# Meta Data Management
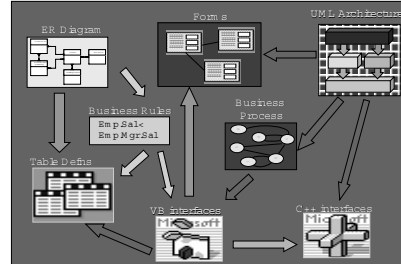
Philip A. Bernstein
Sergey Melnik
{philbe, melnik}@microsoft.com

Microsoft Research

Modified version of the seminar presented at ICDE,
Boston, April 1, 2004 – for presentation at CSEP 544
© 2004 Microsoft Corporation. All rights reserved.

---

# Meta Data Management

1 Meta data = structural information
   DB schema, interface defns, web site map, form defns,…

---

# Meta Data Problems

1 Many data management applications primarily involve transformations of structured data

| | |
|---|---|
| 1 Data translation | 1 UI/4GL generation |
| 1 Schema evolution | 1 Dependency tracking |
| 1 XML message translation | 1 Lineage tracing |
| 1 Application integration | 1 Info resource mgmt |
| 1 Data warehouse loading | 1 Binding, renaming |
| 1 ER/UML design tools | 1 Software build (make) |
| 1 Wrapper generation for SQL | 1 Configuration mgmt |

---

# Outline

1 Introduction
1 Meta data problems
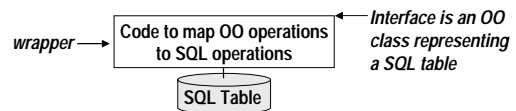1 Design patterns
1 Solution templates
1 Wrap up

---

# Why Meta Data is Important

1 Many DB problems are easier to solve by manipulating meta data
   Instead of writing code
   Instead of manipulating data directly
1 Meta-data-based solutions all involve models (schemas) and mappings
   Mappings - data transformations, queries, dependencies,…
   Model, manipulate, and generate them
   Usually, generate code from them
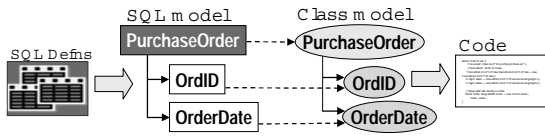
---

# Example: Object-Oriented Wrapper for SQL Tables



1 Manually program a wrapper for each table
1 This is very repetitive work
1 So you write a program to _generate_ a wrapper for each table

---

## OO wrapper for SQL (cont'd)



SQL model    Class model

SQL Defns → PurchaseOrder ⤏ PurchaseOrder → Code
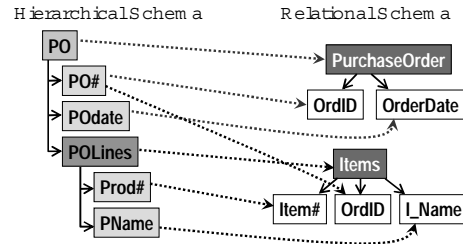OrdID ⤏ OrdID
OrderDate ⤏ OrderDate

1  The wrapper generator does the following:

Imports each table definition into a model

Generates a model for the class wrapper

Generates a mapping from table to class

Generates code from the class model and mapping.

## Example – Data Translation

1  Translate data from one data model to another
1  Either write a program or generate it

Hierarchical Schema    Relational Schema



PO
PO#
POdate
POLines
  Prod#
  PName

PurchaseOrder
OrdID    OrderDate
Items
Item#    OrdID    I_Name

## Meta-data-Speak

| Meta-data-Speak | English Example |
|---|---|
| meta-meta-model = meta-meta-meta data | Built-in types (usually hard-coded) |
| metamodel = meta-meta data | Schema for "Table," "Column", "Key," … |
| model = meta data | Schema for the Employee Table |
| data | Employee Table |

## Outline

1  Introduction
⇨  Meta data problems
1  Design patterns
1  Solution templates
1  Wrap up

## Meta Data Solution Template
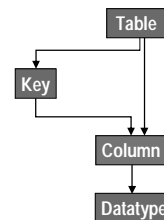
1  Get a data manager for models and mappings

1  Usually, it's an object manager
    OO programming language
    OODB

1  Hence, meta-metamodel is the object manager's built-in types
    Classes, attributes, methods, objects
    Plus operators to manipulate them, such as New Class, New Attribute, New Object, Write Attribute

## Meta Data Solution Template

1  Get a data manager for models and mappings

2  Design metamodel(s) (e.g., for SQL schemas)



Table
Key
Column
Datatype

If the meta-metamodel is OO, then the metamodel consists of class definitions

## Meta Data Solution Template

1 Get a data manager for models and mappings
2 Design metamodel(s) (e.g., for SQL schemas)
3 Build a model importer for each metamodel

SQL Defns → Model Importer → Model

SQL Metamodel
Instance-of

---

## Meta Data Solution Template

1 Get a data manager for models and mappings
2 Design metamodel(s) (e.g., for SQL schemas)
3 Build a model importer for each metamodel
4 Invoke model importer(s)

Model Importer → Model, Model

---

## Meta Data Solution Template

1 Get a data manager for models and mappings
2 Design metamodel(s) (e.g., for SQL schemas)
3 Build a model importer for each metamodel
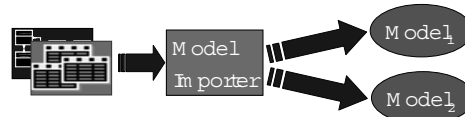4 Invoke model importer(s)

| Problem | $Model_1$ | $Model_2$ |
|---|---|---|
| Data translation | source schema | target schema |
| Msg translation | source format | target format |
| App integration | source interfaces | target interfaces |
| DW loading | source schema | DW schema |

