

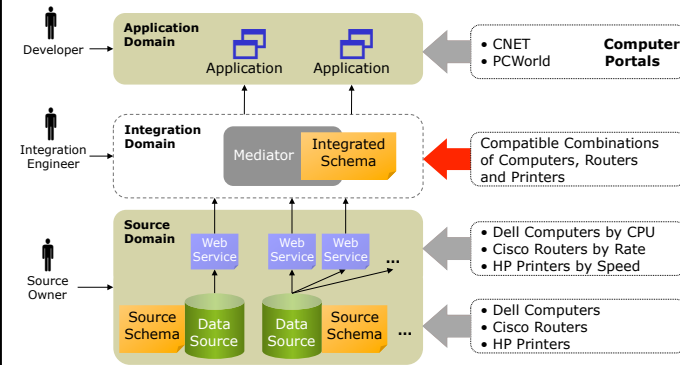
# Exporting and Interactively Querying Web Service-Accessed Sources: The CLIDE System

Michalis Petropoulos



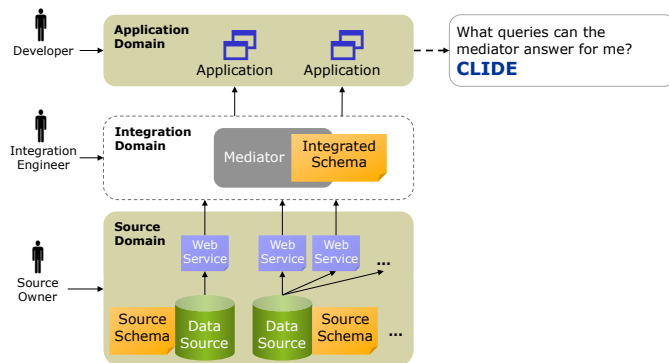
Database Seminar, February 2010

## Large-Scale Data Integration Systems



2

## Large-Scale Data Integration Systems



3

## Running Example

### Parameterized Views

Dell	Cisco
<b>Schema</b> <b>Computers</b> (cid, cpu, ram, price) <b>NetCards</b> (cid, rate, standard, interface)	<b>Schema</b> <b>Routers</b> (rate, standard, price, type)
<b>Views</b> <b>V1 ComByCpu(cpu) → (Computer)*</b> SELECT DISTINCT Com1.* FROM Computers Com1 WHERE Com1.cpu= <i>cpu</i>	<b>Views</b> <b>Wired() → (Router)*</b> SELECT DISTINCT Rou1.* FROM Routers Rou1 WHERE Rou1.type='Wired'
<b>V2 ComNetByCpuRate(cpu, rate) → (Computer, NetCard)*</b> SELECT DISTINCT Com1.*, Net1.* FROM Computers Com1, Network Net1 WHERE Com1.cid=Net1.cid AND Com1.cpu= <i>cpu</i> AND Net1.rate= <i>rate</i>	<b>Wireless Routers</b> <b>Wireless() → (Router)*</b> SELECT DISTINCT Rou1.* FROM Routers Rou1 WHERE Rou1.type='Wireless'

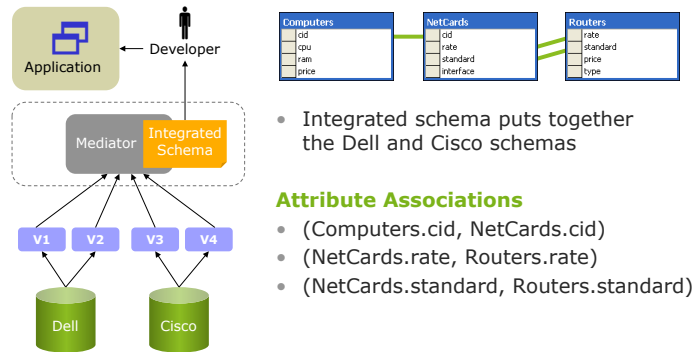
**Conjunctive Queries CQ**

- Equality & Comparison Conditions
- Parameters

4

## Running Example

### Integrated Schema



5

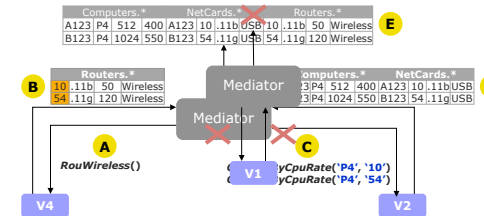
## Sophisticated Mediators Make Feasible Queries Hard to Predict

