

## Kapitel 7

# Wrapper-basierte Datenintegration

### Inhalt:

- Einführung: Wrapper
- Garlic
- SQL/MED (Management of External Data)
- Überblick: OLE (Microsoft)
- Zusammenfassung



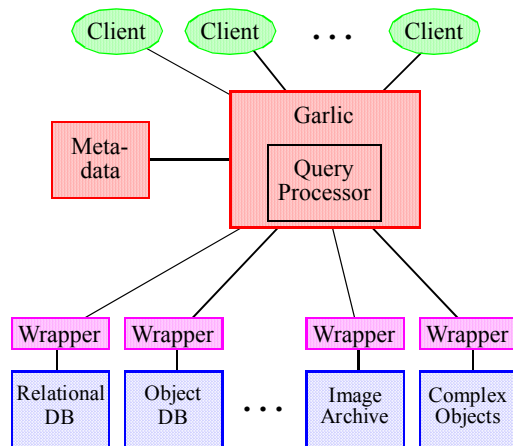
## Einführung: Wrapper

- kapselt die in einer Datenquelle gespeicherten Daten und 'vermittelt' zwischen der Datenquelle und der Middleware
- entscheidend hinsichtlich der Bewältigung der Heterogenität zwischen den Datenquellen:
  - Wrapper-Architektur
  - Wrapper-Interfaces
- (zu bewältigende) Heterogenität, jede Datenquelle hat
  - ihr eigenes Datenmodell und Schema, wobei letzteres sich mit der Zeit ändern kann
  - ihr eigenes API
  - und spezielle Anfragemöglichkeiten
    - Anfragesprache  
(Spektrum: einfache Scans, Sortierung, einfache Prädikate, komplexe Prädikate, Aggregation, binäre Joins, n-Wege-Joins; evtl. sogar 'eigentümlich': spezielle Prädikate nur auf bestimmten Daten erlaubt)
    - Klassen-/Funktionsbibliothek
    - spezifische Programmierschnittstelle



## 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



© Prof.Dr.-Ing. Stefan Deßloch

3

Middleware for Heterogenous and  
Distributed Information Systems -  
WS04/05

## Wrapper-Architektur

- Ziele
  - Geringe Start-up-Kosten für die Erstellung eines Wrappers
  - Möglichkeit der Evolution des Wrappers
  - Flexible Architektur; Möglichkeit des ‘Einhängens’ weiterer Datenquellen (werden in Garlic auch Repositories genannt) und zugehöriger Wrapper
  - Wrapper sollte der Garlic-Anfrageverarbeitung exakt die Anfragemöglichkeiten verfügbar machen, die die zugehörige Datenquelle unterstützt

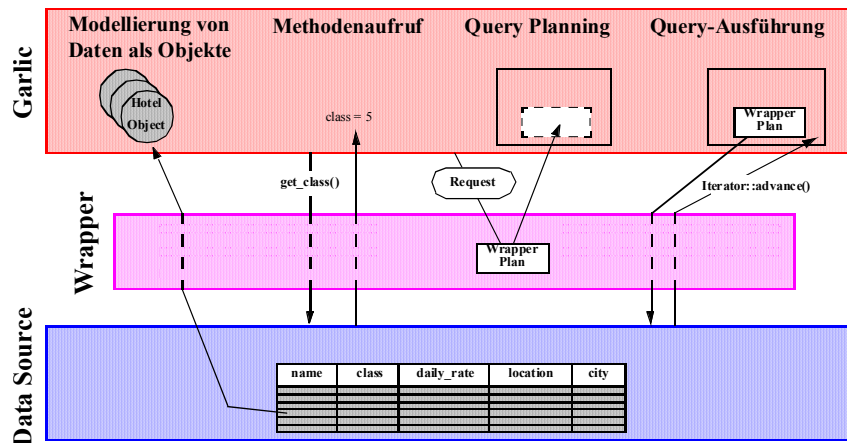


© Prof.Dr.-Ing. Stefan Deßloch

4

Middleware for Heterogenous and  
Distributed Information Systems -  
WS04/05

# Wrapper-Services



# Modellierung von Daten als Objekte

- Registrierung
  - Wrapper liefert Garlic Beschreibung der Datenquelle in GDL (Garlic Data Language, Variante von ODMG-ODL)
  - 'globales Schema' auf der Garlic-Ebene
- Garlic-Objekt:
  - Schnittstelle
  - mind. eine Implementierung (evtl. mehrere, höchstens eine pro Datenquelle)
  - Identität: OID besteht aus
    - IID (implementation identifier)
    - key (identifiziert die Ausprägung innerhalb der Datenquelle)
  - Root-Objekte (üblicherweise Kollektionen) dienen dem Einstieg und können auch über einen externen Namen angesprochen werden