---

## Meta Data Solution Template

1 Get a data manager for models and mappings
2 Design metamodel(s) (e.g., for SQL schemas)
3 Build a model importer for each metamodel
4 Invoke model importer(s)
5 Generate or design mappings

ER Model → generate → SQL Model

Data Source Model ← design → Data Target Model

---

## Meta Data Solution Template

1 Get a data manager for models and mappings
2 Design metamodel(s) (e.g., for SQL schemas)
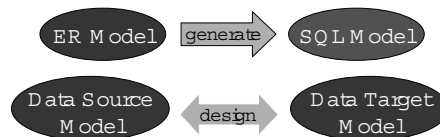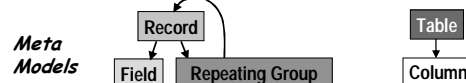3 Build a model importer for each metamodel
4 Invoke model importer(s)
5 Generate or design mappings
6 Generate code: data / msg translatn script, app wrapper, ETL script, view defn's, etc.

---

## Example – Data Translation

**Meta Models**

Record
Field   Repeating Group

Table
Column

---

## Example – Data Translation

**Meta Models**

Record → Field, Repeating Group

Table → Column

**Schemas a.k.a. Models**

PO → PO#, POdate, POLines → Prod#, PName

PurchaseOrder → OrdID, OrderDate

Items → Item#, OrdID, I_Name

---

## Example – Data Translation

**Meta Models**

Record → Field, Repeating Group

Table → Column

**Schemas a.k.a. Models**

PO → PO#, POdate, POLines → Prod#, PName

PurchaseOrder → OrdID, OrderDate

Items → Item#, OrdID, I_Name

---

**Generated data translation script**

```
Foreach [po#,poD,poL] in PO
    Insert [po#,poD] into PurchaseOrder
    Foreach [prod#,pN] in poL
        Insert [prod#,po#,pN] into Items
    End
End
```

**Schemas a.k.a. Models**

PO → PO#, POdate, POLines → Prod#, PName

PurchaseOrder → OrdID, OrderDate

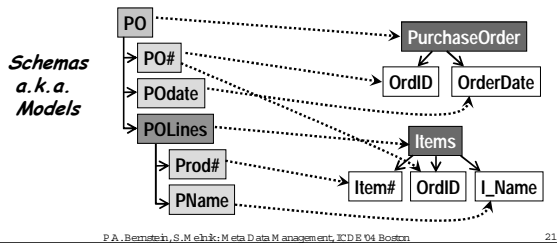Items → Item#, OrdID, I_Name

---

## Meta Data Problems

- Data translation
- OO or XML wrapper generation for SQL DB
- User-Interface / 4GL-program generation
- Design tool support (DB, UML, …)
    - Model generation, reverse engineering
    - Round-trip engineering
- Schema evolution (applies to all scenarios)
- XML message translation for e-commerce
- Integrate custom apps with commercial apps

---

## Meta Data Problems (cont'd)

- Data warehouse loading (clean & transform)
- Lineage tracing (provenance)
- Information resource management
- Dependency tracking
    - Impact analysis
    - Navigation between tools
- Binding, renaming
- Software build (make)
- Version and configuration management
    - Release management
    - Product data management

---

## Meta Data Solutions

- They strongly resemble one another
- We characterize that resemblance
    - Prototypical problems, or design patterns
    - Solution specifications, or solution templates
    - Primitive solution steps, or operators
- Goals
    - A methodology to solve meta data problems
    - Ultimately, operator implementations to turn solution templates into solution programs

---

## Outline

l   Introduction
l   Meta data problems
⇒  Design patterns
l   Solution templates
l   Wrap up

## Meta Data Design Patterns

l  Design pattern – a problem description consisting of
   Input models and mappings
   Output models and mappings
   Criteria for the output to be correct
   An application specializes it to meta models and mapping languages

l  Solution template – a sequence of operators producing the desired output

l  Operators – a single step that computes a model and/or mappings

## Operators

l  $map = Match(M_1, M_2)$
   Return a mapping between the two models

l  $\langle M_2, map_{12} \rangle = ModelGen(M_1, metamodel_2)$
   Return a model $M_2$ that is expressed in metamodel$_2$ and is equivalent to model $M_1$

l  $\langle M_3, map_{13}, map_{23} \rangle = Merge(M_1, M_2, map)$
   Return the union of models $M_1$ and $M_2$

l  $map_3 = Compose(map_1, map_2) = map_1 \circ map_2$
   Return the composition of $map_1$ and $map_2$, which is a mapping from $map_1$'s domain to $map_2$'s range.

