

CS 245: Database System Principles

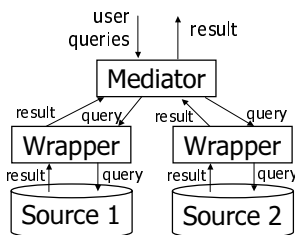
Notes 14: Coping with Limited Capabilities of Sources

Hector Garcia-Molina
(Some modifications by Chris Olston)

Recall ...

- Three approaches to information integration:
 - Federated databases ← did teaser
 - Data warehousing ← did overview
 - Mediation ← next

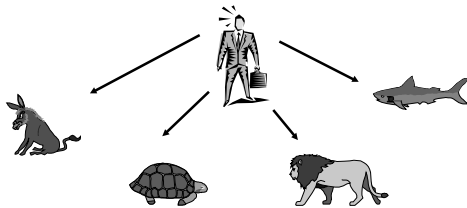
Third Approach to Information Integration: Mediation



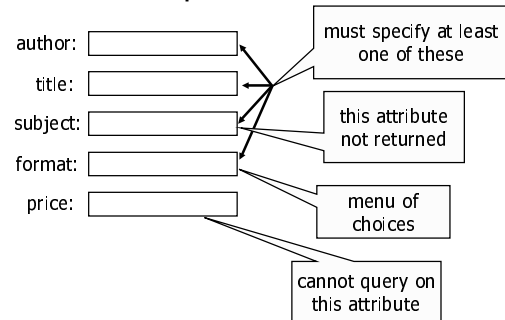
Recall Advantages of Mediation

- No need to copy data
 - less storage
 - no need to purchase data
- More up-to-date data
- Query needs can be unknown
- Only query interface needed at sources
 - Sources may not support "upload all data"

- Worse, sources may provide very limited and differing query capabilities
- One major challenge for mediation: coping with limited source capabilities



Example: Amazon.com



Example: BarnesAndNoble.com

author:

title:

subject:

format:

price:

must specify at least one of these

Menu of choices

can query if one of other attributes specified

CS 245 Notes 14 7

Why Limited Capabilities?

- Search forms
- Security
- Indexes
- Legacy

CS 245 Notes 14 8

Capability vs. Content

- Capability description
 - Can only search for subject = "art," "history," "science"
- Content description
 - Source only contains subject = "art," "history," "science"

CS 245 Notes 14 9

Outline

- Describing source capabilities
- Extending source capabilities
- How mediators cope with limited capabilities
- Mediator capabilities
- Other topics

```

graph TD
    M[mediator] --- S1[source]
    M --- S2[source]
    M --- S3[source]
    P((stick figure)) --- M
  
```

CS 245 Notes 14 10

Describing Query Capabilities

R(X, Y, ... Z)

Adornments:

- **f**: may or may not specify
- **u**: cannot be specified
- **b**: must be specified
- **c[S]**: specified from list S
- **o[S]**: optional, chose from S

CS 245 Notes 14 11

Describing Query Capabilities

R(X, Y, ... Z)

Adornments:

- **f**: may or may not specify
- **u**: cannot be specified
- **b**: must be specified
- **c[S]**: specified from list S
- **o[S]**: optional, chose from S

With output restriction

- **f'**
- **u'**
- **b'**
- **c'[S]**
- **o'[S]**

CS 245 Notes 14 12

Example

- Relation $R(X, Y, Z)$
- Description Templates: $bu'f, uf'c[z_1, z_2]$
- Answerable queries: $R(x_1, Y, Z), R(X, Y, z_1)$
- Unanswerable queries:
 $R(X, y_1, Z), R(X, Y, z_3)$

CS 245

Notes 14

13

Many Other Description Mechanisms Exist ...