## Beispiel: Schema einer Reiseagentur

### 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 class;  
    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);  
}
```



## Methodenaufruf

- Methodenaufrufe können durch die Garlic-Ausführungsmaschine bzw. durch eine Garlic-Applikation, die eine Referenz auf ein entsprechendes Repository-Objekt erhalten hat, veranlasst werden
- Methoden
  - implizit definierte *Get/Set*-Methoden (accessor methods)
  - explizit definierte Methoden



## Arten des Methodenaufrufs

- stub dispatch
  - natürlicher Mechanismus für Repositories, deren Schnittstelle als Klassenbibliothek gegeben ist
  - Beispiel: *display* (siehe oben)  
Wrapper stellt Routine zur Verfügung, die aus dem key-Feld der OID den Dateinamen des angesprochenen Bildes und aus den von Garlic übergebenen Parametern den Gerätenamen extrahiert; zur Wiedergabe wird dann die zugehörige Funktion der Klassenbibliothek aufgerufen;
- generischer Dispatch
  - Wrapper bietet nur genau einen Methodeneinsprungpunkt
  - schema-unabhängig
  - Beispiel: relationaler Wrapper (siehe oben)  
nur accessor methods; jeder Aufruf wird in Anfrage umgesetzt; Methodenname → Attributname,  
IID → Relationenname,  
Werte → Zuweisung (Set);

## Query Planning

- Grundlegende Idee
  - Wrapper nehmen an der Erstellung des Query-Plans teil
- Allgemeiner Ablauf:
  - Garlic-Optimierer identifiziert das größtmögliche, ein Repository betreffende Query-Fragment und schickt es zum zugehörigen Wrapper
  - Wrapper liefert einen oder mehrere Pläne zurück, die einen Teil der durch das Query-Fragment angesprochenen Arbeit bzw. die gesamte Arbeit ausführen können
  - Garlic-Optimierer erstellt und beurteilt die möglichen Pläne zur Abarbeitung der gesamten Query; dabei versucht er, für die Teile, die ein Wrapper nicht übernehmen konnte, alternative Operatoren zu finden

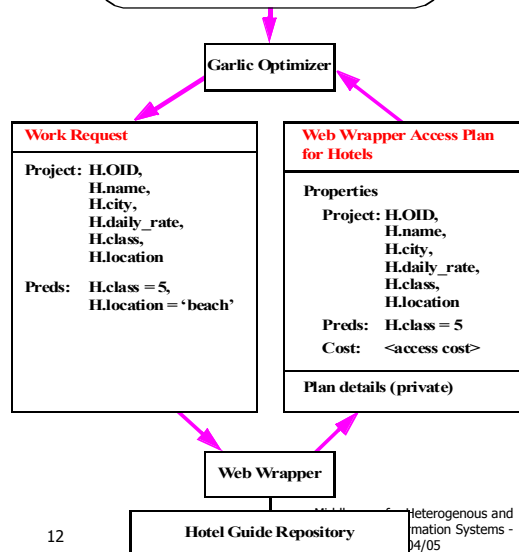
## Query Planning (Forts.)

- Wrapper stellt zu diesem Zweck folgende Methoden bereit, die durch *work requests* von Garlic angesprochen werden können:
  - *plan\_access()*: generiert *single-collection access plans*
  - *plan\_join()*: generiert *multi-way join plans* (Joins können in Anwenderanfragen vorkommen oder von Garlic zur Auflösung von Pfadausdrücken erzeugt werden)
  - *plan\_bind()*: generiert speziellen Plan, der als *inner stream* eines *bind joins* benutzt werden kann
- Ergebnis eines *work requests*:
  - Menge von *plans*
  - jeder *plan* beinhaltet eine Liste von Eigenschaften, die beschreiben, welche Teile des work requests durch den plan implementiert werden und zu welchen Kosten



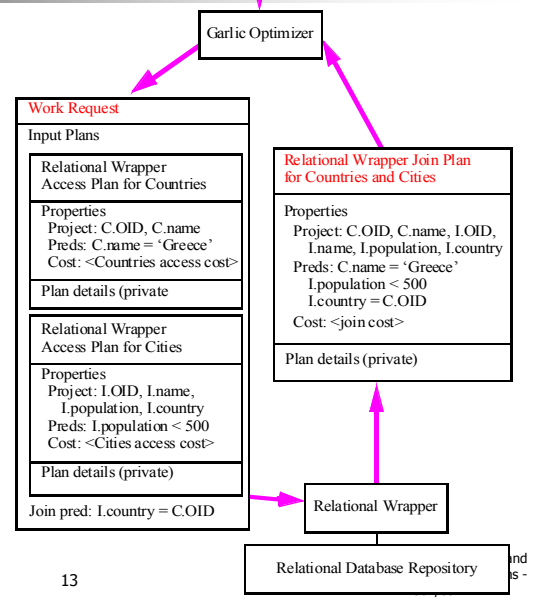
## Single Collection Access Plan

