

Chapter 10 Wrappers and External Data



Data Federation/Integration

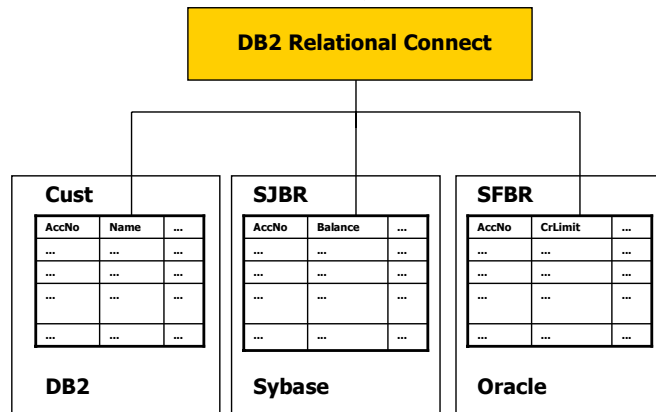
- Goal: homogeneous, integrated view of data from multiple sources
 - a single logical database
 - a single query may collect (or join) data from multiple sources
- requires
 - Wrappers
 - Data and schema integration mechanisms



Example – Federating Relational Sources

```

Select *
From Cust, SJBR, SFBR
Where Cust.Acct No = SJBR.Acct No
And SJBR.Acct No = SFBR.Acct No
    
```



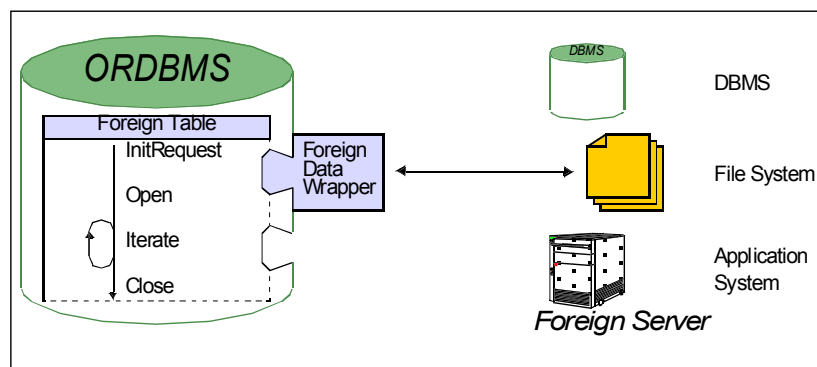
© Prof. Dr.-Ing. Stefan Dießloch

3

Middleware for Heterogenous and Distributed Information Systems - WS05/06

SQL – Management of External Data (MED)

- 'Foreign Data Wrapper' in 'SQL/MED'



© Prof. Dr.-Ing. Stefan Dießloch

4

Middleware for Heterogenous and Distributed Information Systems - WS05/06

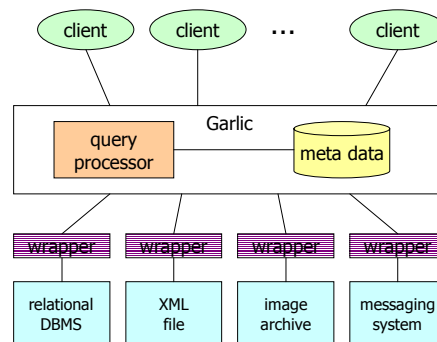
Foreign Data Wrapper

- Encapsulates a data source and 'mediates' between FDBMS and data source
- Provides infrastructure for overcoming heterogeneity among data sources:
 - wrapper architecture
 - wrapper interfaces
- Heterogeneity of data sources (needs to be overcome)
 - data model, schema (may evolve)
 - data access API
 - query capabilities
 - query language and expressiveness (simple scan, sort, simple predicates, complex predicates aggregation, binary joins, n-way joins, ...)
 - class/function libraries
 - proprietary query APIs



Garlic

- "The wrapper architecture of Garlic ... addresses the challenge of diversity by standardizing how information in data sources is described and accessed, while taking an approach to query planning in which the wrapper and the middleware dynamically determine the wrapper's role in answering a query"



