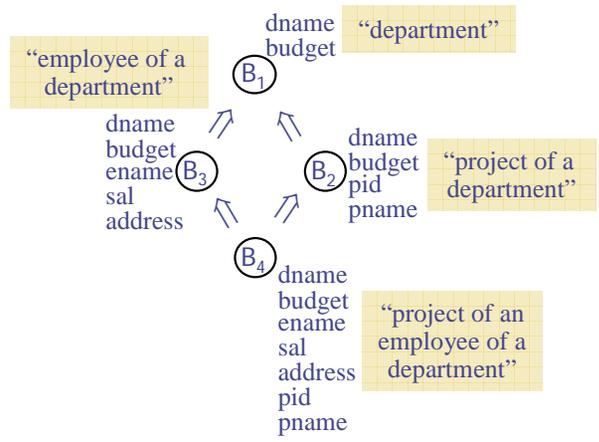
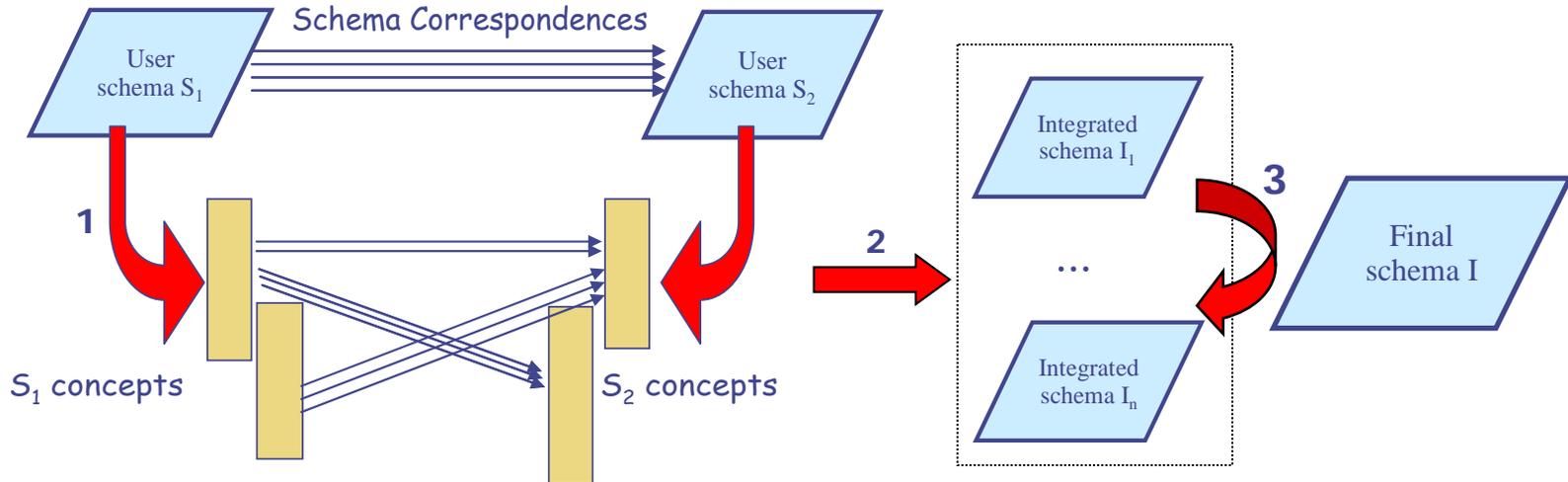


Schema Integration

- Problem
 - Given multiple overlapping schemas in the same domain, consolidate them into one.
- Applications
 - Provide a standard representation of the data (for unified querying, data warehousing, reference point, etc.)
 - “Metadata chaos reduction”
- Schema integration is a hard problem
 - Requires a lot of human interaction/feedback
 - A series of “recipes” that arrive at a single integrated schema
- Related work
 - [Pottinger, Bernstein. VLDB 03]
 - [Chiticariu, Hernandez, Popa, Kolaitis. VLDB 07 demo]



Schema Integration [Chiticariu et al.]



- **Step 1.** Extraction of "concepts" from a hierarchy
- **Step 2.** Consider all possible ways of performing a hierarchical merge of the related concepts
 - Generate initial set of candidate "good" integrated schemas
 - **Duplication-free enumeration algorithm**
 - Avoids duplicates by using constraints (Horn clauses)
 - *Polynomial-delay algorithm* for enumerating satisfying assignments
- **Step 3.** Browse/Search/Refine set of candidate schemas
 - **Combination of partial enumeration with user constraints**
 - Based on the schemas seen so far, users express constraints on the "future" schemas to be generated

