


Planning for Information Gathering

Craig Knoblock
University of Southern California

These slides are based in part on slides from José Luis Ambite and Rao Kambhampati, which are in turn based in part on slides from Alon Halevy.

Craig Knoblock University of Southern California 1



Planning on the Web

- Part I: Planning for Information Gathering
- Part II: Plan Execution for Information Gathering

Craig Knoblock University of Southern California 2



Outline

- Information Gathering
- Planning for Information Gathering
 - View Integration
 - Query Reformulation
 - Source Capabilities
- Optimizing Information Gathering Plans
 - Removing Redundant Sources
 - Optimizing Sources and Queries
- Interleaving Planning and Sensing
 - Sensing to Handle Incomplete Information
 - Sensing to Optimize Plans
- Contingent Planning for Information Gathering
- Planning to Compose Web Sources
- Discussion

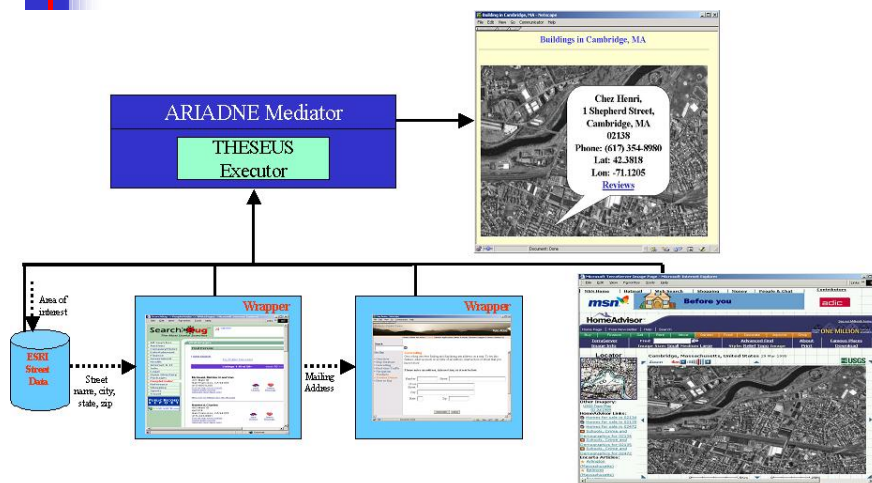
Craig Knoblock

University of Southern California

3



Information Gathering Example



Craig Knoblock

University of Southern California

4