M. T. Roth, P. Schwarz:
 "Don't Scrap It, Wrap It!
 A Wrapper Architecture for
 Legacy Data Sources",
 VLDB'97



Wrapper Architecture

- Garlic and wrappers cooperate for query processing
 - wrapper provides information about its processing capabilities
 - Garlic compensates for (potential) lack of wrapper functionality
- Extensibility
 - add new data sources (accessed using existing wrappers)
 - add new wrappers for supporting new types of data sources
- Wrapper evolution
 - start with simple wrappers (equivalent of a table/collection scan)
 - low cost
 - expand query processing capabilities of the wrapper until it provides full support of the data source functionality

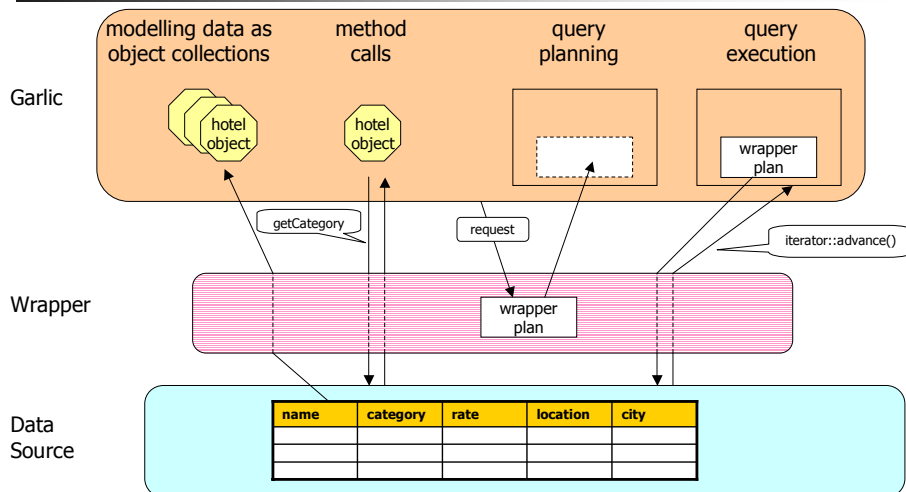


© Prof. Dr.-Ing. Stefan Dießloch

7

Middleware for Heterogenous and Distributed Information Systems - WS05/06

Wrapper-Services



© Prof. Dr.-Ing. Stefan Dießloch

8

Middleware for Heterogenous and Distributed Information Systems - WS05/06

Modelling Data as Object Collections

- Registration
 - wrapper supplies description of data source in GDL (Garlic Data Language, derived from ODMG-ODL)
 - 'global schema' at the garlic level
- Garlic object
 - interface
 - at least one implementation (multiple are possible, but only one per data source)
 - identity: OID consists of
 - IID (implementation identifier)
 - key (identifies instance within a data source)
 - root objects (collections) serve as entry into data source, can be reference using external names

Example: Travel Agency Schema

Relational Repository Schema:

```
interface Country {
    attribute string name;
    attribute string airlines_served;
    attribute boolean visa_required;
    attribute Image scene;
}
```

```
interface City {
    attribute string name;
    attribute long population;
    attribute boolean airport;
    attribute Country country;
    attribute Image scene;
}
```

Web Repository Schema:

```
interface Hotel {
    attribute readonly string name;
    attribute readonly short category;
    attribute readonly double daily_rate;
    attribute readonly string location;
    attribute readonly string city;
}
```

Image Server Repository Schema:

```
interface Image {
    attribute string file_name;
    double matches (in string file_name);
    void display (in string device_name);
}
```

Method Calls

- Method can be called by Garlic query execution engine or by the application, based on an object reference
- Methods
 - implicitly defined get/set-methods (accessor methods)
 - explicitly defined methods
- Invocation mechanisms
 - stub dispatch
 - natural if data source provides object class libraries
 - example: *display* (see previous charts)
wrapper provides routine that extracts file name from OID, receives device name as parameter, calls class library for display operation
 - generic dispatch
 - wrapper provides one entry point
 - schema-independent
 - example: relational wrapper (see previous charts)
access methods only; each call is translated into a query:
 - method name -> attribute
 - IID -> relation name
 - value -> assignment value (SET)