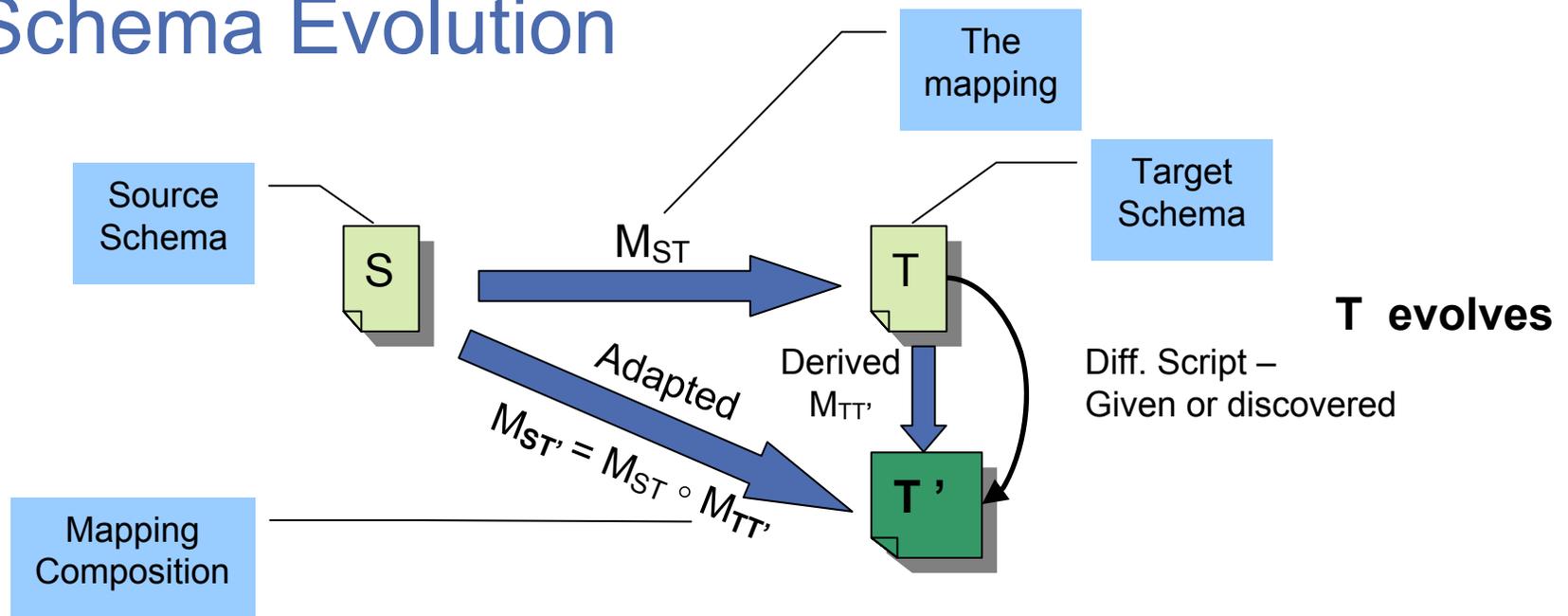


Outline

- Clio: Basic Features of a Schema Mapping System
- MAUI: Advanced Features of a Schema Mapping System
 - Nested Mapping Model
 - Mapping-Based XML Transformation Engine
 - Schema Integration
 - Schema Evolution
- Clio2010: Mapping-Based Authoring of Data Flows
- Conclusions and Future Directions



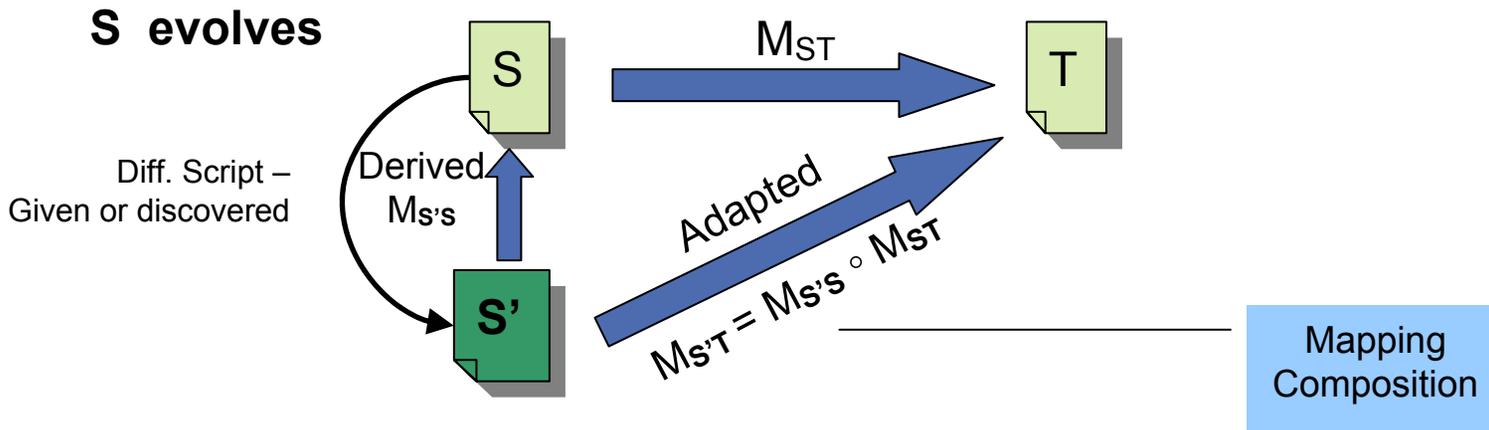
Schema Evolution



- When the target schema evolves, the system can generate an initial **default mapping** for $T \rightarrow T'$ automatically
- The user adds “new mapping lines” for new schema elements in T' to complete the mapping $M_{TT'}$ (for $T \rightarrow T'$)
- The new mapping $M_{ST'}$ ($S \rightarrow T'$) is a composition of M_{ST} ($S \rightarrow T$) and $M_{TT'}$ ($T \rightarrow T'$)



Schema Evolution



- When the source schema evolves, the system can generate an initial default mapping for $S' \rightarrow S$ automatically
- The user adds “new mapping lines” for new schema elements in S' to complete the mapping $M_{S'S}$ (for $S' \rightarrow S$)
- The new mapping $M_{S'T}$ ($S' \rightarrow T$) is a composition of $M_{S'S}$ ($S' \rightarrow S$) and M_{ST} ($S \rightarrow T$)



Related Theory and Algorithms

- Mapping Composition
 - [Madhavan, Halevy. VLDB 03]
 - [Fagin, Kolaitis, Popa, Tan. PODS 04]
 - [Nash, Bernstein, Melnik. PODS 05]
 - [Bernstein, Green, Melnik, Nash. VLDB 06]
- Mapping Inversion
 - [Fagin. PODS 06]
 - [Fagin, Kolaitis, Popa, Tan. PODS 07]
- Schema Evolution
 - [Rahm, Bernstein. SIGMOD Rec. Dec 06]
- Mapping Adaptation under Evolving Schemas
 - [Velegrakis, Miller, Popa. VLDB 03]
 - [Yu, Popa. VLDB 05]
- Query rewrite
 - [Yu, Popa. SIGMOD 04]



Outline

- Clio: Basic Features of a Schema Mapping System
- MAUI: Advanced Features of a Schema Mapping System
- ▶ ■ Clio2010: Mapping-Based Authoring of Data Flows
 - **ETL**: “Mapping ↔ ETL” Conversion
 - **Web-Service Composition**: Mapping for web-service data sources
 - **Mashups**: “Mapping → Mashup” Generation
 - **Reuse**: Mapping Polymorphism
- Conclusions and Future Directions



Mapping-Based Authoring of Data Flows

■ Motivation

- Schema mappings are the **building blocks** for larger data transformation and integration applications
- The **graph of mappings** is a declarative specification of the **flow of data**
 - Similar to ETL & mashups
 - Need to take functional (web-service) data sources
- Need to extend Clio where multiple mappings can be **defined, loaded** (sharing a context), **reused** and **managed**



Clio2010

■ Clio2010

- Automatically assemble a graph of initial, uncorrelated mappings into larger, richer mappings
- Support more complicated
 - **mapping composition** and
 - **mapping merge**
- Support functional data sources (e.g., **web services**)
- Compile the larger mappings into a global execution plan in
 - Queries (XQuery)
 - transformation scripts (XSLT)
 - ETL flows (e.g., IBM WebSphere DataStage)
 - mashups (e.g., IBM DAMIA)
- Support mapping reuse through **mapping polymorphism** framework