Wrappers for Accessing Online Information Sources

- Wrappers provide uniform querying and data extraction

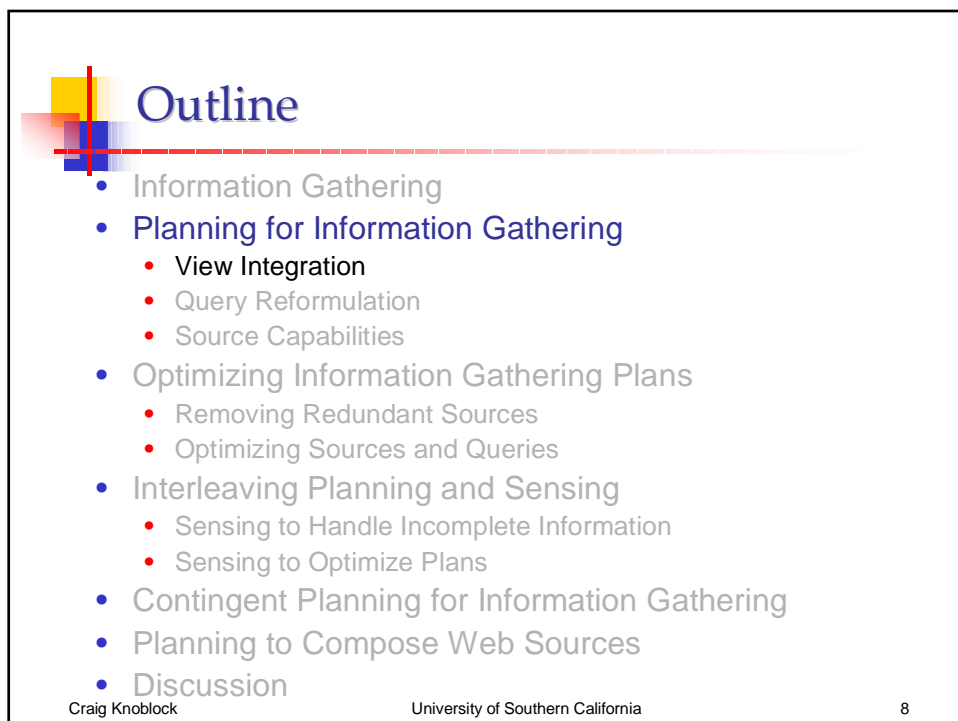
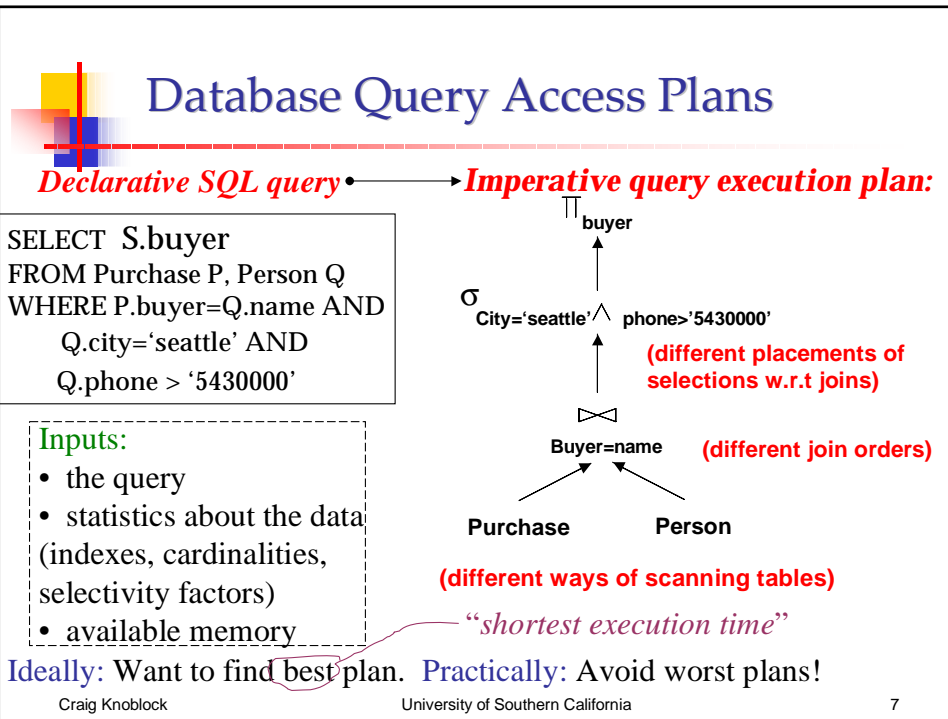


NAME Casablanca Restaurant
STREET 220 Lincoln Boulevard
CITY Venice
PHONE (310) 392-5751

- State of the art in wrapper induction
 - Data extraction is based on Web page layout (Muslea *et al.* 1999, Kushmerick *et al.* 1997)
 - User labels examples of data on pages
 - Induction algorithm learns extraction rules for data

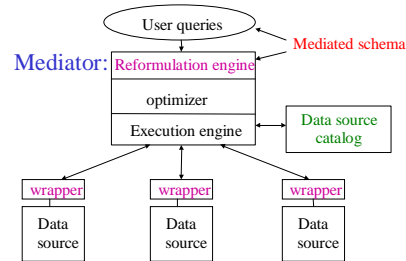
Planning for Information Gathering

- Database query access planning
 - Specialized planner optimized for task
 - Sources are fixed
 - Mappings predefined in global schema
 - Complete plan is generated and then executed
 - Assumes closed-world and complete information
- Distributed, heterogeneous environments:
 - Sources and mappings are not fixed
 - Sources are autonomous
 - Overlapping and redundant sources
 - Sources may be incomplete
 - Sources may be unavailable
 - Additional information may be required to access a source
 - Access to sources may be costly