Query Planning

- Fundamental Idea: wrappers participate in query planning process
- Query planning steps
 - Garlic optimizer identifies for each data source the largest possible query fragment that does not reference other data sources, sends it to the wrapper
 - wrapper returns one or more query plans that can be used to process the full query fragment or parts of the query fragment
 - providing all objects in a collection is minimal requirement
 - optimizer generates alternative plans, estimates execution costs, provides for compensation of fragments not supported by the wrapper

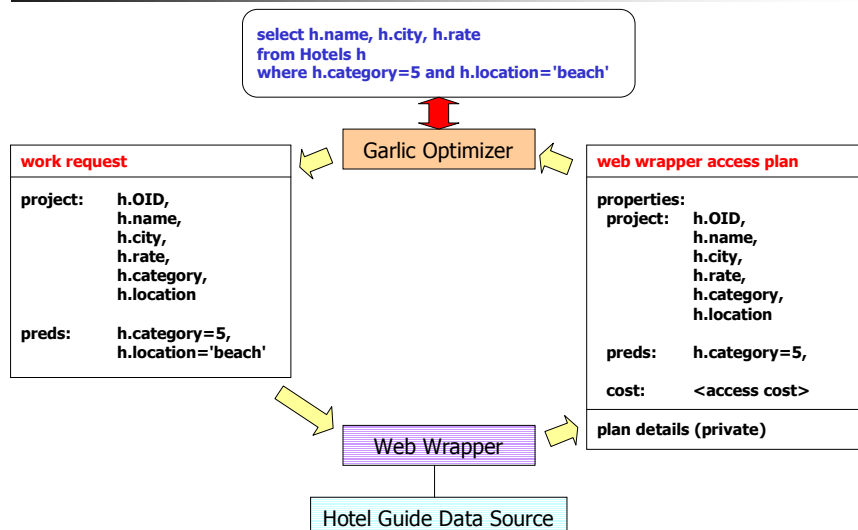


Query Planning (continued)

- Wrapper provide the following methods to be used by Garlic *work requests*:
 - *plan_access()*: generates *single-collection access plans*
 - *plan_join()*: generates *multi-way join plans* (joins may occur in application queries or in the context of resolving path expressions)
 - *plan_bind()*: generates special plan, which can be used as an *inner stream* of a *bind join*
- Result of a *work request*:
 - sets of *plans*
 - each *plan* contains a list of properties that describe which parts of the work request are implemented by the plan and what the costs are

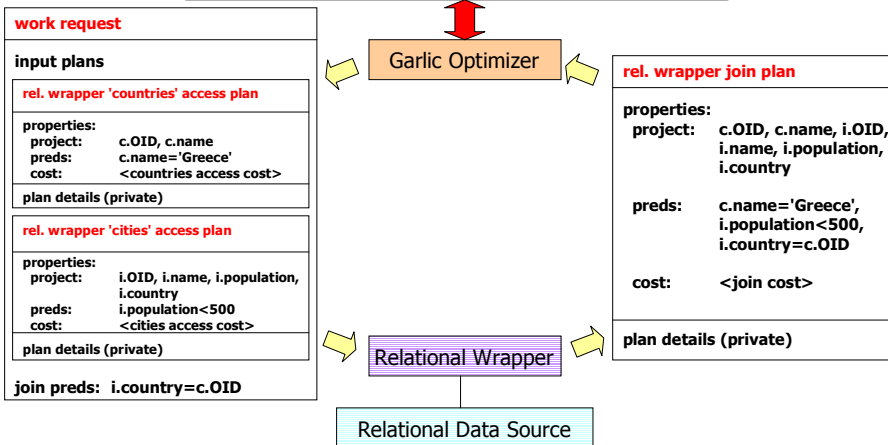


Single Collection Access Plan



Join Plan

```
select i.name
from Countries c, Cities i
where c.name='Greece' and i.population<500 and i.country = c.OID
```