### Feasible Queries FQ

- Equivalent CQ query rewritings using the views
- Might involve more than one views
- Order might matter

Query: **Feasible**

Get all **Computers**, together with their **NetCards** and their compatible **'Wireless' Routers**



6

## Problem

1. Large number of sources
2. Large number of views (web-services)
3. Mediator capabilities

Developer formulates an application query

→ Is an application query feasible?

→ If not, how do I know which ones are feasible?

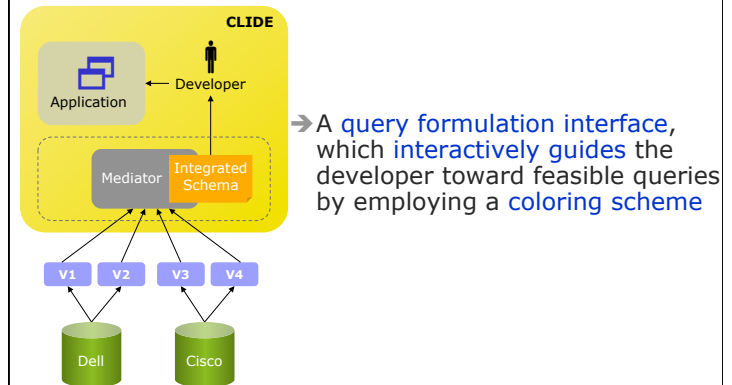
Previous options:

- The developer had to browse the view definitions and somehow formulate a feasible query
- Or formulate queries until a feasible one is found (trial-and-error)

**No system-provided guidance**

7

## The CLIDE Solution



8

## QBE-Like Interfaces

### Microsoft SQL-Server

The screenshot shows the 'Data in Table' window for a query. It displays two tables, 'Computers' and 'NetCards', with their respective columns. A join is established between the 'cid' columns of both tables. The 'Criteria' column shows the join condition:  $Computers.cid = NetCards.cid$ . The 'Output' column shows the selected columns:  $Computers.ram, Computers.price, NetCards.interface$ . The 'Criteria' column also includes a feasibility flag:  $= 'P4'$ . The SQL query is displayed at the bottom:

```
SELECT Computers.ram, Computers.price, NetCards.interface
FROM Computers INNER JOIN
NetCards ON Computers.cid = NetCards.cid
WHERE (Computers.cpu = 'P4') AND (NetCards.rate = '54Mbps')
```

9

## CLIDE Interface

The screenshot shows the CLIDE interface with a query design grid. The grid has columns for 'Table Alias', 'Selection Boxes', 'Table Boxes', 'Projection Box', and 'Feasibility Flag'. The 'Table Alias' column contains 'Computers1' and 'NetCards1'. The 'Selection Boxes' column contains 'cid ='. The 'Table Boxes' column contains 'Computers' and 'NetCards'. The 'Projection Box' column contains 'ram =', 'price =', and 'interface ='. The 'Feasibility Flag' column contains '= "P4"'. The SQL query is displayed at the bottom:

```
SELECT DISTINCT Computers1.ram, Computers1.price, NetCards1.interface
FROM Computers Computers1, NetCards NetCards1
WHERE NetCards1.cid = Computers1.cid AND Computers1.cpu = "P4" AND NetCards1.rate = "54"
```

- Table, selection, projection and join actions
- Feasibility Flag
- Color-based suggestions

10

## Example Interaction

### Snapshot 1

The screenshot shows the CLIDE interface with a query design grid. The 'cpu = "P4"' action is highlighted in yellow. The SQL query is displayed at the bottom:

```
INFEASIBLE
SELECT DISTINCT Computers1.ram, Computers1.price
FROM Computers Computers1
```

- Yellow** → required action
- All feasible queries require this action
- White** → optional action
- Feasible queries can be formulated w/ or w/o these actions