Virtual Integration Architecture

- Leave the data in the sources
- When a query comes in:
 - Determine the relevant sources to the query
 - Break down the query into sub-queries for the sources
 - Get the answers from the sources, and combine them appropriately
- Data is fresh. Approach scalable
- Issues:
 - Relating Sources & Mediator
 - Reformulating the query
 - Efficient planning & execution



Garlic [IBM], Hermes[UMD];Tsimmis, InfoMaster[Stanford]; DISCO[INRIA]; Information Manifold [AT&T]; SIMS/Ariadne[USC];Emerac/Havasu[ASU]

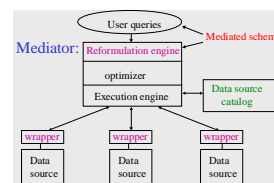
Craig Knoblock

University of Southern California

9

Desiderata for Relating Source-Mediator Schemas

- **Expressive power:** distinguish between sources with closely related data. Hence, be able to prune access to irrelevant sources.
- **Easy addition:** make it easy to add new data sources.
- **Reformulation:** be able to reformulate a user query into a query on the sources efficiently and effectively.
- **Nonlossy:** be able to handle all queries that can be answered by directly accessing the sources



Reformulation

- **Given:**
 - A query Q posed over the mediated schema
 - Descriptions of the data sources
- **Find:**
 - A query Q' over the data source relations, such that:
 - Q' provides only *correct answers* to Q, and
 - Q' provides *all possible answers* to Q given the sources.

Craig Knoblock

University of Southern California

10



Source Descriptions

Elements of source descriptions:

- Contents: source contains movies, directors, cast.
- Constraints: only movies produced after 1965.
- Completeness: contains *all* American movies.
- Capabilities:
 - Negative: source requires movie title or director as input
 - Positive: source can perform selections, joins, ...



Approaches to Specification of Source Descriptions

- Global-as-View (GAV):
Mediator relation defined as a view over source relations
Ex: TSIMMIS (Stanford), HERMES (Maryland)
- Local-as-View (LAV):
Source relation defined as view over mediator relations
Ex: Information Manifold (AT&T), Tukwila(UW), InfoMaster (Stanford), Ariadne (USC)

View ~ named query ~ logical formula



Views and Conjunctive Queries

```
CREATE VIEW Big-LA-buyers AS
SELECT buyer, seller, price
FROM Person, Purchase
WHERE Person.city = "Los Angeles" AND
      Person.name = Purchase.buyer AND
      Purchase.price > 10000
```

big-LA-buyers(Buyer, Seller, Price) :-
 person(Buyer, "Los Angeles"),
 purchase(Buyer, Seller, Product, Price),
 Price > 10000.

Datalog rule ~ view definition

Rule body ~ select-from-where construct of SQL



Outline

- Information Gathering
- Planning for Information Gathering
 - View Integration
 - Query Reformulation
 - Source Capabilities
- Optimizing Information Gathering Plans
 - Removing Redundant Sources
 - Optimizing Sources and Queries
- Interleaving Planning and Sensing
 - Sensing to Handle Incomplete Information
 - Sensing to Optimize Plans
- Contingent Planning for Information Gathering
- Planning to Compose Web Sources
- Discussion



Query Reformulation

Problem: rewrite the user query expressed in the mediated schema into a query expressed in the source schemas

Given a query Q in terms of the mediated-schema relations, and descriptions of the information sources,

Find a query Q' that uses only the source relations, such that

- $Q' \models Q$ (i.e., answers are correct; i.e., $Q' \subseteq Q$) and
- Q' provides all possible answers to Q given the sources



Global-as-View (GAV)

Each mediator relation is defined as a view over source relations.