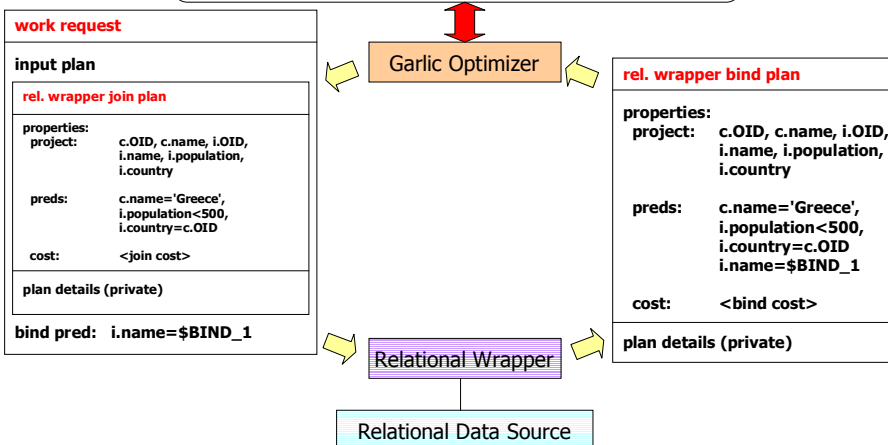
© Prof.Dr.-Ing. Stefan Deißloch

15

Middleware for Heterogenous and Distributed Information Systems - WS05/06

Bind Plan

```
select h.name, h.rate
from Hotels h, Countries, C, Cities I
where h.category=5 and h.location='beach' and c.name='Greece'
and i.population<500 and h.city=i.name and i.country=c.OID
```



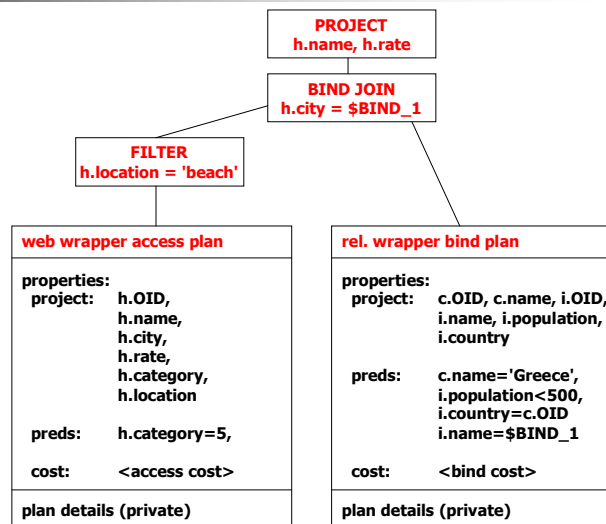
© Prof.Dr.-Ing. Stefan Deißloch

16

Middleware for Heterogenous and Distributed Information Systems - WS05/06

Wrapper Plan Synthesis

- Plan generation needs to be supported by wrapper methods
- Plan execution has to be supported by wrapper as well (Iterator methods)



© Prof.Dr.-Ing. Stefan Dießloch

17

Middleware for Heterogenous and Distributed Information Systems - WS05/06

Wrapper Packaging

- Wrapper programm provides the following wrapper components in a package:
 - interface files
 - GDL definitions
 - environment files
 - support for data-source-specific information
 - libraries
 - schema registration
 - method calls
 - query processing interfaces



© Prof.Dr.-Ing. Stefan Dießloch

18

Middleware for Heterogenous and Distributed Information Systems - WS05/06

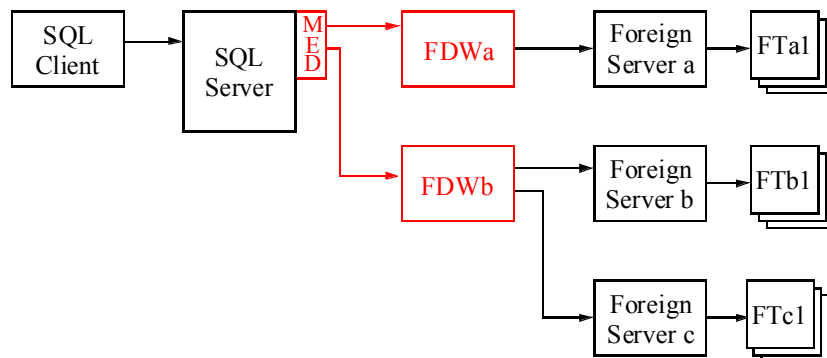
SQL/MED

- Part 9 of **SQL:1999: Management of External Data**
 - extended in SQL:2003
- Two major parts
 - Datalinks
 - Foreign Data Wrapper / Foreign Data Server