Outline

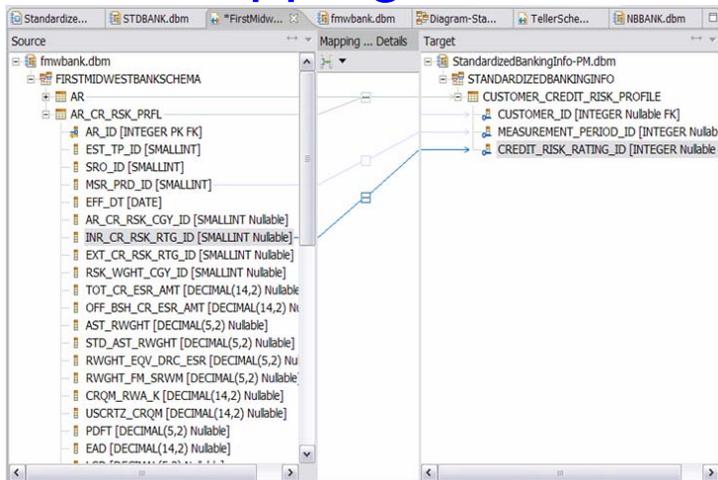
- Clio: Basic Features of a Schema Mapping System
- MAUI: Advanced Features of a Schema Mapping System
- Clio2010: Mapping-Based Authoring of Data Flows
 - **ETL**: “Mapping \leftrightarrow ETL” Conversion
 - [Hernandez et al. 07]
 - **Web-Service Composition**: Mapping for web-service data sources
 - **Mashups**: “Mapping \rightarrow Mashup” Generation
 - **Reuse**: Mapping Polymorphism
- Conclusions and Future Directions



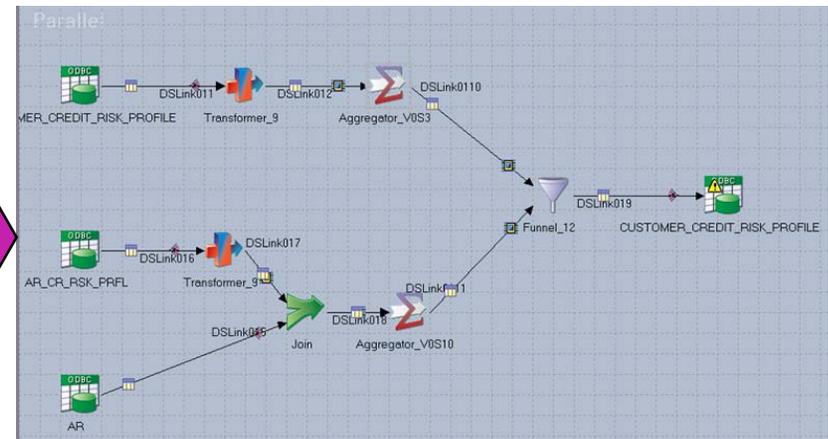
Orchid: “Mapping ↔ ETL” Conversion

- Round-tripping between ETL (Extract, Transform, Load) scripts and declarative mappings
 - Generate ETL script from mappings
 - Extract mapping information from ETL scripts
 - Track and propagate mapping-related modifications
- Optimization of ETL scripts
 - Remove redundant operations
 - Push computation to other runtimes
 - “Rewrite” ETL scripts

EII Mapping in RDA



ETL Flow



Mapping → ETL Generation

The screenshot displays the IBM Rational Software Development Platform interface. The main window is titled "Data - FirstMidwestToStandard.msl". The "Data Project Explorer" on the left shows a project structure with various data sources and targets. The "Mapping ... Details" pane in the center shows a mapping between source and target tables. The source table is "AR_CR_RSK_PRFL" and the target table is "CUSTOMER_CREDIT_RISK_PROFILE". The mapping shows the following connections:

- AR_ID [INTEGER PK FK] maps to CUSTOMER_ID [INTEGER Nullable FK]
- EST_TP_ID [SMALLINT] maps to MEASUREMENT_PERIOD_ID [INTEGER Nullable]
- SRO_ID [SMALLINT] maps to CREDIT_RISK_RATING_ID [INTEGER Nullable]

A dialog box titled "Generate Ascential Datastage Model" is open in the foreground. It prompts the user to select options for generating the script. The "Create a Datastage Designer model" option is selected. The script name is "rd2datastage_model_transform.xml".

The "Properties" pane at the bottom shows the mapping details:

Location	Data type
/fmwbank/FIRSTMIDWESTBANKSCHE...	SMALLINT

Location	Data type
/STDBANK/STANDARDIZEDBANKINGI...	INTEGER

Trans...tion: AVG (INR_CR_RSK_RTG_ID)

The developer inspects the connections between the sources, rules, and data warehouse model...