$\text{MovieActor}(\text{title}, \text{actor}) \leftarrow$

$\text{DB1}(\text{id}, \text{title}, \text{actor}, \text{year})$

$\text{MovieActor}(\text{title}, \text{actor}) \leftarrow$

$\text{DB2}(\text{title}, \text{director}, \text{actor}, \text{year})$

$\text{MovieReview}(\text{title}, \text{review}) \leftarrow$

$\text{DB1}(\text{id}, \text{title}, \text{actor}, \text{year}) \wedge \text{DB3}(\text{id}, \text{review})$

Query Reformulation in GAV

Query reformulation = rule unfolding+simplication

Query: *Find reviews for 'DeNiro' movies*

$q(\text{title}, \text{review}) :- \text{MovieActor}(\text{title}, \text{'DeNiro'}),$
 $\text{MovieReview}(\text{title}, \text{review})$

1. $q'(\text{title}, \text{review}) :- \text{DB1}(\text{id}, \text{title}, \text{'DeNiro'}, \text{year}),$
 ~~$\text{DB1}(\text{id}, \text{title}, \text{actor}, \text{year'}), \text{DB3}(\text{id}, \text{review})$~~

Redundant

2. ~~$q'(\text{title}, \text{review}) :-$
 $\text{DB2}(\text{title}, \text{director}, \text{'DeNiro'}, \text{year}),$
 $\text{DB1}(\text{id}, \text{title}, \text{actor}, \text{year'}), \text{DB3}(\text{id}, \text{review})$~~

Redundant
wrt 1

Local-as-View (LAV)

- Each source relation is defined as a view over mediator relations

$V1(\text{title}, \text{year}, \text{director}) \xrightarrow{\subseteq} \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre})$
 $\wedge \text{American}(\text{director}) \wedge \text{year} \geq 1960 \wedge \text{genre} = \text{'Comedy'}$

$V2(\text{title}, \text{review}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge$
 $\text{year} \geq 1990 \wedge \text{MovieReview}(\text{title}, \text{review})$



Query Reformulation in LAV

Query: *Reviews for comedies produced after 1950*

$q(\text{title}, \text{review}) \text{ :- Movie}(\text{title}, \text{year}, \text{director}, \text{'Comedy'}), \text{ year} \geq 1950, \text{ MovieReview}(\text{title}, \text{review})$

Reformulated query:

$q'(\text{title}, \text{review}) \text{ :- V1}(\text{title}, \text{year}, \text{director}),$
 $\text{V2}(\text{title}, \text{review})$

$q' \subseteq q$

$\text{V1}(\text{title}, \text{year}, \text{director}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge$
 $\text{American}(\text{director}) \wedge \text{year} \geq 1960 \wedge \text{genre} = \text{'Comedy'}$

$\text{V2}(\text{title}, \text{review}) \rightarrow \text{Movie}(\text{title}, \text{year}, \text{director}, \text{genre}) \wedge \text{year} \geq 1990 \wedge$
 $\text{MovieReview}(\text{title}, \text{review})$



Inverse-Rules Algorithm

[Duschka+1997]

Idea: Construct an equivalent logic program which
evaluation yields the answer to the query

- The antecedent of the query and views is in term of mediator predicates
- Would like to have source predicates in antecedent so that program can be evaluated

\Rightarrow Invert the rules

(simply by using standard logical manipulations)

The Inverse-Rules Algorithm: Example

$V1(\text{dept}, \text{course}) \rightarrow \text{Enrolled}(\text{student}, \text{dept}) \wedge \text{Registered}(\text{student}, \text{course})$

$$a \rightarrow b \equiv \neg a \vee b$$


$$\begin{aligned} & \forall D, C [\neg V1(D, C) \rightarrow \exists S [e(S, D) \wedge r(S, C)]] \\ & \equiv \neg V1(D, C) \vee [e(f(D, C), D) \wedge r(f(D, C), C)] \\ & \equiv [\neg V1(D, C) \vee e(f(D, C), D)] \wedge [\neg V1(D, C) \vee r(f(D, C), C)] \\ & \equiv [V1(D, C) \rightarrow e(f(D, C), D)] \wedge [V1(D, C) \rightarrow r(f(D, C), C)] \\ & \equiv \\ & \quad e(f(D, C), D) \leftarrow V1(D, C) \\ & \quad r(f(D, C), C) \leftarrow V1(D, C) \end{aligned}$$

The Inverse-Rules Algorithm: Example

$q(D) \leftarrow \text{Enrolled}(S, D) \wedge \text{Registered}(S, \text{"DB"})$
 $V1(D, C) \rightarrow \text{Enrolled}(S, D) \wedge \text{Registered}(S, C)$