## Operators (cont'd)

l  $map_3 = Confluence(map_1, map_2) = map_1 \, ^- \, map_2$
   Return the "merge" of mappings $map_1$ and $map_2$

l  $\langle M_2, map_{12} \rangle = Extract(M_1, map)$
   Return the sub-model of $M_1$ that participates in the mapping $map$

l  $\langle M_2, map_{12} \rangle = Diff(M_1, map)$
   Return the sub-model of $M_1$ that does not participate in the mapping $map$

## Design Patterns

l  Meta Modeling
l  Model Mapping
l  Model Generation
l  Model Integration ⎱
l  Mapping Composition ⎰ Single Operator Solutions
l  Mapping Alignment
l  Change Propagation
l  Model Reintegration

## Meta Modeling

l  Design pattern – develop a representation (i.e. metamodel) for models and mappings
l  Applications – they all depend on this
l  Solution template
   Design a metamodel
   Write Import & Export functions
   § ImportSQL, ImportXSD, ImportERD, …
   Today, it is manual engineering design
   Design once and reuse often

## Meta Modeling (cont'd)

1 The Import function for models
   - Parse text
   - Copy elements of the parsed form into a model that conforms to its metamodel

1 The Import function for mappings
   - Same as models but may require more semantic analysis
   - E.g., program dependencies, data lineage
   - For some languages and mapping metamodels, Export is hard (e.g., XSLT)

## Model Mapping

1 Design pattern – Design a mapping between two models and generate code from it

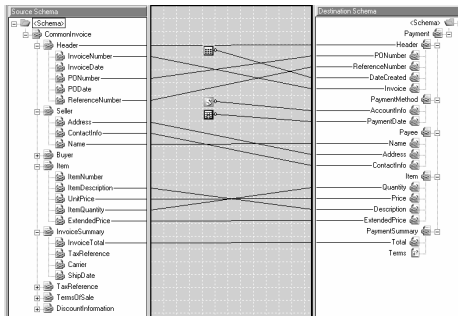$$M_1 \quad \xleftrightarrow{\text{map}} \quad M_2$$

1 Applications
   - Data translation
   - XML message translation for e-commerce
   - Integrate custom and commercial apps
   - Data warehouse extract, transform & load

1 Solution templates
   - map = Match ($M_1$, $M_2$); Export(map)
   - Mapping reuse: Compose, Confluence

## An XML Mapping Tool

## A Data Warehouse Loading Tool

## Model Generation

1 Design pattern – Given a model, generate an equivalent model in another metamodel
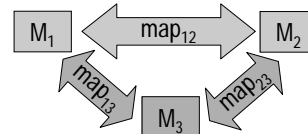
$$M_1 \quad \xleftrightarrow{\text{map}} \quad M_2$$

1 Applications
   - Wrapper generation (SQL → OO or XML)
   - Design tools (ER → SQL, SQL → ER)
   - UI/4GL generation

1 Solution template
   - ⟨$M_2$, map⟩ = ModelGen ($M_1$, metamodel$_2$); Export($M_2$)
   - ModelGen often needs human guidance

## Model Integration

1 Design pattern – Given two models, develop a model that subsumes both of them

$$M_1 \xleftrightarrow{\text{map}_{12}} M_2 \quad \text{map}_{13} \quad M_3 \quad \text{map}_{23}$$

1 Applications
   - View integration
   - Data integration

1 Solution template
   - ⟨$M_3$, map$_{13}$, map$_{23}$⟩ = Merge ($M_1$, $M_2$, map); Export($M_3$, map$_{13}$, map$_{23}$)

## Mapping Composition

- Design pattern – Compose two given mappings



$$M_1 \quad map_{12} \quad M_2 \quad map_{23} \quad M_3$$
$$map_{13}$$
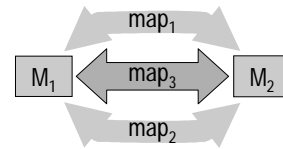
- Applications
  - Processing queries on views
- Solution templates
  - $map_{13} = Compose(map_{12}, map_{23})$
  - Answering queries using views (LaV),
  - Query modification (GaV), GLaV

## Mapping Alignment

- Design pattern – Align two mappings between the same pair of models
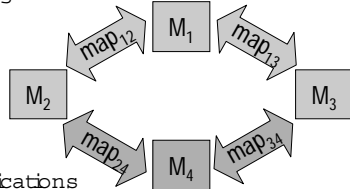


$$map_1$$
$$M_1 \quad map_3 \quad M_2$$
$$map_2$$

- Applications
  - P2P query processing, mapping design
- Solution template
  - $map_3 = Confluence(map_1, map_2)$

## Model Reintegration

- Design pattern – Given a model and mappings to two modified versions of the model, produce a merged model



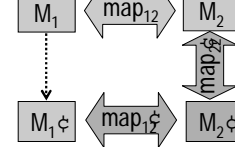$$map_{12} \quad M_1 \quad map_{13}$$
$$M_2 \qquad\qquad M_3$$
$$map_{24} \quad M_4 \quad map_{34}$$

- Applications
  - Parallel development
- Solution template
  - Multistep application of many operators

## Change Propagation

- Design pattern – Given two models and a mapping. One model changes. Fix the mapping and other model.



$$M_1 \quad map_{12} \quad M_2$$
$$M_1¢ \quad map_{12}¢ \quad M_2¢$$

- Applications
  - Schema evolution, interface evolution, ...
  - Required maintenance for all meta data problems
- Solution template
  - Requires all of the operators

## Outline

- Introduction
- Meta data problems
- Design patterns
- → Solution templates
  - Change propagation
  - Model reintegration
  - Change propagation revisited
- Research background
- Wrap up

## Change Propagation

- Given
  - $map_1$ between xsd1 and SQL schema rdb1
  - xsd2, a modified version of xsd1
- Produce
  - rdb2 to store instances of xsd2
  - a mapping between xsd2 and rdb2



$$xsd1 \quad map_1 \quad rdb1$$
$$xsd2 \quad map \quad rdb2$$

7

## Change Propagation

- Given
    - $map_1$ between xsd1 and SQL schema rdb1
    - xsd2, a modified version of xsd1
- Produce
    - rdb2 to store instances of xsd2
    - a mapping between xsd2 and rdb2

xsd1 — $map_1$ — rdb1

1. $map_2$ (xsd1 to xsd2)

xsd2

1. $map_2$ = Match(xsd1, xsd2)