...and generates a DataStage job to move data

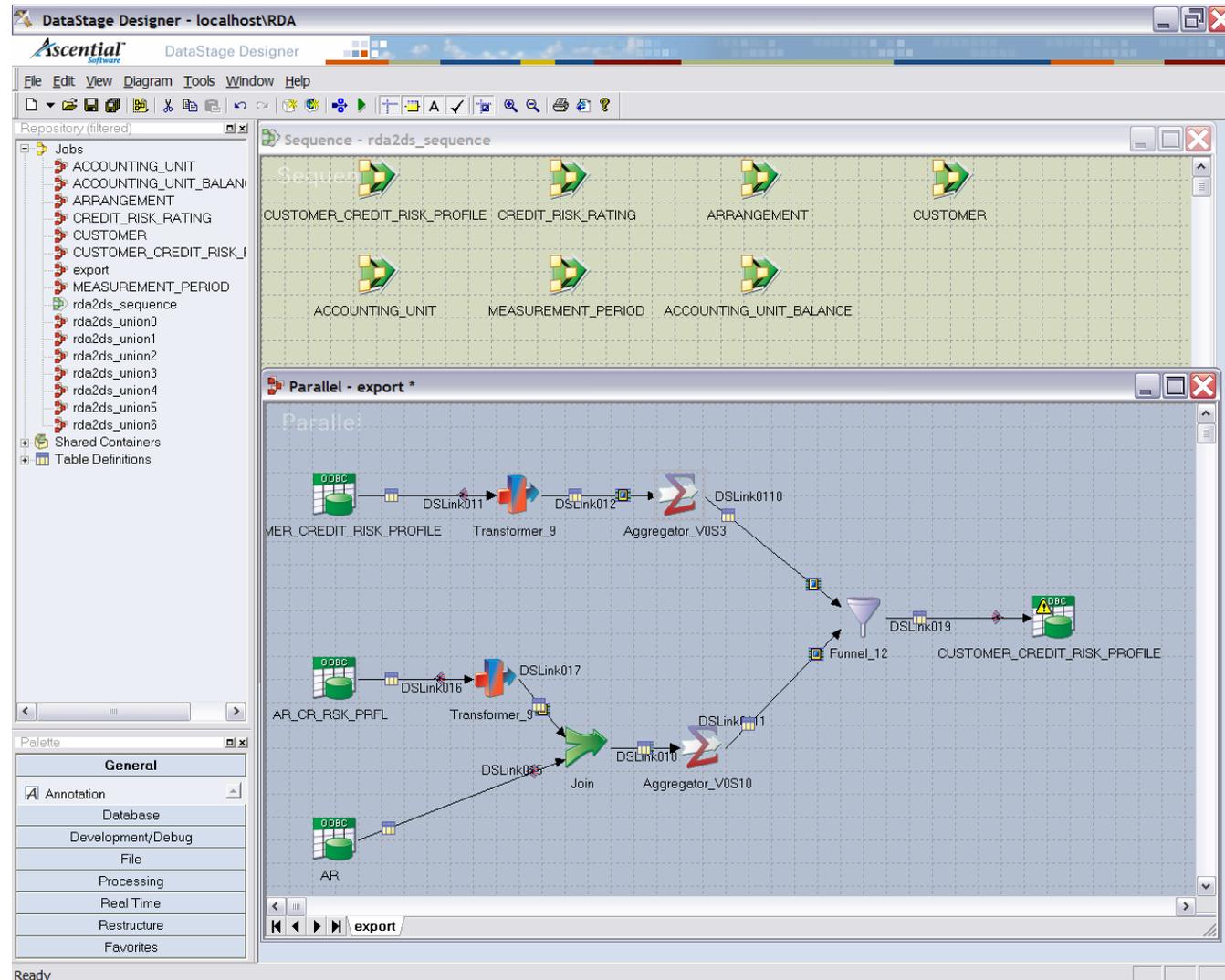


ETL → Mapping Abstraction

The developer refines, tests and deploys the DataStage job to production.

Need ETL → Mapping:

- To reflect changes made in the ETL script back in the mapping.*
- To extract mappings out of ETL scripts.*