```
select H.name, H.city, H.daily_rate
from Hotels H
where H.class = 5 and H.location = 'beach'
```



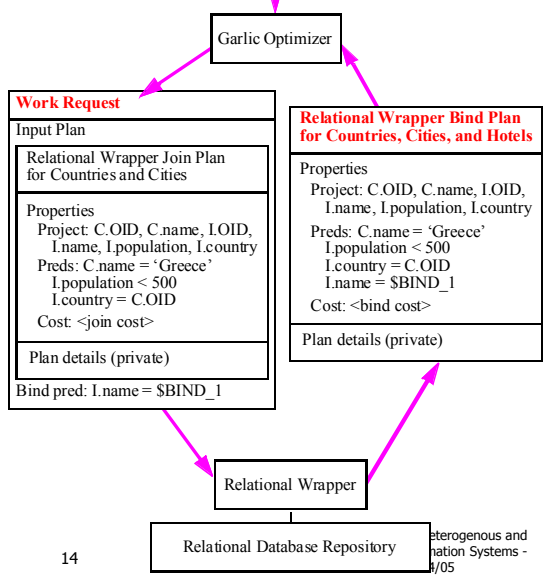
# Join Plan

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



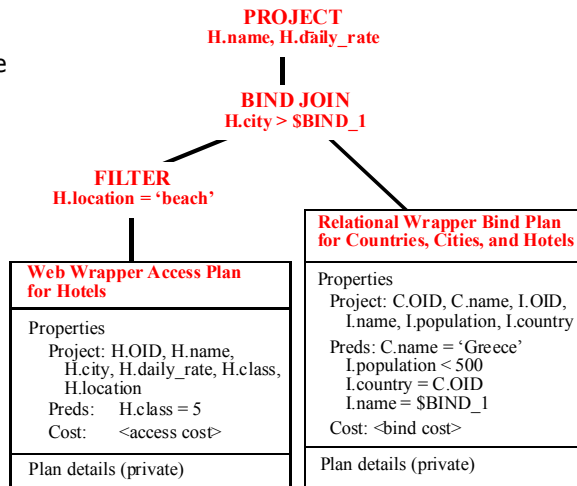
# Bind Plan

