


Querying Heterogeneous Information Sources Using Source Descriptions

VLDB 1996


Alon Y. Levy – AT&T Laboratories
Anand Rajaraman – Stanford University
Joann J. Ordille – Bell Labs

Presentation By: **Mirza Beg**



Outline

- o Problem Description
- o Proposed System
- o System Architecture
- o Description of System Modules
- o Algorithms
- o Experiments & Results
- o Discussion



Problem Statement

- o Increasing number of structured data sources
- o Interrelated data
- o The user interacts with each information source separately and combine data !

Alternatively :

- o How do we extract the relevant data for a given query ?



Solution

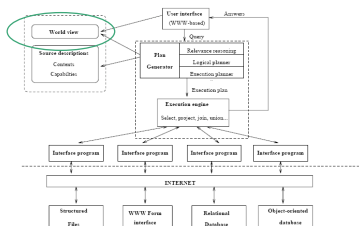
A System that:

- Provides a uniform query interface to distributed structured sources
- Uses source descriptions to describe data sources
- Generates executable query plans
- Returns the merged result set to the user

INFORMATION MANIFOLD



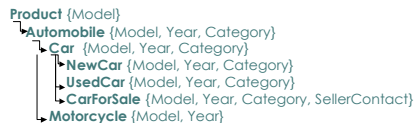
Information Manifold Architecture



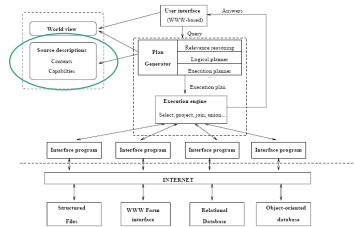


Information Manifold World View

- A virtual global schema on which the user can pose queries



Information Manifold Source Descriptions



Source Descriptions for Auto Sources

<p>Source 1: Used cars for sale. Accepts as input a category or model of car, and optionally a price range and a year range. For each car that satisfies the conditions, gives model, year, price, and seller contact information.</p>
<p>Source 2: Luxury cars for sale. All cars in this database are priced above \$20,000. Accepts as input a category of car and an optional price range. For each car that satisfies the conditions, gives model, year, price, and seller contact information.</p>
<p>Source 3: Vintage cars for sale (cars manufactured before 1950). Accepts as input a model and an optional year range. Gives model, year, price, and seller contact information for qualifying cars.</p>
<p>Source 4: Motorcycles for sale. Accepts as input a model and an optional price range. Gives model, year, price, and seller contact information.</p>
<p>Source 5: Car reviews database. Contains reviews for cars manufactured after 1990. Accepts as input a model and a year. Output is a car review for that model and year.</p>

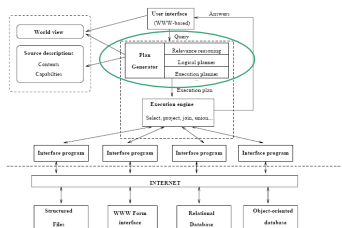
Content Records of Auto Sources

<p>Source 1: Used cars for sale. Contentsets: $V_1(c) \subseteq \text{CarForSale}(c)$, $\text{UsedCar}(c)$</p>
<p>Source 2: Luxury cars for sale. All cars in this database are priced above \$20,000. Contentsets: $V_2(c) \subseteq \text{CarForSale}(c)$, $\text{Price}(c, p)$, $p \geq 20000$</p>
<p>Source 3: Vintage cars for sale (cars manufactured before 1950). Contentsets: $V_3(c) \subseteq \text{CarForSale}(c)$, $\text{Year}(c, y)$, $y \leq 1950$</p>
<p>Source 4: Motorcycles for sale. Contentsets: $V_4(c) \subseteq \text{Motorcycle}(c)$</p>
<p>Source 5: Car reviews database. Contains reviews for cars manufactured after 1990. Contentsets: $V_5(m, p, r) \subseteq \text{Car}(c)$, $\text{Model}(c, m)$, $\text{Year}(c, y)$, $\text{ProductReview}(m, p, r)$</p>