Foreign Data Wrapper/Server

- Concept based on Garlic idea
 - data provided as tables instead of object collections
- Model:



Foreign Data Server

- Manages data stored outside the SQL server
- SQL server and SQL client use foreign server descriptors (catalog elements) to communicate with foreign servers
- Catalog (implementation-specific):
 - SQL schemas
 - Foreign server descriptors
 - Foreign table descriptors
 - Foreign wrapper descriptors
- Foreign table
 - stored in a (relational) foreign server or dynamically generated by foreign wrapper capabilities
- Modes of interaction
 - *Decomposition*
 - SQL query is analyzed by SQL server, communicating with foreign data wrapper using *InitRequest*
 - *Pass-Through* (see discussion of *TransmitRequest*)



Foreign Data Wrapper Interface

- Handle routines
- Initialization routines
 - AllocDescriptor
 - AllocWrapperEnv
 - ConnectServer
 - GetOps: request meta data about
 - foreign data wrapper/server capabilities
 - foreign table (columns)
 - InitRequest: initializes processing of a request (query)
- Access routines
 - Open
 - Iterate: for delivering foreign data to SQL server
 - ReOpen
 - Close
 - GetStatistics
 - TransmitRequest: „pass-through“ of a query/request using the proprietary language of the foreign server



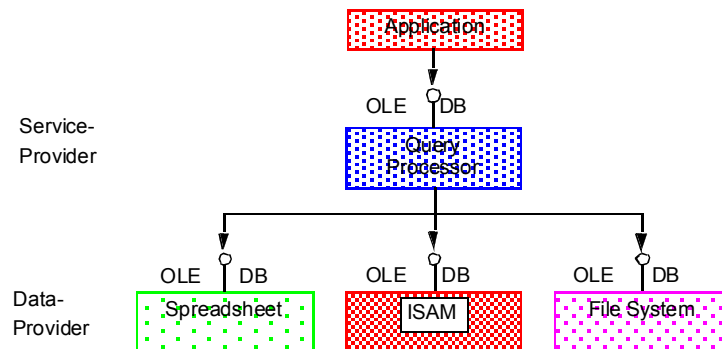
Security, Updates, and Transactions

- User Mapping
 - defines mapping of SQL server user-ids to corresponding concept of a foreign server
 - example:
 - ```
CREATE USER MAPPING FOR dssloch
SERVER myforeignserver
OPTIONS
 (user_id 'SD',
 user_pw 'secret')
```
- Updates, transactions on external data
  - not supported in SQL/MED
    - goal for future version of the standard
  - supported as product extensions in a limited form
    - usually, updates on non-relational data sources are not supported
      - distributed TAs are useful, read-only optimization can be used for foreign data source
    - updates on relational data sources
      - pass-through
      - transparent



## Microsoft OLE-DB

- Overview



J.A. Blakeley: "Universal Data Access with OLE DB",  
Proc. IEEE Comcon'97, San Jose, IEEE Computer Society Press, Feb. 1997