11

## Example Interaction

### Snapshot 2

The screenshot shows the CLIDE interface with a query design grid. The 'cpu = "P4"' and 'price = 400' actions are highlighted in blue. The SQL query is displayed at the bottom:

```
FEASIBLE
SELECT DISTINCT Computers1.ram, Computers1.price
FROM Computers Computers1
WHERE Computers1.cpu = "P4"
```

- Blue** → required choice of action
- At least on **Mediator** query cannot be formulated unless this **query** is performed
- Diagram illustrating the Mediator concept:
- ```

    graph TD
      A[ComByCpu("P4")] --> V1[V1]
      B[ram price 1024 550] --> V1
      V1 --> Mediator[Mediator]
      Mediator --> C[query cannot be formulated]
  
```
- Table 1:
- | cid  | cpu | ram  | price |
|------|-----|------|-------|
| A123 | P4  | 512  | 400   |
| B123 | P4  | 1024 | 550   |
- Table 2:
- | cid  | cpu | ram  | price |
|------|-----|------|-------|
| A123 | P4  | 512  | 400   |
| B123 | P4  | 1024 | 550   |

12

## Example Interaction

Snapshot 3

### Join Lines:

- Only yellow and blue are displayed
- Must appear in Attribute Associations

13

## Example Interaction

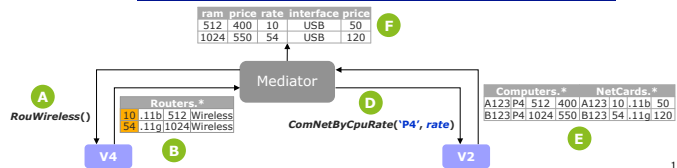
Snapshot 4

- \* → any other constant
- Red → prohibited action
  - Does not appear in any feasible query
  - Lead to "Dead End" state

14

## Example Interaction

Snapshot 5



15

## Demo

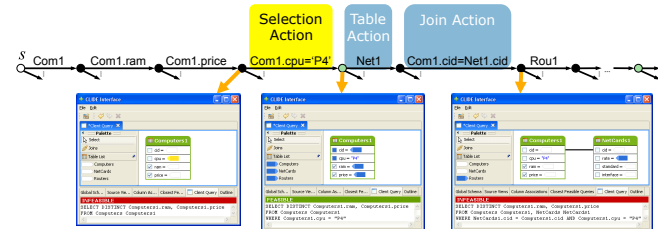
16

## CLIDE Properties

- **Completeness of Suggestions**
  - Every feasible query can be formulated by performing yellow and blue actions at every step
- **Summarization of Suggestions**
  - At every step, only a minimal number of actions is suggested, i.e., the ones that are needed to preserve completeness
- **Rapid Convergence By Following Suggestions**
  - The shortest sequence of actions from a query to any feasible query consists of suggested actions

17

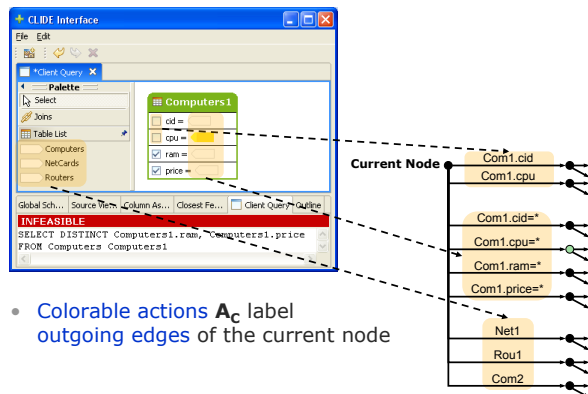
## Interaction Graph



- Nodes are queries: One for each  $q \in CQ$
- Edges are actions: Table, selection, projection and join actions
- Green nodes are feasible queries
- Infinitely big structure
  - All CQ queries
  - All possible combinations of actions formulating them