---

## Change Propagation

- Given
    - $map_1$ between xsd1 and SQL schema rdb1
    - xsd2, a modified version of xsd1
- Produce
    - rdb2 to store instances of xsd2
    - a mapping between xsd2 and rdb2

xsd1 — $map_1$ — rdb1

1. $map_2$ ; 2. $map_3$

xsd2

1. $map_2$ = Match(xsd1, xsd2)

2. $map_3 = map_2 \bullet map_1$

---

## Change Propagation

- Given
    - $map_1$ between xsd1 and SQL schema rdb1
    - xsd2, a modified version of xsd1
- Produce
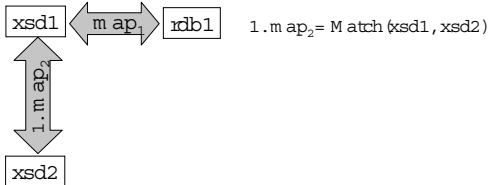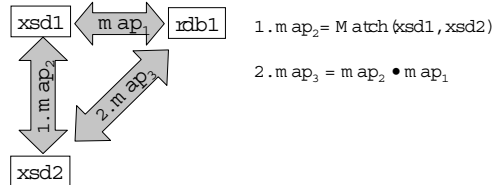    - rdb2 to store instances of xsd2
    - a mapping between xsd2 and rdb2

xsd1 — $map_1$ — rdb1

1. $map_2$ ; 2. $map_3$

xsd2 — $map$ — rdb2

1. $map_2$ = Match(xsd1, xsd2)

2. $map_3 = map_2 \bullet map_1$

3. $<map_4, rdb3> = Copy(map_3)$

Now we need to merge
Diff(xsd2, $map_4$) into rdb3

---

## Change Propagation (cont'd)

xsd1 — $map_1$ — rdb1

1. $map_2$ ; 2. $map_3$

xsd2 — 3. $map_4$ — rdb3

4. $map_5$

xsd2¢

4. $<xsd2¢, map_5> =$ Diff(xsd2, $map_4$)

---

## Change Propagation (cont'd)

xsd1 — $map_1$ — rdb1

1. $map_2$ ; 2. $map_3$

xsd2 — 3. $map_4$ — rdb3

4. $map_5$

xsd2¢ — 5. $map_6$ — rdb4

4. $<xsd2¢, map_5> =$ Diff(xsd2, $map_4$)

5. $<rdb4, map6> =$ ModelGen(xsd2¢, SQL)

---

## Change Propagation (cont'd)

xsd1 — $map_1$ — rdb1

1. $map_2$ ; 2. $map_3$

xsd2 — 3. $map_4$ — rdb3

4. $map_5$ ; 6. $map_7$ ; 7. $map_9$

xsd2¢ — 5. $map_6$ — rdb4 ; 7. $map_8$ — rdb2

4. $<xsd2¢, map_5> =$ Diff(xsd2, $map_4$)

5. $<rdb4, map6> =$ ModelGen(xsd2¢, SQL)

6. $map7 = map4 \bullet map5 \bullet map6$

7. $<rdb2, map8, map9> =$ Merge(rdb3, rdb4, map7)

## Complete Script in Rondo

OperatorDefinition: PropagateChanges(s1,d1,s1_d1,s2,c,s2_c)

1. s1_s2 = Match(s1,s2);
2. ⟨d1¢,d1¢_d1⟩ = Delete(d1,Traverse(All(s1) - Domain(s1_s2)),s1_d1));
3. ⟨c¢,c¢_c⟩ = Extract(c,Traverse(All(s2) - Range(s1_s2)),s2_c));
4. c¢_d1¢= c¢_c * Invert(s2_c) * Invert(s1_s2) * s1_d1 * Invert(d1¢_d1);
5. ⟨d2,c¢_d2,d1¢_d2⟩ = Merge(c¢,d1¢,c¢_d1¢;
6. s2_d2 = s2_c * Invert(c¢_c) * c¢_d2 +
              Invert(s1_s2) * s1_d1 * Invert(d1¢_d1) * d1¢_d2;
7. return ⟨d2,s2_d2⟩;

OperatorUse:

SQL_XSD:PropagateChanges(s1,d1,s1_d1,s2,ModelGen(s2,XSD));

---

## Model reintegration

- Design pattern
  - Reconcile independent changes
  - All changes of each model
  - No "duplicate additions"
- Simplified example
  - "Additions" = add model element
    (also: drop constraints, reorg. model)
  - "Deletions" = delete model element
    (also: add constraints, reorg. model)
  - Mappings shown as lines betw. elements

---

## Design pattern

---

## Naive solution



- Direct model integration "loses" either deletions or additions
- Need $m$, $m\_m_A$, $m\_m_B$

---

## Solution template (begin)



- Propagate deletions before integrating $m_A$ and $m_B$

---



- Composition produces a partial mapping

Invert($m\_m_A \circ m_A\_m_A$¢) $\circ m\_m_B \circ m_B\_m_B$¢

---

**Slide 55:**

m — OrderID / Product / Brand

$m\_m_A \circ m_A\_m_A¢$        $m\_m_B \circ m_B\_m_B¢$

$m_{Ax}$        $m_{Bx}$

Rebate        Discount / ShipDate

Diff        Diff

1 Identify additions (Diff)

OrderID / Product / Rebate — $m_A¢$        OrderID / Product / Discount / ShipDate — $m_B¢$