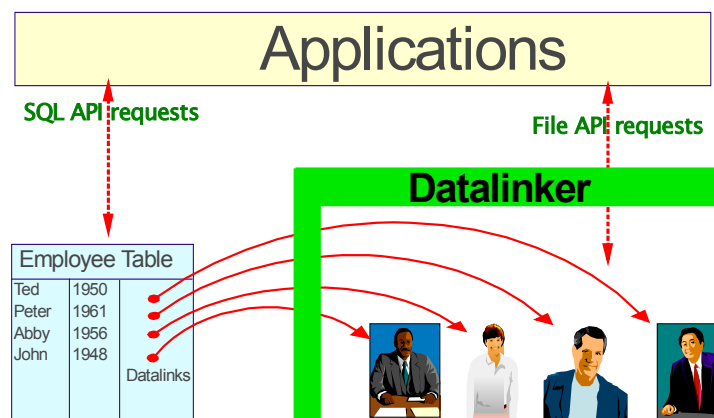


## Concepts

- Data provider
  - simple wrapper
  - encapsulates data source access
  - provides a rowset abstraction to allow iteration over a stream of data values
- Service provider
  - can provide view over heterogeneous sources by combining rowsets from different providers
    - union, join, aggregation, ...
  - special OLE-DB protocols for service provider implementation
- Recent enhancements for going beyond "flat" rowsets
- Garlic vs. OLE-DB service provider
  - Garlic queries use Object-SQL
  - Garlic wrapper and query processor interact dynamically for determining query plan



## Managing External Data: Datalinks



## DataLinks in SQL/MED

- Goal
  - preserve external storage, manipulation of files
  - synchronize integrity control, recovery, and access control of files and associated SQL data
- Concepts
  - datalink is an instance of the DATALINK data type
    - references a file (URL) that is not stored by the SQL server, but maintained by an external file manager
  - datalink options (per column)
    - define the amount of management and control the SQL server has over the datalink values of a column
      - integrity, read/write access, recovery
    - specifies the semantics of link/unlink behavior
  - datalinker
    - implementation-dependent
    - implements a number of mechanisms for guaranteeing datalink properties such as integrity control, recovery, access control



## Functions and Operations

- New SQL functions for datalinks
  - constructor: DLVALUE, ...
  - (components of) URLs: DLURLCOMPLETE, ...
- SQL statements (examples)
  - insert ("link")

```
INSERT INTO Movies (Title, Minutes, Movie)
VALUES ('My Life', 126,
DLVALUE('http://my.server.de/movies/mylife.avi'))
```
  - select (incl. URL access token)

```
SELECT Title, DLURLCOMPLETE(Movie)
FROM Movies
WHERE Title LIKE '%Life%'
```



## Data Link Options

- Link control (NO, FILE)
  - NO LINK CONTROL
    - URL-Format of datalink
    - no further control, file is not "linked"
  - FILE LINK CONTROL
    - file is "linked", file has to exist!
    - level of control can be specified using further options
- Integrity control option (ALL, SELECTIVE, NONE)
  - INTEGRITY ALL
    - linked files cannot be deleted or renamed
  - INTEGRITY SELECTIVE
    - linked files can only be deleted or modified using file manager operations, if no datalinker is installed
  - INTEGRITY NONE
    - referenced files can be deleted or modified using file manager operations
      - not compatible with FILE LINK CONTROL



## Data Link Options (continued)

- Read permission option (FS, DB)
  - READ PERMISSION FS
    - read access is determined by file manager
  - READ PERMISSION DB
    - read access is controlled by SQL server, based on access privileges to the datalink value
    - involves read access tokens
      - encoded into the URL by the SQL server
      - verified by external file manager/data linker
- Write permission option (FS, ADMIN, BLOCKED)
  - WRITE PERMISSION FS
    - write access controlled by file manager
  - WRITE PERMISSION BLOCKED
    - linked files cannot be modified
  - WRITE PERMISSION ADMIN [NOT] REQUIRING TOKEN FOR UPDATE
    - write access governed by SQL server (and datalinker)
      - requires READ PERMISSION DB
    - involves write access token for modifying file content
      - may have to be presented to the SQL server again



## Functions and Operations (continued)

- “Update-in-place”

```
SELECT Title, DLURLCOMPLETWRITE(Movie)
INTO :t, :url ...
```

open using URL, modify ...

```
UPDATE Movies SET Movie = DLNEWCOPY(:url, 1)
WHERE Title = :t
```

- DLNEWCOPY

- indicates to the SQL server that the file content has changed and should be managed appropriately
- alternative: DLPREVIOUSCOPY – file content may have changed, but the application is not interested in keeping the changes, original file is restored



## Data Link Options (continued)

- RECOVERY YES/NO

- indicates whether SQL server coordinates recovery (jointly with datalinker) or not