18

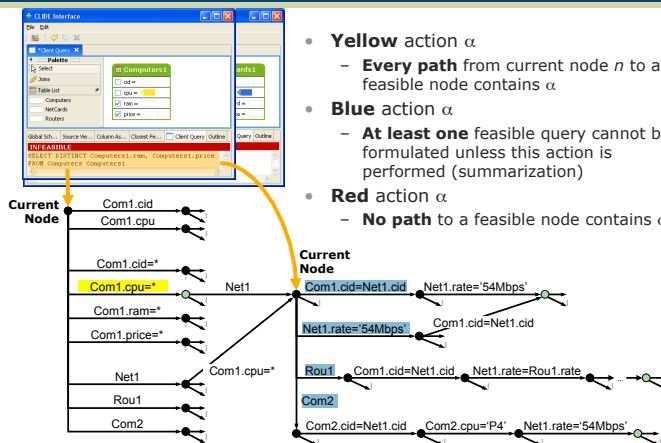
## Interaction Graph: Colorable Actions



- Colorable actions  $A_c$  label outgoing edges of the current node

19

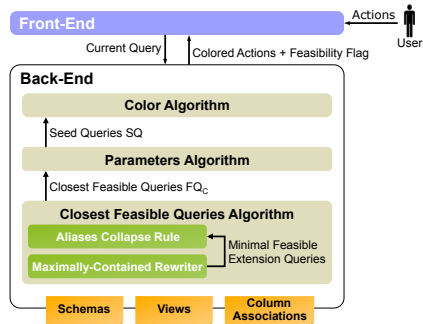
## Interaction Graph: Colors



- **Yellow** action  $\alpha$ 
  - **Every** path from current node  $n$  to a feasible node contains  $\alpha$
- **Blue** action  $\alpha$ 
  - **At least one** feasible query cannot be formulated unless this action is performed (summarization)
- **Red** action  $\alpha$ 
  - **No path** to a feasible node contains  $\alpha$

20

## CLIDE Architecture

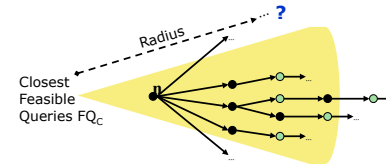


- Back-End invoked every time the user performs an action
  - i.e., the user arrives at a new node in the interactions graph

21

## Color Determined By a Finite Set of Feasible Queries

**Challenge: Infinitely Many Feasible Queries**



**Solution: Closest Feasible Queries  $FQ_C$**

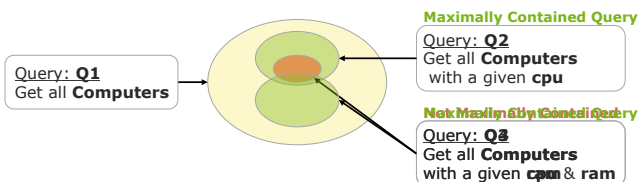
- $FQ_C$  is sufficient to color actions in  $A_C$
- **Theorem: Set of Closest Feasible Queries is Finite**

**Challenge: How far can the Closest Feasible Queries  $FQ_C$  be?**

**Solution: Based on Maximally Contained Queries  $FQ_{MC}$**

22

## Maximally Contained Queries $FQ_{MC}$



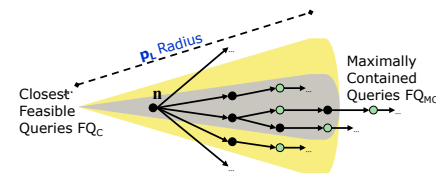
- Assuming fixed SELECT clause (projection list)
- Covered extensively in literature
  - MiniCon, Bucket, InverseRules Algorithms
- $FQ_{MC}$  is finite

23

## Closest Feasible Queries $FQ_C$ Algorithm

**Challenge: How far can the Closest Feasible Queries  $FQ_C$  be?**

**Solution: Maximally Contained Queries  $FQ_{MC}$**

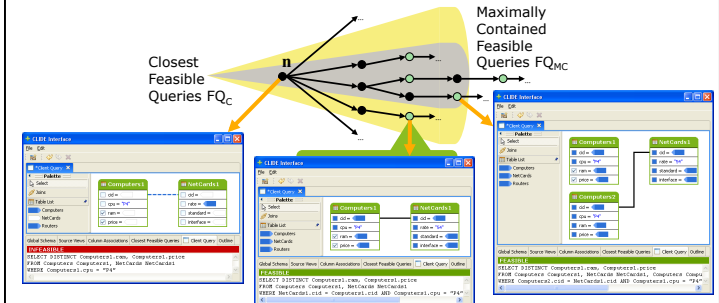


- Compute maximally contained queries  $FQ_{MC}$
- **Theorem: All  $FQ_C$  queries are reachable via a path of length  $p \leq p_L$**
- The radius  $p_L$  is the longest path to a maximally contained query