- Tsimmis
- Information Manifold
- Disco
- Garlic
- Context-free grammars

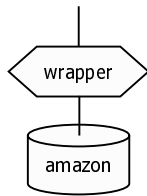
CS 245

Notes 14

14

Extending Source Capabilities

Query: $author="Freud" \text{ AND } price > 10$



Source: $R(author, price, \dots)$
 Template: b, u, \dots

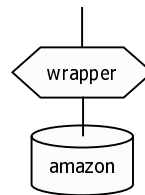
CS 245

Notes 14

15

Extending Source Capabilities

Query: $author="Freud" \text{ AND } price > 10$



Wrapper Filter: $price > 10$

Source Query: $author="Freud"$

Source: $R(author, price, \dots)$
 Template: b, u, \dots

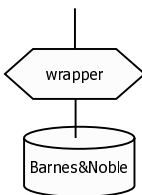
CS 245

Notes 14

16

Another Example

Query: $(author = "Freud" \text{ OR } author = "Jung") \text{ AND } price < 10$



$R(author, price, \dots)$
 No disjunctive conditions;
 Price can only be specified with author

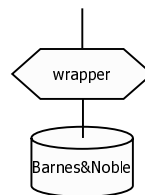
CS 245

Notes 14

17

Another Example

Query: $(author = "Freud" \text{ OR } author = "Jung") \text{ AND } price < 10$



Union Operation

Q1: $author = "Freud" \text{ AND } price < 10$
 Q2: $author = "Jung" \text{ AND } price < 10$

$R(author, price, \dots)$
 No disjunctive conditions;
 Price can only be specified with author

CS 245

Notes 14

18

Extending Source Capabilities

- General scheme:
 - try many query rewritings
 - check if query fragments supported by source
 - check if wrapper can combine answer fragments
 - do all this very efficiently!! [See ICDE99 paper]
- Tsimmis, Info Manifold: no disjunctive queries
- DISCO: no query splitting
- Garlic: only CNF queries

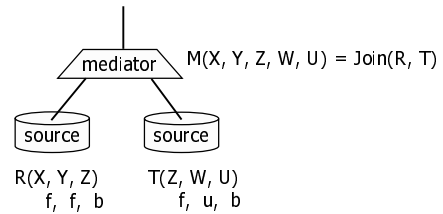
CS 245

Notes 14

19

Mediator Processing

Query: $M(5, Y, Z, W, 3)$

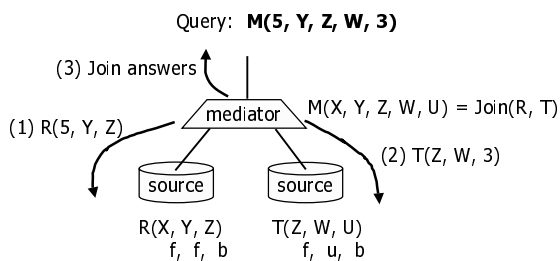


CS 245

Notes 14

20

Plan 1

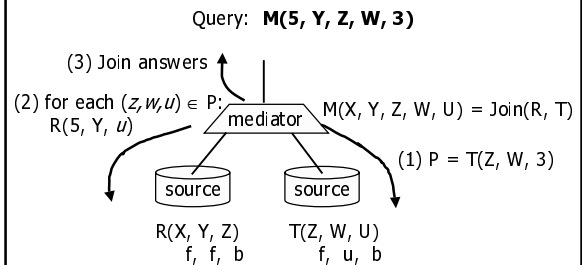


CS 245

Notes 14

21

Plan 2



CS 245

Notes 14

22

Mediator Plan Generation

- Need feasible and efficient plan
- Search space is huge
- Tsimmis, Info Manifold, Garlic:
 - exponential algorithms
- Polynomial algorithms:
 - often find optimal or near-optimal plan
 - bounded performance
 - [See ICDT99 Paper]

CS 245

Notes 14

23

Conclusion

- Not all sources are created equal!
- Need to
 - describe what sources can do
 - efficiently process queries even though sources have limited capabilities
 - describe what mediators can do
 - exploit content information
 - deal with unavailable sources

CS 245

Notes 14

24