**Slide 56:**

m — OrderID / Product / Brand

$m\_m_A \circ m_A\_m_A¢$        $m\_m_B \circ m_B\_m_B¢$

$m_{Ax}$   Match   $m_{Bx}$

Rebate —— Discount / ShipDate

Diff        Diff

1 Identify additions (Diff)

1 Match $m_{Ax}$ and $m_{Bx}$

OrderID / Product / Rebate — $m_A¢$        OrderID / Product / Discount / ShipDate — $m_B¢$

**Slide 57:**

m — OrderID / Product / Brand

$m\_m_A \circ m_A\_m_A¢$        $m\_m_B \circ m_B\_m_B¢$

$m_{Ax}$   Match   $m_{Bx}$

Rebate —— Discount / ShipDate

Diff        Compose        Diff

1 Identify additions (Diff)

1 Match $m_{Ax}$ and $m_{Bx}$

1 Compose

OrderID / Product / Rebate — $m_A¢$        OrderID / Product / Discount / ShipDate — $m_B¢$

**Slide 58:**

m — OrderID / Product / Brand

$m\_m_A \circ m_A\_m_A¢$        $m\_m_B \circ m_B\_m_B¢$

$m_{Ax}$   Match   $m_{Bx}$

Rebate —— Discount / ShipDate

Diff        Compose        Diff

1 Identify additions (Diff)

1 Match $m_{Ax}$ and $m_{Bx}$

1 Compose

1 Align partial mappings (Confl.)

OrderID / Product / Rebate — $m_A¢$        OrderID / Product / Discount / ShipDate — $m_B¢$

$m_A¢\_m_B¢$

**Slide 59:**

m — OrderID / Product / Brand

OrderID / Product / Rebate — $m_A¢$        OrderID / Product / Discount / ShipDate — $m_B¢$

$m_A¢\_m_B¢$

**Slide 60:**

m — OrderID / Product / Brand

OrderID / Product / Rebate / Discount / ShipDate

$m_R$   Merge

1 Merge $m_A¢$ and $m_B¢$ using $m_A¢\_m_B¢$

OrderID / Product / Rebate — $m_A¢$        OrderID / Product / Discount / ShipDate — $m_B¢$

$m_A¢\_m_B¢$

## Slide 61

m

OrderID
Product
Brand

Solution template
(end)

OrderID
Product
Rebate/Discount
ShipDate

1 Merge $m_A\mathcal{C}$ and $m_B\mathcal{C}$ using $m_A\underset{\sim}{\subseteq}m_B\mathcal{C}$

$m_R$   Merge

1 Compose Align partial mappings (Confl.)

OrderID
Product
Rebate

OrderID
Product
Discount
ShipDate

$m_A\mathcal{C}$        $m_B\mathcal{C}$

$m_A\underset{\sim}{\subseteq}m_B\mathcal{C}$

## Slide 62

Solution script

Mapping reuse pattern

script Reintegrate($m$, $m_A$, $m\_m_A$, $m_B$, $m\_m_B$)
1. $\langle m_A\mathcal{C}, m_A\_m_A\mathcal{C}\rangle$ = PropagateDeletions($m$, $m_A$, $m\_m_A$, $m_B$, $m\_m_B$);
2. $\langle m_B\mathcal{C}, m_B\_m_B\mathcal{C}\rangle$ = PropagateDeletions($m$, $m_B$, $m\_m_B$, $m_A$, $m\_m_A$);
3. $\langle m_{Ax}, m_A\mathcal{C}\_m_{Ax}\rangle$ = Diff($m_A\mathcal{C}$, Invert($m\_m_A \circ m_A\_m_A\mathcal{C}$));
4. $\langle m_{Bx}, m_B\mathcal{C}\_m_{Bx}\rangle$ = Diff($m_B\mathcal{C}$, Invert($m\_m_B \circ m_B\_m_B\mathcal{C}$));
5. $m_{Ax}\_m_{Bx}$ = Match($m_{Ax}$, $m_{Bx}$);
6. $m_A\mathcal{C}\_m_B\mathcal{C}$= ($m_A\mathcal{C}\_m_{Ax} \circ m_{Ax}\_m_{Bx} \circ$ Invert($m_B\mathcal{C}\_m_{Bx}$)) $^-$ (Invert($m\_m_A \circ m_A\_m_A\mathcal{C}$) $\circ$ $m\_m_B \circ m_B\_m_B\mathcal{C}$);
7. $\langle m_R, m_R\_m_A\mathcal{C}, m_R\_m_B\mathcal{C}\rangle$ = Merge($m_A\mathcal{C}$, $m_B\mathcal{C}$, $m_A\mathcal{C}\_m_B\mathcal{C}$);
8. $m\_m_R$ = ($m\_m_A \circ m_A\_m_A\mathcal{C} \circ$ Invert($m_R\_m_A\mathcal{C}$)) $^-$ ($m\_m_B \circ m_B\_m_B\mathcal{C} \circ$ Invert($m_R\_m_B\mathcal{C}$));
9. return $\langle m_R, m\_m_R\rangle$;

## Slide 63

### Sub-pattern: propagate deletions

(not deleted by way of $m_A$)

m

OrderID
Product
Brand

$m_B$

OrderID
Product
Brand
Discount
ShipDate

$m_A$

OrderID
Product
Rebate

(added to $m_B$)

## Slide 64

### Sub-pattern: propagate deletions

(not deleted by way of $m_A$)

m

OrderID
Product
Brand

$m_B$

OrderID
Product
Brand
Discount
ShipDate

$m_A$

OrderID
Product
Rebate

Compose

(added to $m_B$)