24

# Closest Feasible Queries $FQ_C$ Algorithm

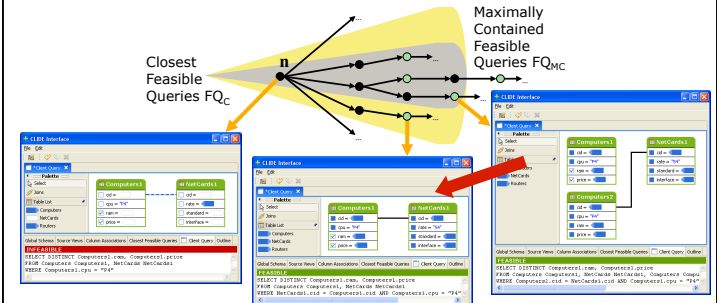
Challenge: Find the Closest Feasible Queries



- **Theorem:** All queries in  $FQ_{MC}$  are in  $FQ_C$
- But not all queries in  $FQ_C$  are in  $FQ_{MC}$

# Closest Feasible Queries $FQ_C$ Algorithm

Solution: Collapse Aliases



- Collapse Aliases to compute  $FQ_C \setminus FQ_{MC}$
- Check satisfiability

# Color Algorithm

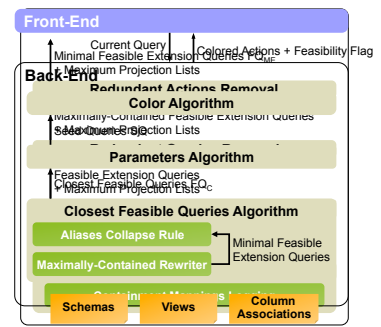
## Yellow and Blue

- An action  $\alpha$  is colored based on which closest feasible queries it appear in
- Yellow, if  $\alpha$  appears in all queries in  $FQ_C$
- Blue, if  $\alpha$  appears in at least one (but not all) query in  $FQ_C$