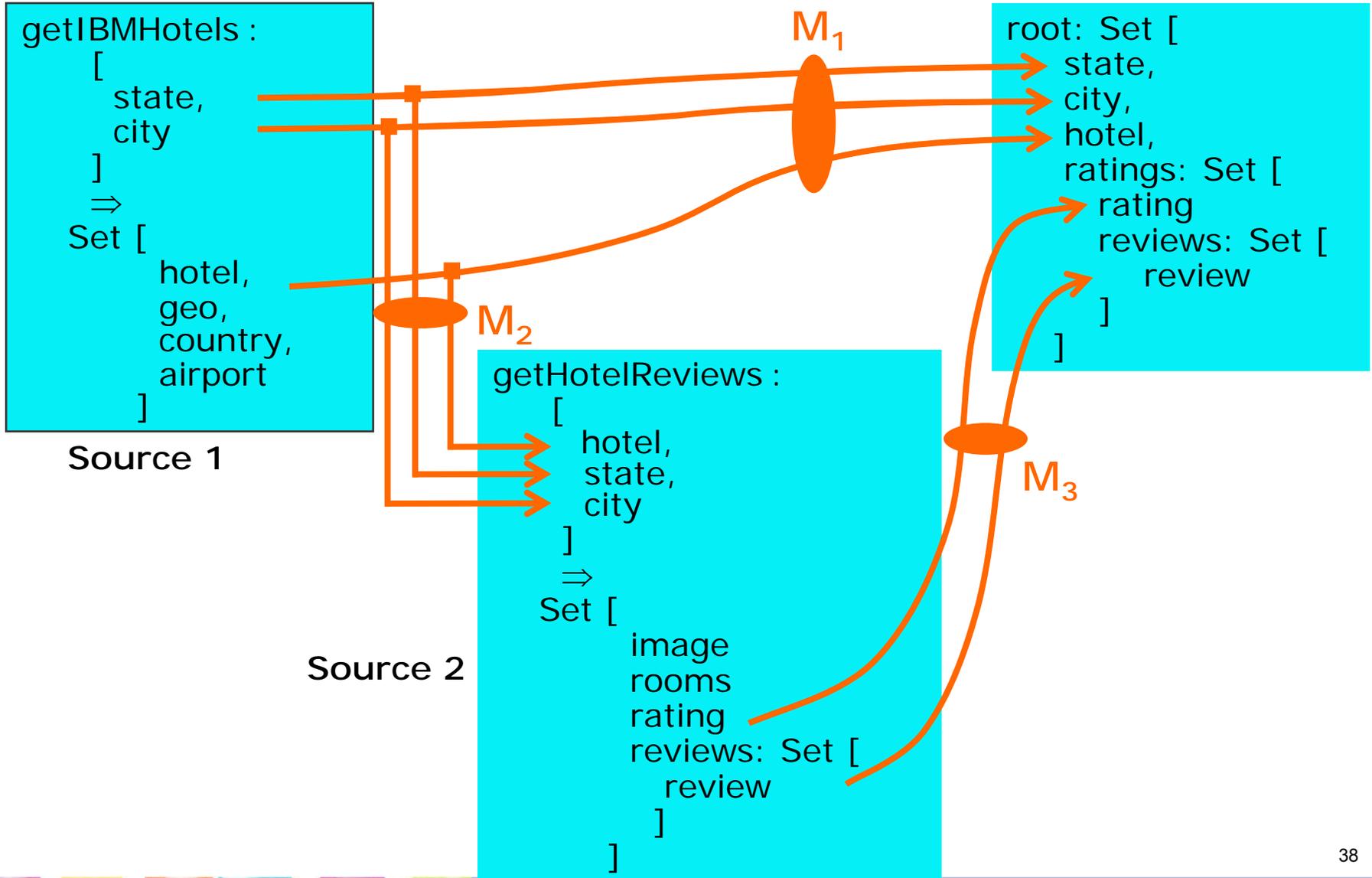


Outline

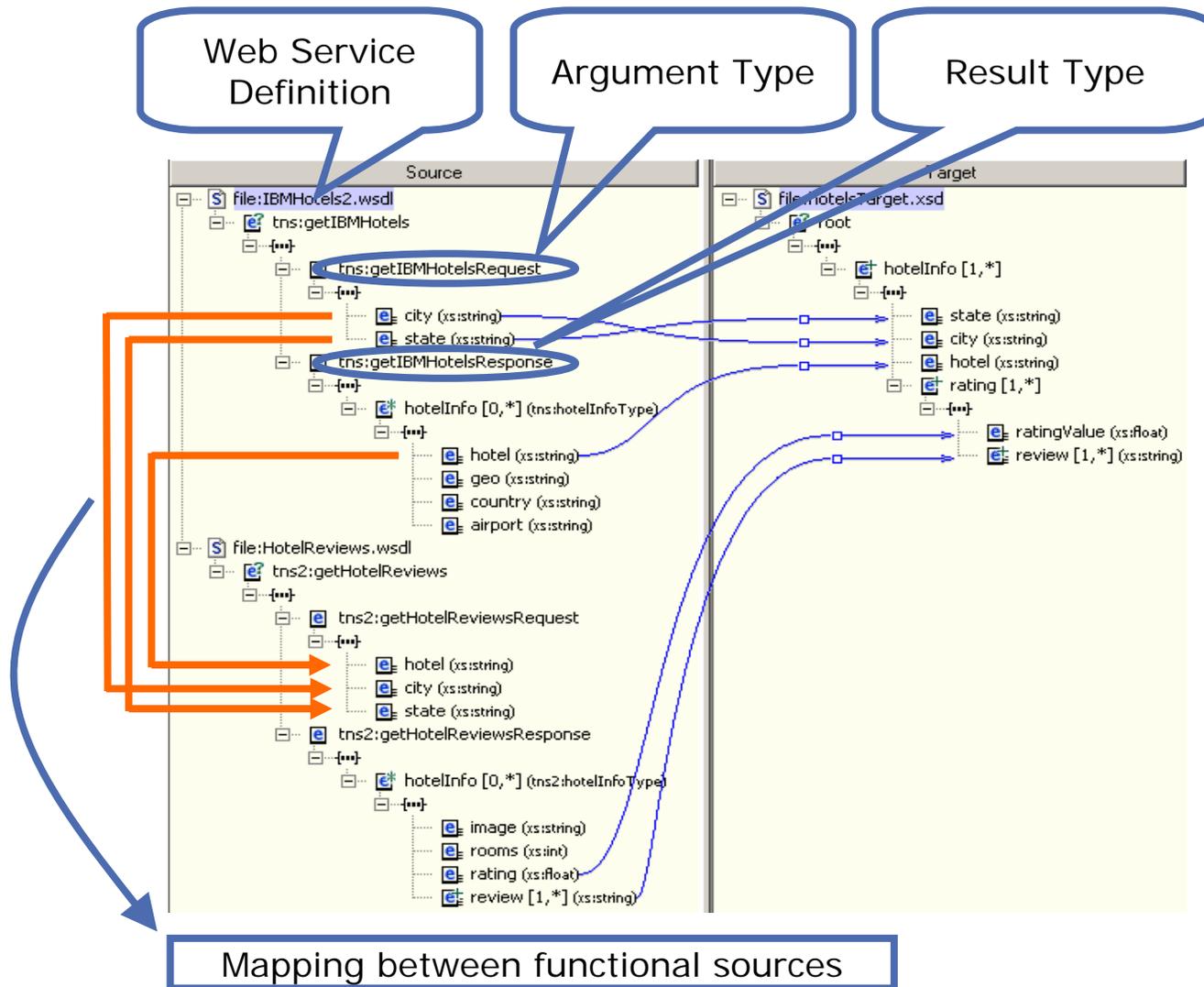
- Clio: Basic Features of a Schema Mapping System
- MAUI: Advanced Features of a Schema Mapping System
- Clio2010: Mapping-Based Authoring of Data Flows
 - ETL: “Mapping ↔ ETL” Conversion
 - Web-Service Composition: Mappings for web-service data sources
 - [Alexe et al. 07]
 - Mashups: “Mapping → Mashup” Generation
 - Reuse: Mapping Polymorphism
- Conclusions and Future Directions



Flows of Mappings



Functional Data Sources



Outline

- Clio: Basic Features of a Schema Mapping System
- MAUI: Advanced Features of a Schema Mapping System
- Clio2010: Mapping-Based Authoring of Data Flows
 - **ETL**: “Mapping ↔ ETL” Conversion
 - **Web-Service Composition**: Mapping for web-service data sources
 - **Mashups**: “Mapping → Mashup” Generation
 - [Camacho et al. 07]
 - **Reuse**: Mapping Polymorphism
- Conclusions and Future Directions



IBM DAMIA: A Mashup Tool

■ IBM DAMIA

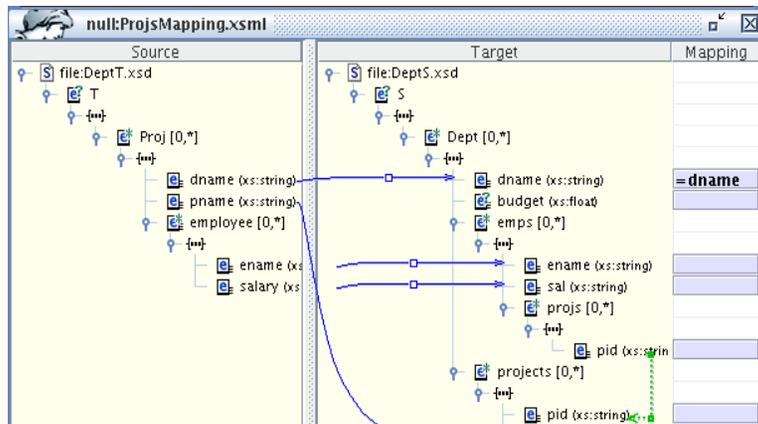
- Produces easily customizable flows
- Mashup fabric for data aggregation and transformation
- Web-based tool with user-friendly interface
- No scripting/programming skills required
- XQuery-like model

The screenshot displays the IBM DAMIA web interface within a Mozilla Firefox browser window. The title bar reads "IBM DAMIA - weather - Mozilla Firefox". The interface includes a menu bar (File, Edit, View, History, Bookmarks, Tools, Help) and a header with the "w3" logo and "IBM DAMIA" text. Below the header are tabs for "Sources", "Operators", and "Primitive Operators". A toolbar contains buttons for "Augment", "Merge", "Filter", "Sort", "Group", "Union", and "Transform". The main workspace shows a workflow diagram with components: a green "\$Import2" component, a blue "\$Transform3" component, a blue "\$Merge4" component, and a blue "\$Publish5" component. Below the diagram is a configuration panel for the "\$Import2" component. The panel has tabs for "Params" and "Preview". The "Params" tab is active, showing fields for "Name" (set to "\$Import2"), "Source URL" (set to "http://www.weather.gov/alerts/tx.rss"), "Parameters in URL Query string" (with a table for Name and Value), and "Feed Type" (set to "RSS"). A "Done" button is at the bottom. To the right of the configuration panel is a help panel with text: "The **Import** source is used to extract the repeating XML elements from an XML source." and "Parameters" section with definitions for Name, Source URL, Feed type, and Repeating element.