```
select H.name, H.daily_rate
from Hotels H, Countries C, Cities I
where H.class = 5 and H.location = 'beach' and C.name = 'Greece'
and I.population < 500 and H.city = I.name and I.country = C.OID
```



## Synthese der Wrapper Plans

- Übersetzung muss unterstützt werden durch entsprechende Wrapper-Methoden
- Ausführung wird ebenfalls durch Wrapper unterstützt (Iterator-Methoden)



## Wrapper Packaging

- Wrapper-Programmierer stellt folgende Wrapper-Komponenten in einem Wrapper-Package zur Verfügung:
  - Interface Files
    - GDL-Definitionen
  - Environment Files
    - Kodierung Repository-spezifischer Information
  - Libraries (dynamisch ladbar)
    - Schema-Registrierung
    - Methodenaufruf
    - Anfrageschnittstellen





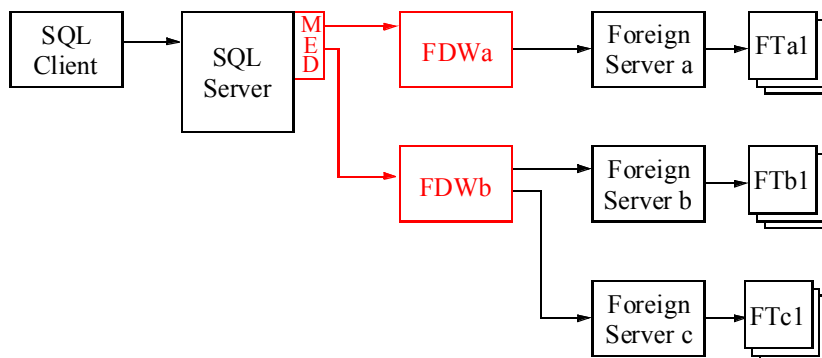
# SQL/MED

- Part 9 von **SQL:1999: Management of External Data**
  - in SQL:2003 weiterentwickelt
- zwei wesentliche Teile
  - Datalinks
  - Foreign Data Wrapper / Foreign Data Server



# Foreign Data Wrapper/Server

- Konzept basiert auf Garlic-Idee
- Modell



## Foreign-Data-Server

- verwaltet außerhalb der eigentlichen SQL-Umgebung liegende Daten
- SQL-Server und SQL-Client nutzen Foreign-Server-Descriptors (Katalogelemente), um mit dem Foreign-Server kommunizieren zu können
- Katalog (implementierungsspezifisch):
  - SQL-Schemata
  - Foreign-Server-Deskriptors
  - Foreign-Table-Descriptors
  - Foreign-Wrapper-Descriptors
- Foreign-Table
  - tatsächlich im Foreign-Server (relational) liegende Tabelle oder durch Wrapper-Funktionalität dynamisch zu erzeugende Tabelle
- (Interaktions-)Modi
  - *Decomposition*
    - Analyse der SQL-Anfrage durch SQL-Server und Kommunikation mit dem Foreign-Data-Wrapper über InitRequest
  - *Pass-Through* (siehe Beschreibung TransmitRequest)



## Foreign-Data-Wrapper-Interface

- Handle-Routinen
- Initialisierungsroutinen
  - AllocDescriptor
  - AllocWrapperEnv
  - ConnectServer
  - GetOps: Abfragen von Möglichkeiten des Foreign-Data-Wrappers/Servers, einer (Spalte einer) Foreign-Table
  - InitRequest: Initiierung der Vorbereitung einer Anfrage
- Zugriffsroutinen
  - Open
  - Iterate: Übertragung von Tupeln
  - ReOpen
  - Close
  - GetStatistics
  - TransmitRequest: „Durchreichen“ einer Anfrage in der Sprache des Foreign-Data-Servers (Pass-Through-Modus)



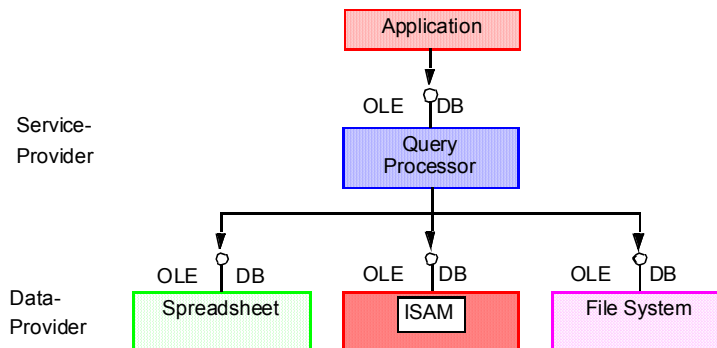
# Sicherheit, Updates und Transaktionen

- User Mapping
  - Definiert Abbildung einer User-ID des SQL-Servers auf ein entsprechendes Konzept des Foreign Server
  - Beispiel:
    - `CREATE USER MAPPING FOR desloch  
SERVER myforeignserver  
OPTIONS  
(user_id 'SD',  
user_pw 'secret')`
- Updates, Transaktionen auf Externen Daten
  - in SQL/MED nicht unterstützt
    - zukünftige Entwicklung
  - als produktspezifische Erweiterungen in eingeschränkter Form verfügbar
    - üblicherweise Änderungen auf nicht-relationalen Datenquellen nicht möglich
      - verteilte TA sinnvoll, nur 1PC (read-only) für foreign data source
    - Änderungsoperationen auf relationalen Datenquellen