$q(D) \leftarrow \text{Enrolled}(S, D) \wedge \text{Registered}(S, \text{"DB"})$
 $\text{Enrolled}(f(D, C), D) \leftarrow V1(D, C)$
 $\text{Registered}(f(D, C), C) \leftarrow V1(D, C)$
 $q(D) \leftarrow V1(D, \text{"DB"})$


$\text{Ext}(V1) = \{(\text{"CS"}, \text{"DB"}), (\text{"EE"}, \text{"DB"}), (\text{"CS"}, \text{"AI"})\}$
 $\text{Ext}(q) = \{(\text{"CS"}), (\text{"EE"})\}$



GAV vs. LAV

<ul style="list-style-type: none"> • Not modular <ul style="list-style-type: none"> • Addition of new sources changes the mediated schema • Can be awkward to write mediated schema without loss of information • Query reformulation easy <ul style="list-style-type: none"> • reduces to view unfolding (polynomial) • Can build hierarchies of mediated schemas • Best when <ul style="list-style-type: none"> • Few, stable, data sources • well-known to the mediator (e.g. corporate integration) <ul style="list-style-type: none"> • Garlic, TSIMMIS, HERMES 	<ul style="list-style-type: none"> • Modular--adding new sources is easy • Very flexible--power of the entire query language available to describe sources • Reformulation is hard <ul style="list-style-type: none"> • Involves answering queries only using views (can be intractable) • Best when <ul style="list-style-type: none"> • Many, relatively unknown data sources • possibility of addition/deletion of sources <ul style="list-style-type: none"> • Information Manifold, InfoMaster, Emerac
--	--

Craig Knoblock
University of Southern California
23



Outline

- Information Gathering
- Planning for Information Gathering
 - View Integration
 - Query Reformulation
 - Source Capabilities
- Optimizing Information Gathering Plans
 - Removing Redundant Sources
 - Optimizing Sources and Queries
- Interleaving Planning and Sensing
 - Sensing to Handle Incomplete Information
 - Sensing to Optimize Plans
- Contingent Planning for Information Gathering
- Planning to Compose Web Sources
- Discussion

Craig Knoblock
University of Southern California
24



Modeling Source Capabilities

Negative capabilities:

- A web-site may require certain inputs (in an HTML form) to answer a query
- Need to consider only valid query execution plans

Positive capabilities:

- A source may be database (understands SQL)
- Need to decide the placement of operations according to capabilities

Problem: how to describe and exploit source capabilities



Negative Capabilities: Binding Patterns

Sources:

$\text{AAAIdb}^f(X) \rightarrow \text{AAAI Papers}(X)$

$\text{CitationDB}^{bf}(X,Y) \rightarrow \text{Cites}(X,Y)$

$\text{AwardDB}^b(X) \rightarrow \text{Award Paper}(X)$

Query: find all the award winning papers:

$q(X) \leftarrow \text{Award Paper}(X)$



Recursive Rewritings

$q(X) \leftarrow \text{AwardPaper}(X)$

- Problem: *Unbounded* union of conjunctive queries

$q_1(X) \leftarrow \text{AAAIdb}(X), \text{AwardDB}(X)$

$q_1(X) \leftarrow \text{AAAIdb}(X_1), \text{CitationDB}(X_1, X), \text{AwardDB}(X)$

...

$q_1(X) \leftarrow \text{AAAIdb}(X_1), \text{CitationDB}(X_1, X_2), \dots, \text{CitationDB}(X_n, X), \text{AwardDB}(X)$

- Solution: Recursive Rewriting

$\text{papers}(X) \leftarrow \text{AAAIdb}(X)$

$\text{papers}(X) \leftarrow \text{papers}(Y), \text{CitationDB}(Y, X)$

$q'(X) \leftarrow \text{papers}(X), \text{AwardDB}(X)$

$\text{AAAIdb}^f(X) \rightarrow \text{AAAI Papers}(X)$

$\text{CitationDB}^{bf}(X, Y) \rightarrow \text{Cites}(X, Y)$

$\text{AwardDB}^b(X) \rightarrow \text{AwardPaper}(X)$



Inverse-Rules Algorithm Binding Patterns

Sources:

$\text{AAAIdb}^f(X) \rightarrow \text{AAAI Papers}(X)$

$\text{CitationDB}^{bf}(X, Y) \rightarrow \text{Cites}(X, Y)$

$\text{AwardDB}^b(X) \rightarrow \text{AwardPaper}(X)$

Query: find all the award winning papers:

$q(X) \leftarrow \text{AwardPaper}(X)$



Inverse-Rules Algorithm

Inverse + Domain Rules (1)

Inverted Rules:

$\text{AAAPapers}(X) \leftarrow \text{AAAIdb}(X)$

$\text{Cites}(X, Y) \leftarrow \text{dom}(X) \wedge \text{CitationDB}(X, Y)$

$\text{AwardPaper}(X) \leftarrow \text{dom}(X) \wedge \text{AwardDB}(X)$

Domain Rules:

$\text{dom}(Y) \leftarrow \text{dom}(X) \wedge \text{CitationDB}(X, Y)$

$\text{dom}(X) \leftarrow \text{AAAIdb}(X)$

Query:

$q(X) \leftarrow \text{AwardPaper}(X)$



Inverse-Rules Algorithm


Inverse + Domain Rules (2)

Simplifying the program:

$q(X) \leftarrow \text{paper}(X) \wedge \text{AwardDB}(X)$

$\text{paper}(Y) \leftarrow \text{paper}(X) \wedge \text{CitationDB}(X, Y)$

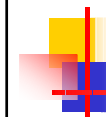
$\text{paper}(X) \leftarrow \text{AAAIdb}(X)$



Outline

- Information Gathering
- Planning for Information Gathering
 - View Integration
 - Query Reformulation
 - Source Capabilities
- Optimizing Information Gathering Plans
 - Removing Redundant Sources
 - Optimizing Sources and Queries
- Interleaving Planning and Sensing
 - Sensing to Handle Incomplete Information
 - Sensing to Optimize Plans
- Contingent Planning for Information Gathering
- Planning to Compose Web Sources
- Discussion

Craig Knoblock University of Southern California 31



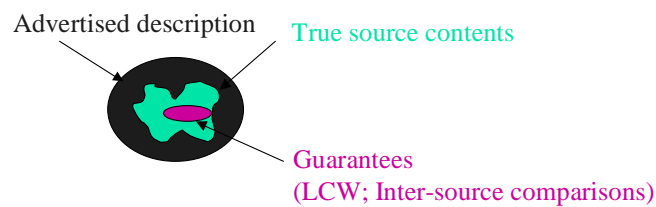
Managing Source Overlap

- Often, sources on the Internet have overlapping contents
 - The overlap is not centrally managed (unlike DDBMS—data replication etc.)
- Reasoning about overlap is important for plan optimality
 - We cannot possibly call all potentially relevant sources!
- Qns: How do we characterize and exploit source overlap?

Craig Knoblock University of Southern California 32

Local Completeness Information

- If sources are incomplete, we may need to look at all of them
- Often, sources are *locally complete*
- Movie(title, director, year) complete for years after 1960, or for American directors
- **Question:** given a set of local completeness statements, is a query Q' a complete answer to Q?



Craig Knoblock

University of Southern California

33

Using LCW rules to minimize plans

Basic Idea:


- If reformulation of Q leads to a union of conjunctive plans
 - $P_1 \vee P_2 \vee \dots P_k$
- Then, if P_1 is “complete” for Q (under the given LCW information), then we can minimize the reformulation by pruning $P_2 \dots P_k$
 - $[P_1 \wedge \text{LCW}]$ contains $P_1 \vee P_2 \vee \dots P_k$ [Duschka, AAAI-97]
- For Recursive Plans (obtained when the sources have access restrictions)
 - We are allowed to remove a rule r from a plan P , if the “complete” version of r is already contained in $P-r$

Emerac [Lambrecht & Kambhampati, 99]