## Slide 65

### Sub-pattern: propagate deletions

(not deleted by way of $m_A$)

Extract

m

OrderID
Product
Brand

$m_B$

OrderID
Product
Brand
Discount
ShipDate

OrderID
Product

$m_A$

OrderID
Product
Rebate

Compose

(added to $m_B$)

## Slide 66

### Sub-pattern: propagate deletions

(not deleted by way of $m_A$)

Extract

m

OrderID
Product
Brand

$m_B$

OrderID
Product
Brand
Discount
ShipDate

OrderID
Product

$m_A$

OrderID
Product
Rebate

Compose

Diff

Discount
ShipDate

(added to $m_B$)

## Sub-pattern : propagate deletions



m    $m_B$    Extract    (not deleted by way of $m_A$)

| OrderID |
| Product |
| Brand |

| OrderID |
| Product |
| Brand |
| Discount |
| ShipDate |

| OrderID |
| Product |

$m_B$¢

| OrderID |
| Product |
| Discount |
| ShipDate |

$m_A$

| OrderID |
| Product |
| Rebate |

Compose, Confluence

Merge

Compose

Diff

| Discount |
| ShipDate |

(added to $m_B$)

---

## Solution script



```
script PropagateDeletions(s_1, d_1, s_1_d_1, s_2, s_1_s_2)
1.  d_1_s_2 = Invert(s_1_d_1) ○ s_1_s_2;
2.  ⟨m, d_1_m⟩ = Extract(d_1, d_1_s_2);
3.  ⟨n, d_1_n⟩ = Diff(d_1, Invert(s_1_d_1));
4.  ⟨d_2, d_2_m, d_2_n⟩ = Merge(m, n, Invert(d_1_m) ○ d_1_n);
5.  d_1_d_2 = (d_1_m ○ Invert(d_2_m)) ⁻ (d_1_n ○ Invert(d_2_n));
6.  return ⟨d_2, d_1_d_2⟩;
```
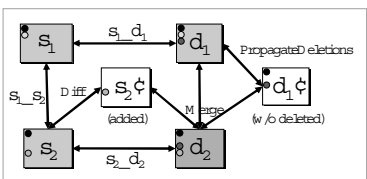
---

## Outline

l   Introduction

l   Meta data problems

l   Design patterns

l   Solution templates

    Change propagation

    Model reintegration

⇨   Change propagation revisited

l   Research background

l   Wrap up

---

## Change propagation

l   Solution template

    Propagate deletions

    Include additions

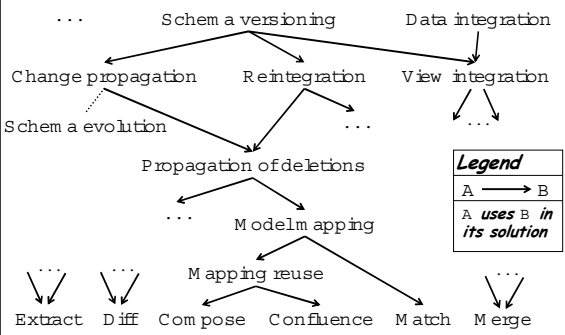    Merge result

---

## Change propagation



```
script PropagateChanges(s_1, d_1, s_1_d_1, s_2, s_1_s_2)
1. ⟨d_1¢, d_1_d_1¢⟩ = PropagateDeletions(s_1, d_1, s_1_d_1, s_2, s_1_s_2);
2. ⟨s_2¢, s_2_s_2¢⟩ = Diff(s_2, Invert(s_1_s_2));
3. ⟨d_2, d_2_s_2¢, d_2_d_1¢⟩ = Merge(s_2¢, d_1¢, Invert(s_1_s_2 ○ s_2_s_2¢) ○ s_1_d_1 ○ d_1_d_1¢);
4. d_1_d_2 = (Invert(s_1_d_1) ○ s_1_s_2 ○ s_2_s_2¢ ○ Invert(d_2_s_2¢)) ⁻ (d_1_d_1¢ ○ Invert(d_2_d_1¢));
5. s_2_d_2 = (Invert(s_1_s_2) ○ s_1_d_1 ○ d_1_d_1¢ ○ Invert(d_2_d_1¢)) ⁻ (s_2_s_2¢ ○ Invert(d_2_s_2¢));
6. return ⟨d_2, d_1_d_2, s_2_d_2⟩;
```

---

## First-cut taxonomy of patterns



...   Schema versioning     Data integration

Change propagation    Reintegration    View integration

Schema evolution    ...    ...

Propagation of deletions

...   Model mapping

Mapping reuse

Extract   Diff   Compose   Confluence   Match   Merge

**Legend**

| A | ⟶ | B |

A **uses** B *in its solution*

---

## Outline

l  Introduction
l  Meta data problems
l  Design patterns
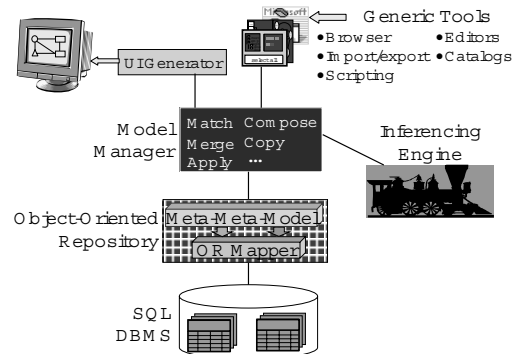l  Solution templates
l  Research background
➭ Wrap up