Capability Records of Auto Sources

	Desired Inputs	Possible Outputs
Source 1: Used cars for sale.	$\{Model(c), Category(c)\}$	$\{Model(c), Category(c), Year(c), Price(c), SellerContact(c)\}$
Source 2: Luxury cars for sale. All cars in this database are priced above \$20,000.	$\{Year(c), Price(c)\}$	$\{Year(c), Price(c)\}$
Source 3: Vintage cars for sale (cars manufactured before 1950).	$\{Model(c), Category(c), Year(c), Price(c), SellerContact(c)\}$	$\{Price(c)\}$
Source 4: Motorcycles for sale.	$\{Model(c), Category(c), Year(c), Price(c), SellerContact(c)\}$	$\{Price(c)\}$
Source 5: Car reviews database. Contains reviews for cars manufactured after 1990.	$\{Model(c), Year(c), Price(c), SellerContact(c)\}$	$\{Price(c)\}$
	$\{(m, y), (m, y, r), \{1, 2, 2\}\}$	

Information Manifold Plan Generator



Query Reformulation Steps

$$Q(\bar{X}) \leftarrow R_1(\bar{Z}_1), \dots, R_n(\bar{Z}_n), C_Q$$

- Prune irrelevant sources
- Split query into sub goals
- Generate conjunctive query plans
- Find an executable ordering of sub goals

Step 1. Bucket Algorithm

Algorithm CreateBuckets(V, Q)
Inputs: V is a set of content descriptions, and Q is a conjunctive query of the form
 $Q: Q(X) \leftarrow R_1(X_1), \dots, R_n(X_n), C_Q$.

Set $Bucket_i$ to \emptyset for $1 \leq i \leq n$.
 For $i = 1, \dots, n$ do:
 For each $V \in \mathcal{V}$
 Let V be of the form:
 $V(Y) \subseteq S_1(Y_1), \dots, S_n(Y_n), C_V$
 For $j = 1, \dots, n$ do
 If $R_i = S_j$ or R_i and S_j are nondisjoint classes
 Let ψ be the mapping defined on the variables of V as follows:
 If y is the j 'th variable in Y and $y \in Y$
 then $\psi(y) = x_j$, where x_j is the j 'th variable in X_i .
 else $\psi(y)$ is a new variable that does not appear in Q or V .
 Let Q' be the 0-ary query:
 $Q' \leftarrow R_1(X_1), \dots, R_n(X_n), C_Q, S_1(\psi(Y_1)), \dots, S_n(\psi(Y_n)), \psi(C_V)$
 If Satisfiable(Q') then add $\psi(V)$ to $Bucket_i$.

End.

Step 1. Bucket Algorithm

Given a query Q :

- o Find a relevant source
- o Create a bucket for this sub-goal
- o Check source for Satisfiability
- o Add information source to bucket for this sub-goal

Example: Contents and Capabilities

<p>Source 1: Used cars for sale. Contents: $V_1(e) \subseteq CarForSale(e), IsUsedCar(e)$ Capabilities: $\{(Model(e), Category(e)), \{Model(e), Category(e), Year(e), Price(e), SellerContact(e)\}, \{Year(e), Price(e)\}, 1, 4)\}$</p>
<p>Source 2: Luxury cars for sale. All cars in this database are priced above \$30,000 Contents: $V_2(e) \subseteq CarForSale(e), Price(e, p), p \geq 30000$ Capabilities: $\{(Model(e), Category(e)), \{Model(e), Category(e), Year(e), Price(e), SellerContact(e)\}, \{Price(e)\}, 1, 3)\}$</p>
<p>Source 3: Vintage cars for sale (cars manufactured before 1950) Contents: $V_3(e) \subseteq CarForSale(e), Year(e, y), y \leq 1950$ Capabilities: $\{(Model(e)), \{Model(e), Category(e), Year(e), Price(e), SellerContact(e)\}, \{Price(e)\}, 1, 2)\}$</p>
<p>Source 4: Motorcycles for sale. Contents: $V_4(e) \subseteq Motorcycle(e)$ Capabilities: $\{(Model(e)), \{Model(e), Year(e), Price(e), SellerContact(e)\}, \{Price(e)\}, 1, 2)\}$</p>
<p>Source 5: Car reviews database. Contains reviews for cars manufactured after 1990. Contents: $V_5(m, y, r) \subseteq Car(e), Modl(e, m), Year(e, y), ProductReview(m, y, r)$ Capabilities: $\{(m, y), \{m, y, r\}\}, 2, 2)\}$</p>