      - pass-through
      - transparent



# Microsoft OLE-DB

- Überblick



J.A. Blakeley: "Universal Data Access with OLE DB",  
Proc. IEEE Comcon'97, San Jose, IEEE Computer Society Press, Feb. 1997  
Middleware for Heterogenous and  
Distributed Information Systems -  
WS04/05

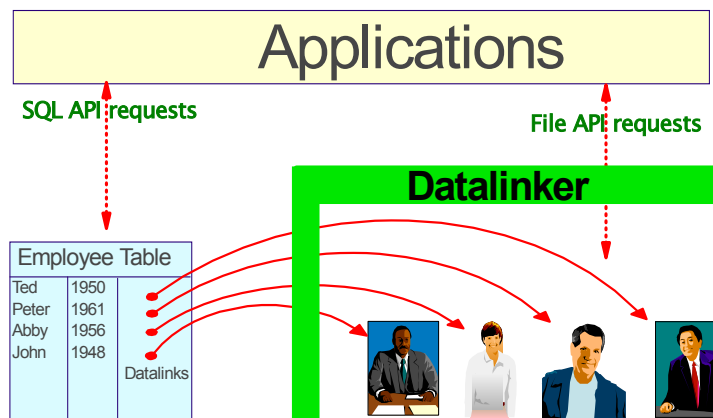


## Eigenschaften

- ein *Data-Provider* (einfacher Wrapper) kapselt den Zugriff auf Datenquelle und bietet als Abstraktion einen rowset als Strom von Datenwerten mit einem Iterator
- Daten, die mehrere Data-Provider in Tabellenform verfügbar machen, können durch einen *Service-Provider* zu einer heterogenen Sicht verknüpft werden (Union, Join, Agg. etc.); zum Erstellen eines Service-Providers bietet OLE DB spezielle Protokolle; die Middleware-Komponente ist hier wenig generisch
- neuere Versionen sollen komplexere Abstraktionen als flache Tabellen, z. B. objektorientierte und semistrukturierte Daten erlauben
- Unterschiede zu Garlic
  - Garlic-Anfragefragmente sind in Object-SQL formuliert
  - Garlic-Wrapper kann dynamisch seine Beiträge ermitteln



## Managing External Data: Datalinks



## DataLinks in SQL/MED

- Konzept
  - ein Datalink ist eine Ausprägung des DATALINK-Datentyps
  - ein Datalink referenziert eine Datei (URL), die nicht Teil der SQL-Umgebung ist, sondern durch einen externen File-Manager verwaltet wird
  - Datalink Type Descriptor, Link-Control-Options:
    - link control (NO, FILE)
    - integrity control option (ALL, SELECTIVE, NONE)
    - read permission option (FS, DB)
    - write permission option (FS, ADMIN, BLOCKED)
    - recovery option (NO, YES)
    - unlink option (RESTORE, DELETE, NONE)
  - Datalinker
    - implementierungsabhängig
    - Menge der Mechanismen, die für durch Datalinks repräsentierte Dateien Integritätskontrolle, Recovery und Zugriffskontrolle durchführen



## Funktionen und Operationen

- Neue SQL-Funktionen für Datalinks
  - Konstruktor: DLVALUE, ...
  - (Komponenten von) URLs: DLURLCOMPLETE,
- SQL Operationen (Beispiele)
  - Einfügen ("Link")

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

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



## Funktionen und Operationen (Forts.)

- SQL Operationen (Beispiele)
  - "Unlink/Replace"

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

RESTORE oder DELETE für ".../movies/mylife.avi"
  - "Update-in-place"

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

Öffnen über URL, Modifizieren ...

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



## Link-Control-Options

- NO LINK CONTROL
  - URL-Format des Datalinks
  - keine weitere Kontrolle
- FILE LINK CONTROL
  - existierende Datei muss referenziert werden
  - Art der Kontrolle durch die weiteren Optionen bestimmt
- INTEGRITY ALL
  - referenzierte Dateien können nur über SQL gelöscht oder umbenannt werden
- INTEGRITY SELECTIVE
  - referenzierte Dateien können mittels File-Manager-Operationen gelöscht oder umbenannt werden, solange kein Datalinker vorhanden ist
- INTEGRITY NONE