## The Commercial World

l  Books for IT professionals
   A. Tanenbaum: Metadata Solutions, Addison-Wesley, 2001
   D. Marco: Building and Managing the Meta Data Repository, Wiley, 2000
l  Standards-
   UML, MOF, CWM (OMG)
   XML, RDF, XML Schema, OWL (W3C)
l  Products and tools
   Modeling: IBM Rational Rose, Visio, CA AllFusion, Borland Together
   General meta data managers: CA Advantage, Microsoft Meta Data Services, MetaIntegration
   Meta data services in data warehousing ETL tools: Informatica, Ascential, ETI, Data Advantage, …

## The Research World

l  Model Management
   A computational meta data framework based on models, mappings, and the operators described here (Match, Merge, Compose, … )
l  Meta Data is a very active research area
   Papers coming from many DB research groups
   Some are problem-focused (e.g. data integration)
   Some are operator-focused (e.g. Match, Merge)

## MM System Architecture

## Summary

l  Many DB problems are easier to solve by manipulating meta data
l  Meta data problems and solutions strongly resemble one another
l  Methodology: Use design patterns, solution templates, and operators to simplify development of meta data applications
l  There is much research to be done

## References

l  http://research.microsoft.com/db/ModelMgt
l  Overview
   Bernstein, CIDR 2003
   Bernstein, Halevy, Pottinger, SIGMOD Record, Dec. 2000
l  Implementation
   Melnik, Rahm, & Bernstein, SIGMOD 2003
   and J. Web Semantics 1, 2003
l  Data Warehouse Examples
   Bernstein & Rahm, ER 2000
l  Match Operation
   Survey: Rahm & Bernstein, VLDB J., Dec. 2001
l  Merge Operation
   Pottinger & Bernstein, VLDB 2003

# The Match "Operator"

l Schema matching (mapping discovery)

Given two schemas, return correspondences that specify pairs of related elements

l Semantic Mapping (query discovery)

Given correspondences between two schemas, return an expression that translates instances of one schema into instances of the other.

P.A.Bernstein, S.Melnik: Meta Data Management, ICDE '04 Boston 80
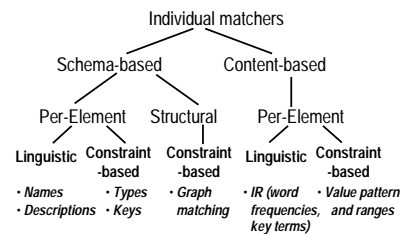
---

# Schema Matching Problem

l Input

Schemas $S_1$ and $S_2$

Possibly data instances for $S_1$ and $S_2$

Background knowledge – thesauri, validated matches, standard schemas, constraints (keys, data types), ontologies, NL glossaries, etc.

l Output

Correspondences between elements of $S_1$ and $S_2$

Used by permission of Erhard Rahm P.A.Bernstein, S.Melnik: Meta Data Management, ICDE '04 Boston 81

---

# Schema Matching Approaches

l Many good ideas

Rahm & Bernstein, VLDB J, Dec '01

l But none are robust fi combine ideas
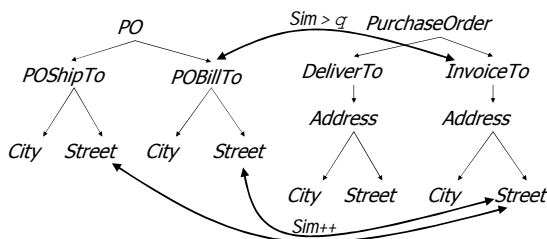


P.A.Bernstein, S.Melnik: Meta Data Management, ICDE '04 Boston 82

---

# The Cupid Algorithm

l Computes linguistic similarity of element pairs
l Computes structural similarity of element pairs
l Generates a mapping



P.A.Bernstein, S.Melnik: Meta Data Management, ICDE '04 Boston 83
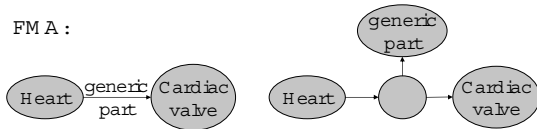
---

# Matching Anatomy Ontologies

l Match two human anatomy ontologies

FMA – Univ. of Washington

Galen CRM – Univ. of Manchester (UK)

By Peter Mork (Univ. of Washington)

Both models are big

l Ultimate goal was finding differences

l Like most match algorithms, ours calculates a similarity score for the m·n pairs of elements

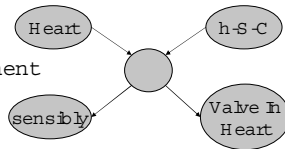P.A.Bernstein, S.Melnik: Meta Data Management, ICDE '04 Boston 84

## Aligning Representations

FMA:



CRM:
Heart sensibly
hasStructuralComponent
ValveInHeart

---

## Anatomy Matching Algorithm

1. Lexical Match
   - Normalize string, UMLS dictionary lookup, convert to concept-ID from thesaurus

---

## Anatomy Matching Algorithm
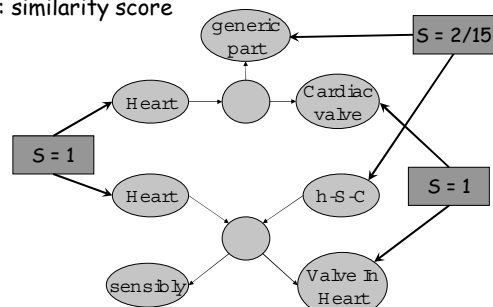
1. Lexical Match
   - Normalize string, UMLS dictionary lookup, convert to concept-ID from thesaurus