Generate DAMIA Flows from Clio

Clio



High level mappings

Compiler

SQL

XSLT

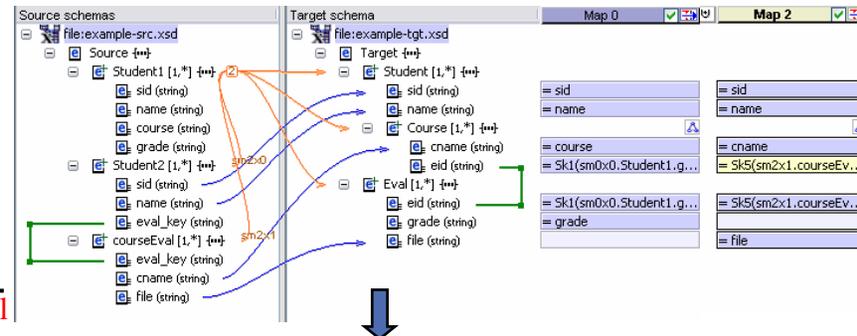
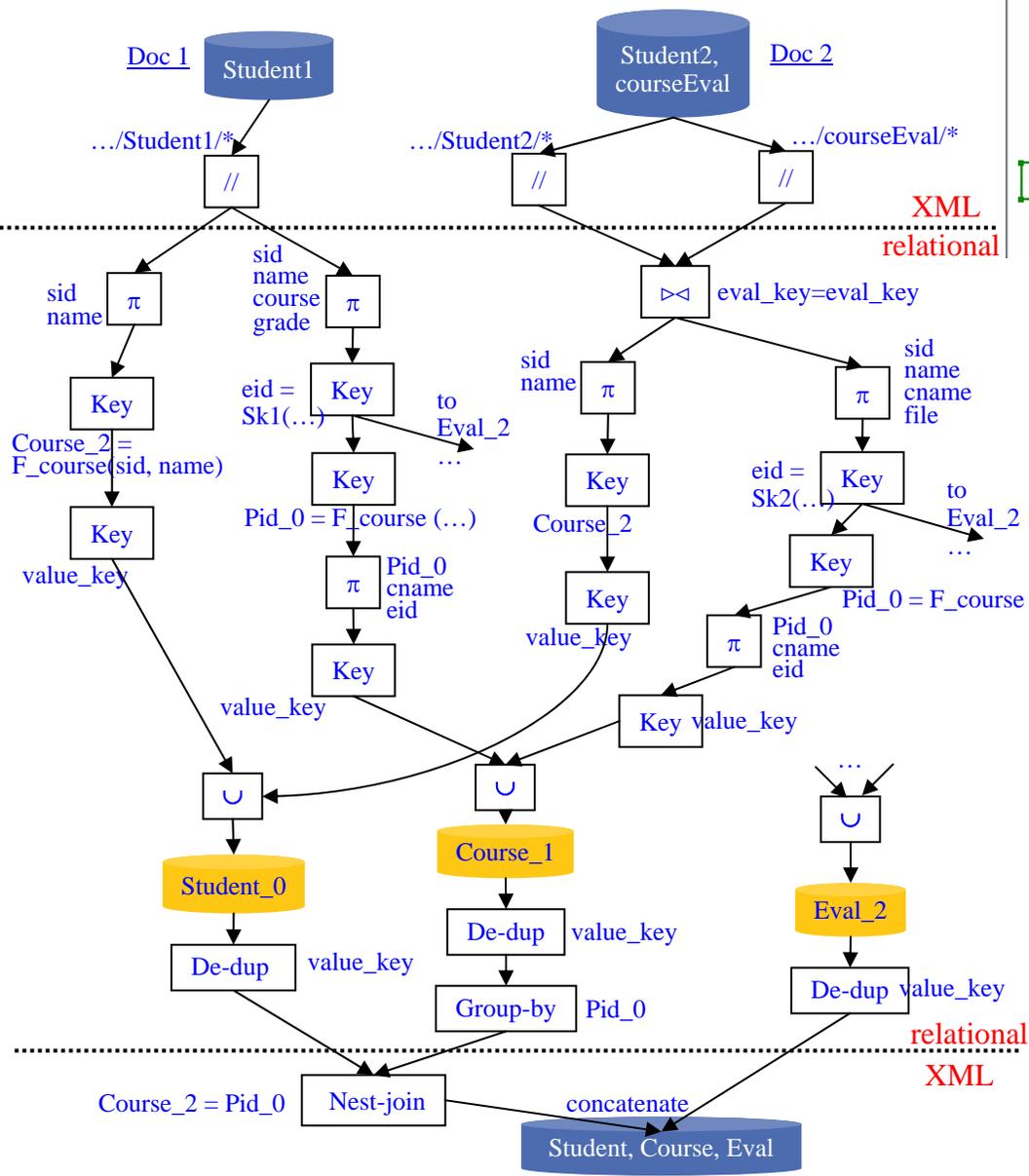
XQuery