  - referenzierte Dateien können ausschließlich mittels File-Manager-Operationen gelöscht oder umbenannt werden
  - nicht verträglich mit FILE LINK CONTROL



## Link-Control-Options (Forts.)

- READ PERMISSION FS
  - Leserecht für referenzierte Dateien wird durch den File-Manager bestimmt
- READ PERMISSION DB
  - Leserecht für referenzierte Dateien wird über SQL bestimmt
- WRITE PERMISSION FS
  - Schreibrecht für referenzierte Dateien wird durch den File-Manager bestimmt
- WRITE PERMISSION BLOCKED
  - kein Schreibzugriff auf referenzierte Dateien, es sei denn, es existiert implementierungsabhängiger Mechanismus
- WRITE PERMISSION ADMIN [NOT] REQUIRING TOKEN FOR UPDATE
  - Schreibrecht für referenzierte Dateien durch SQL bestimmt



## Link-Control-Options (Forts.)

- RECOVERY YES
  - mit SQL-Server koordinierte Recovery (Datalinker-Mechanismus)
- RECOVERY NO
  - keine Recovery auf referenzierten Dateien
- ON UNLINK RESTORE
  - vor der Herstellung des Links bestehende Rechte (Ownership, Permissions) werden durch den File-Manager bei Auflösung des Links (Unlink) wiederhergestellt
- ON UNLINK DELETE
  - Löschung bei Unlink
- ON UNLINK NONE
  - keine Auswirkungen auf die Rechte bei Unlink



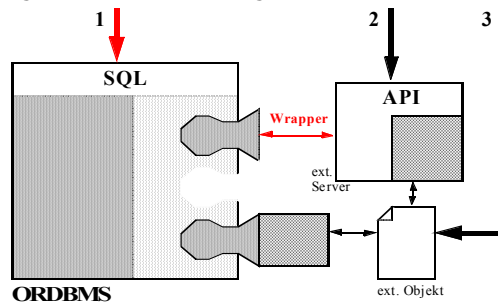
## Gültige Kombinationen

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



## Art der Anbindung über Foreign-Data-Wrapper

- allg. Möglichkeiten der Anbindung externer Daten



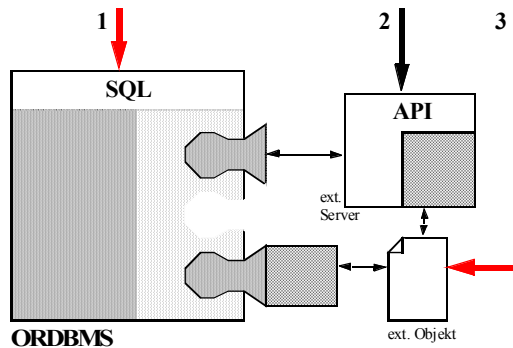
- Einordnung von Foreign-Data-Wrapper/Server
  - externe Objekte = Foreign-Tables/Values
  - Manipulation der Daten: 1, 2
  - Wrapper implementiert vom SQL-Server erwartete Routinen zur Delegation von Anfrage(teile)n an externen Server





## Art der Anbindung mittels Datalinks

- allg. Möglichkeiten der Anbindung externer Daten
- Einordnung von Datalinks:
  - externes Objekt = referenzierte Datei
  - 1: Manipulation von URLs (Datalink-Werte), „Beschaffung“ von Berechtigungen (Token)
  - 3: „überladene“ BS-Zugriffe
    - Rechte- und Integritätskontrolle durch ORDBVS
  - „Einklinken“ von Recovery- und Backup-Mechanismen implementierungsabhängig und durch Standard nicht geregelt



## Zusammenfassung

- Wrapper als Mediator zwischen Datenquelle und Middleware
- Beispiel: Garlic (IBM)
  - nahezu jede Art von Datenquelle integrierbar
  - globale Anfrageoptimierung
    - Middleware (Garlic) und Wrapper entscheiden dynamisch, welche Teile der Anfrageverarbeitung durch Wrapper übernommen werden kann, wie diese ausgeführt werden (Planalternativen)
    - spezifische Möglichkeiten der Datenquellen können genutzt werden
- SQL/MED
  - Part 9 des SQL:1999-Standards
  - folgt der Garlic-Idee
  - Foreign-Data-Wrapper/Server
- Vorteile
  - Bewältigung der Heterogenität (bzgl. DM, API)
  - Verteilungstransparenz
    - globales Schema, Anfrage kann mehrere, verteilte Datenquellen überspannen
- Einschränkungen
  - Überwindung von struktureller und semantischer Heterogenität noch notwendig
  - Änderungsoperationen auf externen Datenquellen nicht standardisiert