- Unlink option (RESTORE, DELETE, NONE)

- ON UNLINK RESTORE
  - original properties (ownership, permissions) restored as well
- ON UNLINK DELETE
  - file is deleted when unlinked
- ON UNLINK NONE
  - ownership and permissions are not restored

- SQL statement (example)

- “Unlink/Replace”

```
UPDATE Movies SET Movie =
DLVALUE('http://my.newserver.de/mylife.avi')
WHERE Title = 'My Life'
```

RESTORE or DELETE for “.../movies/mylife.avi”





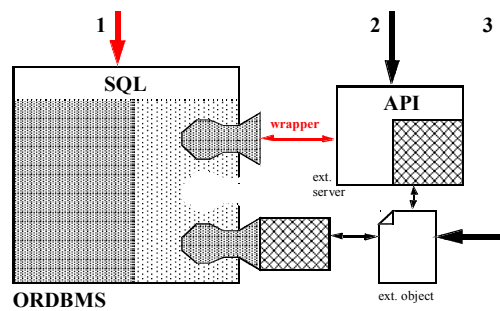
## Valid Combinations

| Integrity | Read permission | Write permission | Recovery | Unlink  |
|-----------|-----------------|------------------|----------|---------|
| ALL       | FS              | FS               | NO       | NONE    |
| ALL       | FS              | BLOCKED          | NO       | RESTORE |
| ALL       | FS              | BLOCKED          | YES      | RESTORE |
| ALL       | DB              | BLOCKED          | NO       | RESTORE |
| ALL       | DB              | BLOCKED          | NO       | DELETE  |
| ALL       | DB              | BLOCKED          | YES      | RESTORE |
| ALL       | DB              | BLOCKED          | YES      | DELETE  |
| ALL       | DB              | ADMIN            | NO       | RESTORE |
| ALL       | DB              | ADMIN            | NO       | DELETE  |
| ALL       | DB              | ADMIN            | YES      | RESTORE |
| ALL       | DB              | ADMIN            | YES      | DELETE  |
| SELECTIVE | FS              | FS               | NO       | NONE    |



## Comparing Wrappers and Data Links

- Choices for accessing/manipulating external data

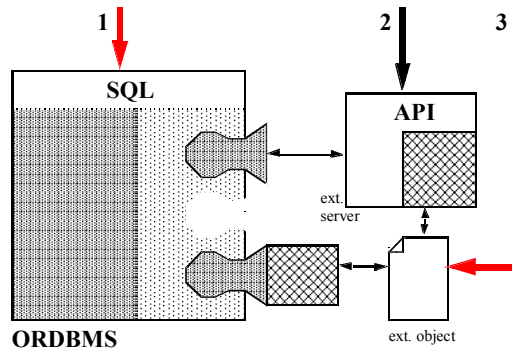


- For foreign data wrapper/server
  - external object = foreign tables/values
  - manipulation of data: 1, 2
  - wrapper implements routines needed to interact with SQL server for delegation of query fragments to external server



## Comparing Wrappers and Data Links (cont.)

- For datalinks:
  - external object = referenced file
  - 1: manipulation of URLs (datalink values), obtaining permissions (token)
  - 3: "overloaded" file system access
    - access and integrity control through ORDBMS
  - Synchronization of recovery- and backup-mechanisms implementation-dependent, not defined by standard



## Summary

- Wrapper as a mediator between data source and federated DBMS middleware
  - Example: Garlic (IBM)
    - almost any data source can be integrated
    - global query optimization
      - middleware (Garlic) and wrapper decide dynamically which query fragments are processed by the wrapper
      - specific capabilities of data sources can be utilized
  - SQL/MED
    - part 9 of SQL:1999 Standards
    - follows the Garlic idea
    - Foreign-Data-Wrapper/Server
  - Advantages
    - overcomes heterogeneity regarding data model, API
    - location transparency
      - global schema, query can reach across multiple, distributed data sources
  - Limitations
    - structural and schematic heterogeneity remain to be addressed
    - no standardized update operations
- Datalinks as a means for supporting referential integrity for external files, synchronizing recovery and access control of files and SQL data