**DAMIA
flow**

From XQuery to ETL

- The XQuery has a specific flow pattern, with several phases:
 - Source extraction queries (projection, navigation, join)
 - Result: sets of flat tuples
 - Key generation
 - Union
 - Duplicate elimination
 - Computation of target groups
 - Generate target output (hierarchy)
 - Back to hierarchical data
- We can visualize all this as a graph of operators
 - we get the equivalent of an ETL job for XML





```

(===== Phase-1: shredding into relational views =====)
declare function local:shred($doc0 as element(root)) as element(root) {
  <root>
  {
    for $sm0x0 in $doc0/Source/Student1
    return
      <Student_0>
      <sid_0>{ $sm0x0/sid/text() }</sid_0>
      <name_1>{ $sm0x0/name/text() }</name_1>
      <Course_2>{ fn:concat("F_course(", $sm0x0/sid/text(), " ", $sm0x0/name/text(), ")") }</Course_2>
      <value-key>{ fn:string-join(($sm0x0/sid/text(), $sm0x0/name/text(), fn:concat("F_course(", $sm0x0/sid/text(), " ",
      $sm0x0/name/text(), ")"), "T") ) }</value-key>
    }
  }
  for $sm2x0 in $doc0/Source/Student2,
  $sm2x1 in $doc0/Source/courseEval
  where $sm2x0/eval_key = $sm2x1/eval_key
  return
    <Student_0>
    <sid_0>{ $sm2x0/sid/text() }</sid_0>
    <name_1>{ $sm2x0/name/text() }</name_1>
    <Course_2>{ fn:concat("F_course(", $sm2x0/sid/text(), " ", $sm2x0/name/text(), ")") }</Course_2>
    <value-key>{ fn:string-join(($sm2x0/sid/text(), $sm2x0/name/text(), fn:concat("F_course(", $sm2x0/sid/text(), " ",
    $sm2x0/name/text(), ")"), "T") ) }</value-key>
  }
  for $sm0x0 in $doc0/Source/Student1
  return
    <Course_1>
    <Pid_0>{ fn:concat("F_course(", $sm0x0/sid/text(), " ", $sm0x0/name/text(), ")") }</Pid_0>
    <cname_1>{ $sm0x0/course/text() }</cname_1>
    <eid_2>{ fn:concat("Sk0(", $sm0x0/grade/text(), " ", $sm0x0/course/text(), " ", $sm0x0/name/text(), " ", $sm0x0/sid/
    text(), ")") }</eid_2>
    <value-key>{ fn:string-join((fn:concat("F_course(", $sm0x0/sid/text(), " ", $sm0x0/name/text(), ")"), $sm0x0/course/
    text(), fn:concat("Sk0(", $sm0x0/grade/text(), " ", $sm0x0/course/text(), " ", $sm0x0/name/text(), " ", $sm0x0/sid/text(),
    ")"), "T") ) }</value-key>
  }
}

(===== main query =====)
let $p0 := local:ldgen(),
    $sp1 := local:dupelim(local:shred($p0)),
    $sp2 := local:nest($sp1)
return $sp2

```



Outline

- Clio: Basic Features of a Schema Mapping System
- MAUI: Advanced Features of a Schema Mapping System
- Clio2010: Mapping-Based Authoring of Data Flows
 - **ETL**: “Mapping ↔ ETL” Conversion
 - **Web-Service Composition**: Mapping for web-service data sources
 - **Mashups**: “Mapping → Mashup” Generation
 - **Reuse**: Mapping Polymorphism
 - [Wisnesky et al. 07]
- Conclusions and Future Directions

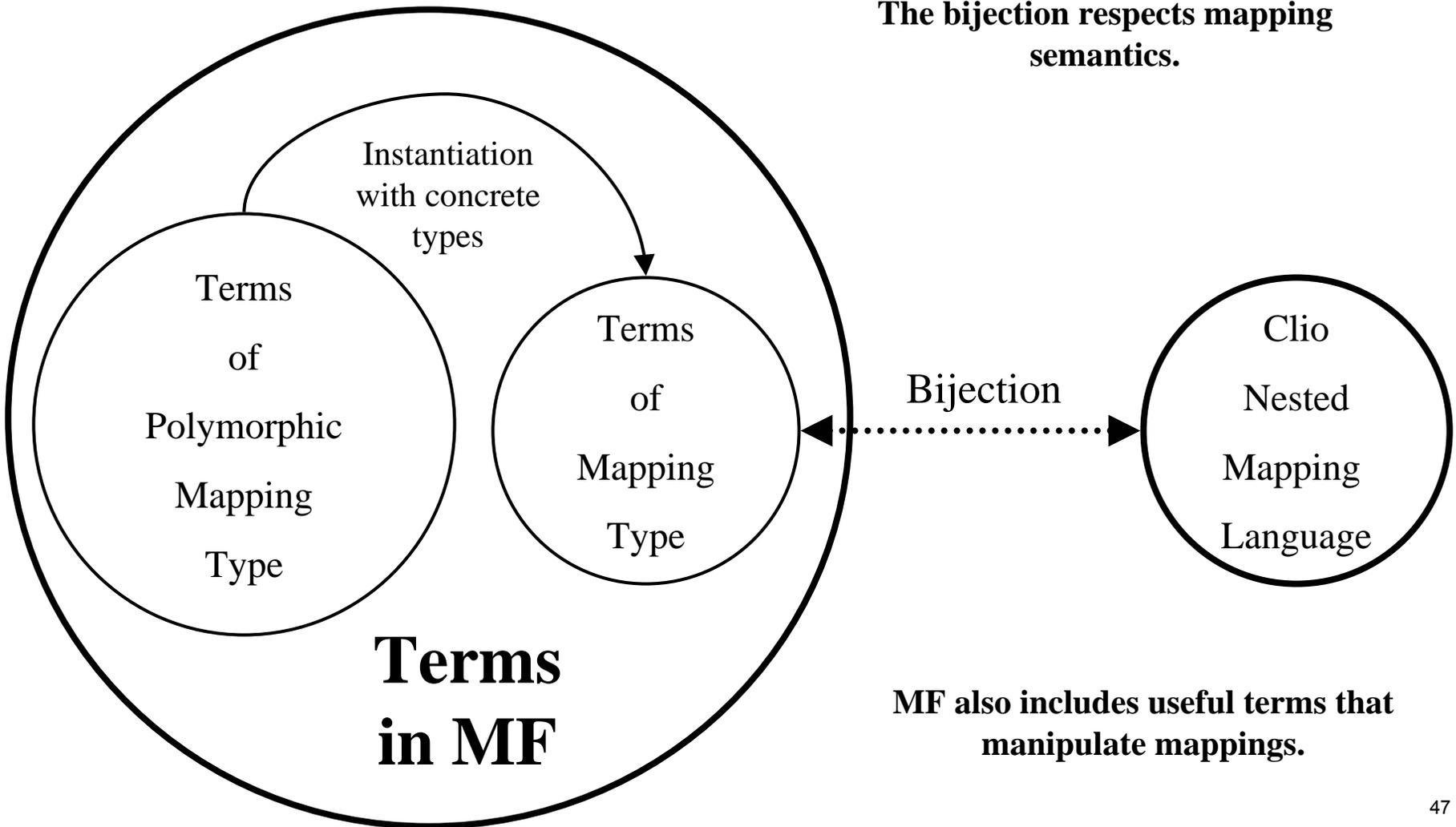


Mapping Reuse

- Motivation:
 - Existing mapping formalisms *implicitly* contain information that we need to be *explicit*
- Example:
 - A Clio nested mapping expression *implicitly* defines a class of schemas for which it has an interpretation
 - But in Clio mapping representation (XSML), every mapping instance requires *concrete* source and target *schemas*
 - The flexibility to *re-use* mappings is lost
- Approach:
 - We need a *formal language* for talking about *mappings*
 - Require theory and tools to manipulate *mappings* as *blocks*
- Solution:
 - Mapping Polymorphism