   - String comparison **fi** 306 matches
   - Adding spaces, ignoring case **fi** 1834 matches
   - Lexical tools **fi** 3503 matches

---

## Anatomy Matching Example

S: similarity score

---

## Anatomy Matching Algorithm

1. Lexical Match
   - Normalize string, UMLS dictionary lookup, convert to concept-ID from thesaurus
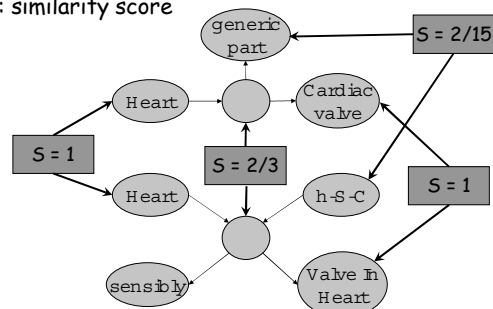2. Structure Match
   - Similarity (reified nodes) = Average (neighbors)
   - Back-propagate to neighbors

---

## Anatomy Matching Example

S: similarity score

---

## Anatomy Matching Algorithm

1. Lexical Match
   - Normalize string, UMLS dictionary lookup, convert to concept-ID from thesaurus
2. Structure Match
   - Similarity (reified nodes) = Average (neighbors)
   - Back-propagate to neighbors

- Adds 64 matches (to previous 3503)
  - Implies 875 reified relationship matches

## Anatomy Matching Algorithm

1. Lexical Match
   - Normalize string, UMLS dictionary lookup, convert to concept-ID from thesaurus
2. Structure Match
   - Similarity (reified nodes) = Average (neighbors)
   - Back-propagate to neighbors
3. Align Super-classes
   - Super-class similarity = average similarity of children, grandchildren, great-grandchildren
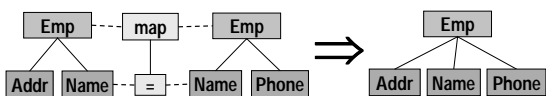- Adds 213 matches (to 3567)

## Some Lessons

- A common encoding of models is hard and involves compromises
  - Different styles of reifying relationships
  - CRM stores transitive relationships
- Match needs to invent generalizations
  - In FMA, *arterial supply, venous drainage, nerve supply, lymphatic drainage*
  - In CRM, these all map to *isServedBy*
- On big models, Match is expensive
  - Some steps required days to execute
  - Cross-product filled 80 GB (< 1GB input).

## Outline

Introduction to Model Management
Using MM to solve meta data problems
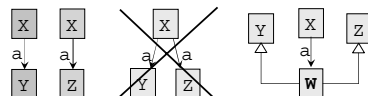Matching anatomy ontologies
- Model merging
- Wrap-up

## Merge ($M_1$, $M_2$, map)

- Return the union of models $M_1$ and $M_2$
  - Use map to guide the Merge
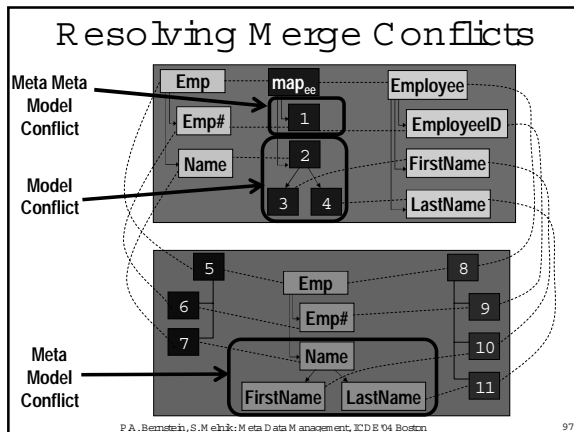  - If elements x = y in map, then collapse them into one element

## Merge ($M_1$, $M_2$, map)

- [Buneman, Davidson, Kosky, EDBT 92]
  - Meta-model has aggregation & generalization only
  - Union, and collapse objects having the same name
  - Fix-up step for inconsistencies created by merging



- Successive fixups lead to different results L
- Batch them at the end, to get a unique minimal result
- Now enrich the meta-model (containment, complex mappings, …) & merge semantics (conflicts, deletes)

## Resolving Merge Conflicts



Meta Meta Model Conflict

Model Conflict

Meta Model Conflict

Emp | map_ee | Employee
Emp# | 1 | EmployeeID
Name | 2 | FirstName
| 3 4 | LastName

Emp
Emp#
Name
FirstName  LastName

## Contributions to Merge
### [Pottinger & Bernstein, VLDB 03]

l  Generic correctness criteria for Merge
l  Use of first-class input mapping (not just correspondences)
l  Taxonomy of conflicts & resolution strategies
l  Characterize when Merge can be automatic
l  A merge algorithm for an EER representation
l  Experimental evaluation

## An Approach to ModelGen
### [Atzeni & Torlone, EDBT '96]

l  Meta-models are made of patterns

Object has △       Aggregation ⊗       Aggregation ⊗
sub-object △       has attributes □       has key ↓
   (a)                   (b)                   (c)

l  Define pattern transformations as rules

For XSD → SQL, △ ⟹ ⊗ + ⊗

l  To translate $M_s$ into meta-model ($MM_t$),
   Apply rules that replace patterns in $M_s$ that are
   not in $MM_t$ by patterns that are in $MM_t$

## ModelGen Research

l  More complete repertoire of patterns
l  Make patterns more generic
l  Integrate with rules engine (avoid cycles, control search)
l  Implement it

*Microsoft*®