● ● ● | **Bucket Algorithm:
Example**

$q(m_1, p_1, r_1) \leftarrow \text{CarForSale}(c_1), \text{Category}(c_1, \text{sportscar}), \text{Year}(c_1, y_1), y \geq 1992,$
 $\text{Price}(c_1, p_1), \text{Model}(c_1, m_1), \text{ProductReview}(m_1, y_1, r_1)$

$V_1^i(m, t, y, p, s) \subseteq \text{CarForSale}(c), \text{UsedCar}(c), \text{Model}(c, m), \text{Category}(c, t), \text{Year}(c, y),$
 $\text{Price}(c, p), \text{SellerContact}(c, s)$

find the mapping $e \rightarrow c_1$

$\text{CarForSale}(c_1), \text{Category}(c_1, \text{sportscar}), \text{Year}(c_1, y_1), y_1 \geq 1992,$
 $\text{Price}(c_1, p_1), \text{Model}(c_1, m_1), \text{ProductReview}(m_1, y_1, r_1), \text{UsedCar}(c_1), \text{SellerContact}(c_1, s)$

Source 1 is added to *bucket*₁
 Source 2 is added to *bucket*₂
 Source 3 does not get added because $(y \leq 1950, y \geq 1992)$ is not satisfiable

● ● ● | **Step 2. Finding an
Executable Ordering**

- Considering all possible combinations of information sources, enumerate semantically correct plans

● ● ● | **Step 2. Algorithm for finding
an Executable Ordering**

- Maintain a list of available parameters
- At every point add to the ordering any sub-goal whose input requirements are satisfied
- Push as many selections as possible to the sources

Step 3. Checking Containment

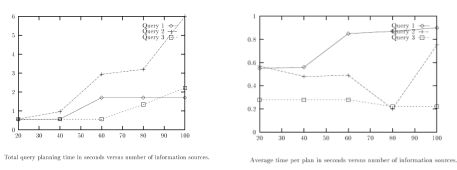
- Minimize each plan by removing redundant sub-goals

Experimental Results

Query	Number of sources	Max. bucket size	Plans enumerated	Plans generated	Time per plan (sec.)	Total time (sec.)
1	20	1	7	1	0.55	0.55
	40	1	7	1	0.56	0.56
	60	2	26	2	0.85	1.70
	80	2	26	2	0.85	1.70
	100	2	26	2	0.85	1.70
2	20	2	7	1	0.57	0.57
	40	3	11	2	0.48	0.96
	60	5	35	6	0.49	2.95
	80	6	44	8	0.40	3.20
	100	7	72	8	0.75	6.00
3	20	2	8	2	0.28	0.56
	40	2	8	2	0.28	0.56
	60	2	8	2	0.28	0.56
	80	6	49	6	0.22	1.32
	100	10	120	10	0.22	2.20

Query 1: Find titles and years of movies featuring Tom Hanks
 Query 2: Find titles and reviews of movies featuring Tom Hanks
 Query 3: Find telephone number(s) for Alaska Airlines

Experimental Results (cont.)





Conclusions

- A novel system that provides a DB-like query interface to distributed structured information sources
- Frees the user from interacting with each information source individually
- Integrates data from multiple sources and filters information
- Information Manifold applicable to WWW and company-wide d-DB's



Open Questions

- How to automatically extract contents and capabilities from sources ?
- Are there better algorithms to determine the relevant sources ?
- Scalability ?
- Overall Performance issues ?



Discussion Points

- A foundational paper in web-data mining.
- Substantial impact on current integration systems.
- Contents & capabilities at the core of the system yet no proposed generation algorithm.
- Experiments carried out on a very small set of queries.

••• | Questions ?

?