Clio Language Embedding into MF

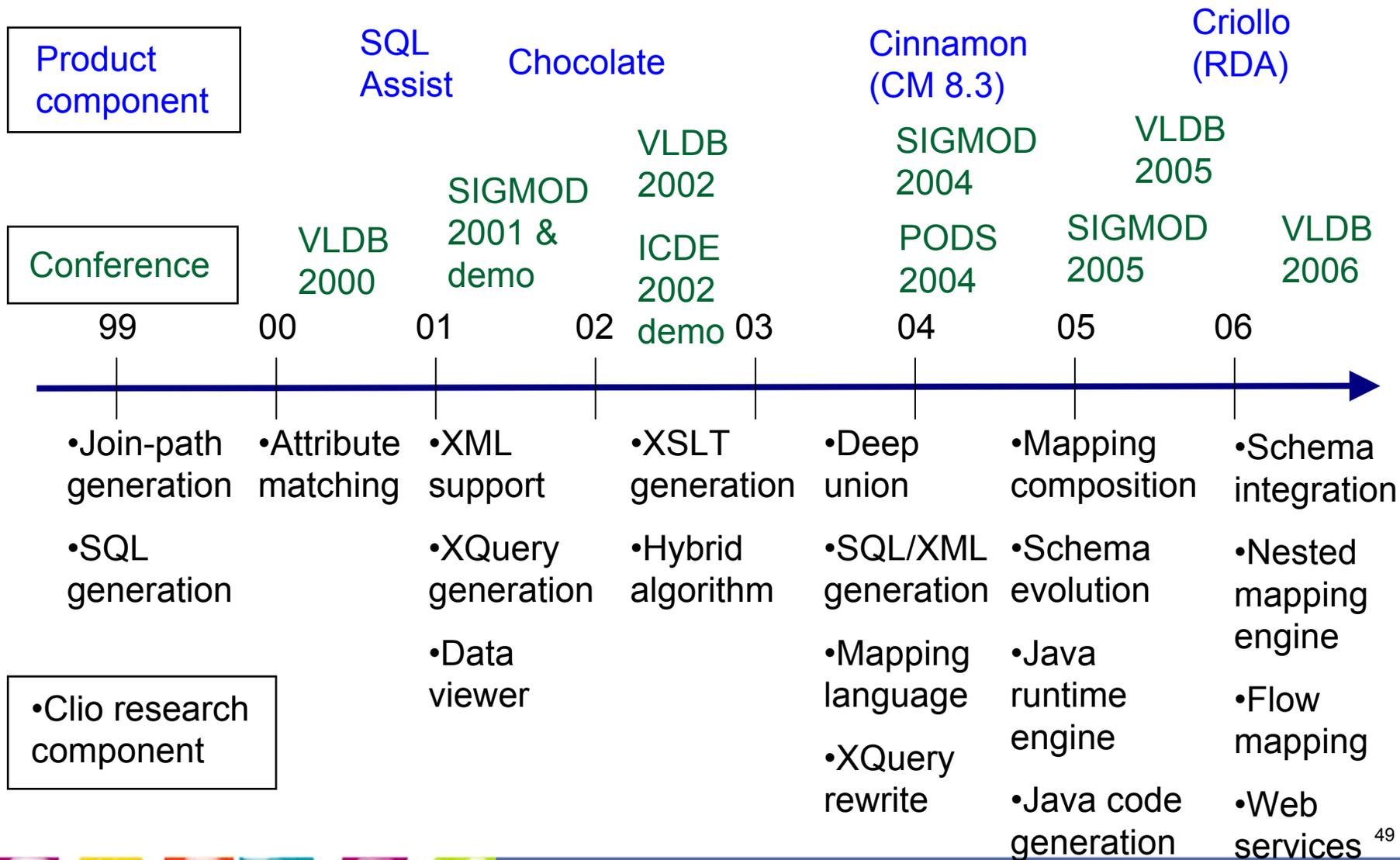


Outline

- Clio: Basic Features of a Schema Mapping System
- MAUI: Advanced Features of a Schema Mapping System
- Clio2010: Mapping-Based Authoring of Data Flows
- ▶ ■ Conclusion and Related Work



Clio Innovations Over Time



Conclusion

- Mappings define relationships between schemas of data sources
- Mappings and transformation queries are the “new” metadata
- Customers want “I” (for integration), not EAI, EII, ETL, etc.
 - Consumability for information integration



Related Work: Very Small Subset

- Information Integration Survey
 - [Haas ICDT 07] “Beauty and the Beast: The Theory and Practice of Information Integration”
 - [Kolatis PODS 05] “Schema mappings, data exchange, and metadata management”
 - [Halevy, Rajaraman, Ordille VLDB 06] “Data Integration: The Teenage Years” as part of their 10-year Best Paper Award on “Information Manifold” paper
- Schema Matching
 - Survey by Erhard Rahm and Philip Bernstein in VLDB J. ‘01
 - Much work by Philip Bernstein, Anhai Doan, Alon Halevy, Jayant Madhavan
 - Information Discovery project at IBM Almaden (Berthold Reinwald et al.)
- Data cleansing, conflict resolution, data quality, duplicate detection, entity resolution, data provenance
- Orchestra project at U Penn
- Schema Mapping Debugger
 - [Alexe, Chiticariu, Tan. VLDB 06]
- A Benchmark for Schema Mapping System
 - [Alexe, Tan, Velegarakis. 07]
- A new schema mapping GUI based on XQBE ideas
 - [Ceri, Hernandez, Raffio et al. 07]
- Deep Web, peer-to-peer, ontology, etc



Acknowledgements: Clio Research Team

- The Project Founders:
 - Laura Haas and Renee Miller (around 1999)
- Regulars:
 - Clio Group: Mauricio Hernandez, Howard Ho, Lucian Popa, Ioana Stanoi
 - Theory Group: Ron Fagin, Phokion Kolaitis, Alan Nash
- Past Visiting scientists:
 - Stefan Dessloch, Paolo Papotti, Wang-Chiew Tan, Yannis Velegrakis, Cathy Wyss
- Past Interns and Post-docs:
 - Bogdan Alexe, Alfredo Camacho, Laura Chiticariu, Ariel Fuxman, Michael Gubanov, Haifeng Jiang, Felix Naumann, Paolo Papotti, Ahmed Radwan, Alessandro Raffio, Michael Richmond, Shivkumar Shivaji, Yannis Velegrakis, Melanie Weis, Ryan Wisnesky, Lingling Yan, Cong Yu, Jindan Zhou