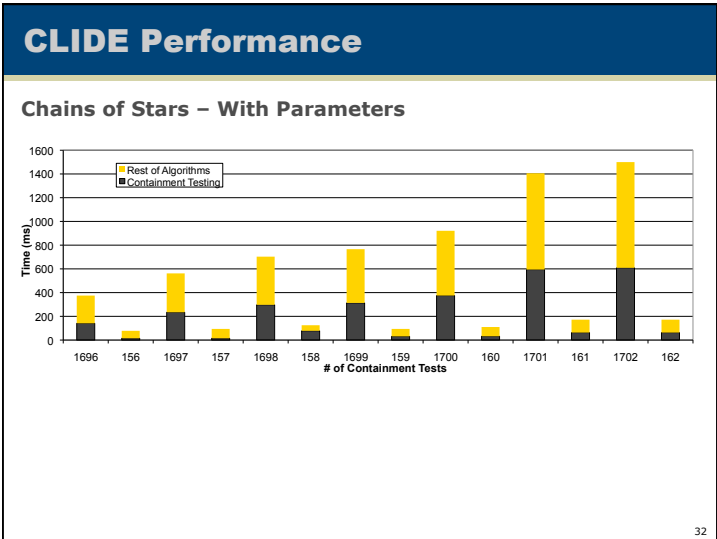
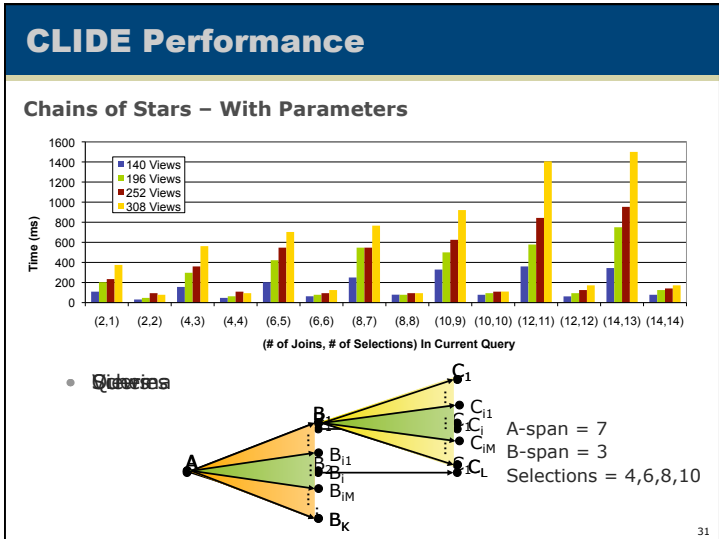
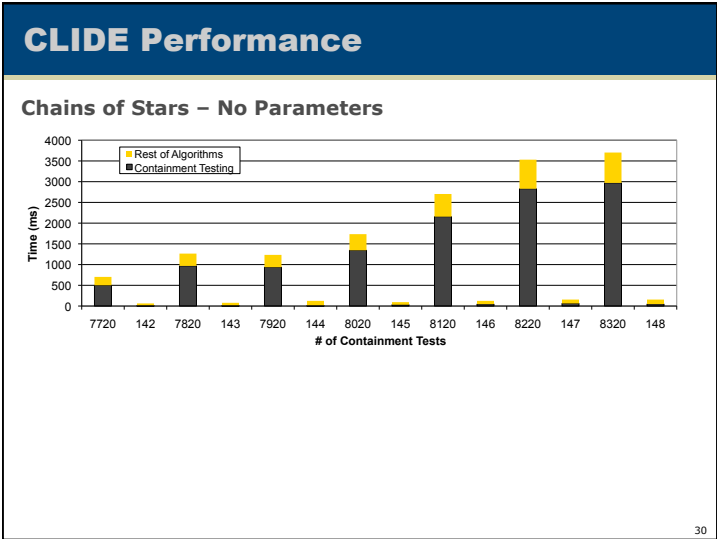
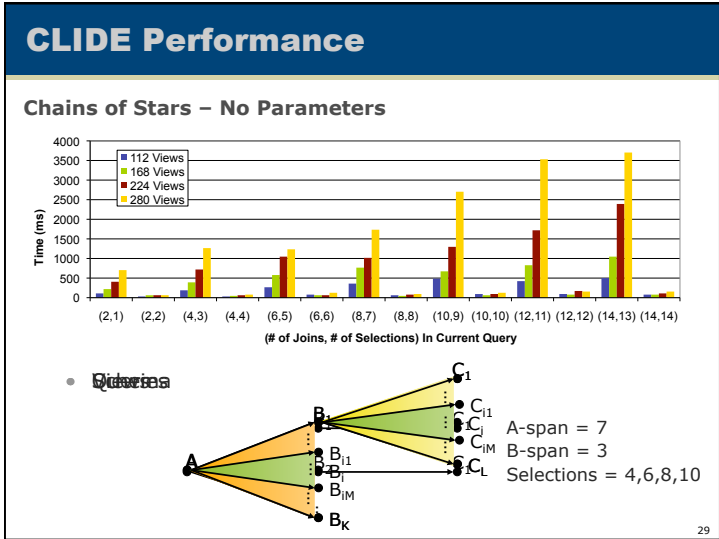
## White and Red

- Attach *Maximum Projection Lists* to Closest Feasible Queries
  - Projections that can be added to a feasible query, without compromising feasibility
- Projection  $\alpha$  is white if in the maximum projection list
- Color selections based on projections

# CLIDE Implementation & Optimizations



- Views expansion introduce redundancy
  - Affects CLIDE's rapid convergence and summarization
- Efficient containment test crucial to redundancy removal





## CLIDE Summary

First interactive query formulation interface based on source and mediator capabilities

### **Applicability**

- Service-Oriented Architectures
- Privacy-Preserving Services

### **Contributions**

- Interaction Guarantees: Rapid Convergence, Completeness, Summarization of Suggestions
- Interaction Graph
- Back-End Algorithms
  - Closest Feasible Queries, Colors, Parameters
- Modular, Customizable Architecture

**<http://www.clide.info>**