Craig Knoblock

University of Southern California

34

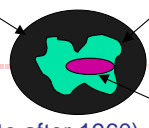


Example

Advertised description

True source contents

Guarantees




- **S1: Movie(title, director, year) (complete after 1960)**
 $S1(T,D,Y) \rightarrow M(T,D,Y)$
- **S2: Show(title, theater, city, hour)(complete for Seattle)**
 $S2(T,Th,C,H) \rightarrow Sh(T,Th,C,H)$
 LCW: $S2(T,Th,C,H) \leftarrow Sh(T,Th,C,H) \ \& \ C = \text{Seattle}$
- **S3: Show(title, theater, city, hour)**
 $S3(T,Th,C,H) \rightarrow Sh(T,Th,C,H)$
- **Query: Find movies and directors playing in Seattle**
 $Q(T,D) \leftarrow M(T,D,Y) \ \& \ Sh(T,Th,C,H) \ \& \ C = \text{"Seattle"}$
- **Plan: Combine S1 with S2 or S3**
 $Q(T,D) \leftarrow S1(T,D,Y) \ \& \ S2(T,Th,C,H) \ \& \ C = \text{"Seattle"}$
 $Q(T,D) \leftarrow S1(T,D,Y) \ \& \ S3(T,Th,C,H) \ \& \ C = \text{"Seattle"}$
- **Optimized Plan: Use LCW to prune S3**
 $Q(T,D) \leftarrow S1(T,D,Y) \ \& \ S2(T,Th,C,H) \ \& \ C = \text{"Seattle"}$

Craig Knoblock

University of Southern California

35



Outline

- Information Gathering
- Planning for Information Gathering
 - View Integration
 - Query Reformulation
 - Source Capabilities
- **Optimizing Information Gathering Plans**
 - Removing Redundant Sources
 - Optimizing Sources and Queries
- Interleaving Planning and Sensing
 - Sensing to Handle Incomplete Information
 - Sensing to Optimize Plans
- Contingent Planning for Information Gathering
- Planning to Compose Web Sources
- Discussion

Craig Knoblock

University of Southern California

36

Planning by Rewriting

[Ambite & Knoblock, 1998]

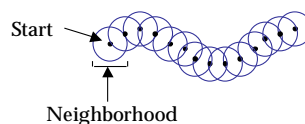
- Efficiently generate an initial solution plan (possibly of low quality)
- Iteratively rewrite the current plan
 - using a set of declarative plan rewriting rules
 - improving plan quality
 - until an acceptable solution or resource limit reached



Efficient High-Quality Planning

Planning by Rewriting as Local Search

- PbR: efficient high-quality planning using local search
- Main issues:
 - Selection of initial feasible point: Initial plan generation
 - Generation of a local neighborhood: Set of plans obtained from application of the plan rewriting rules
 - Cost function to minimize: Measure of plan quality
 - Selection of next point: Next plan to consider -- determines how the global space is explored



Planning by Rewriting for Query Planning in Mediators

- Initial plan generation: random parse of the query
- Plan rewriting rules: based on properties of:
 - relational algebra,
 - distributed environment,
 - integration axioms
- Plan quality: query execution time (size estimation)
- Search Strategies: gradient descent+restart, simulated annealing, variable-depth rewriting, ...

Craig Knoblock

University of Southern California

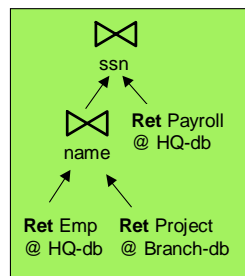
39

Query Planning in PbR

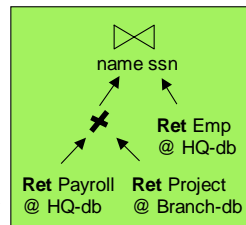
$a(\text{name sal proj}) :- \text{Emp}(\text{name ssn}) \wedge \text{Payroll}(\text{ssn sal}) \wedge \text{Projects}(\text{name proj})$

HQ-db
Emp(name ssn)
Payroll(ssn sal)

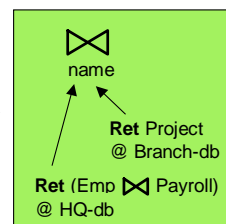
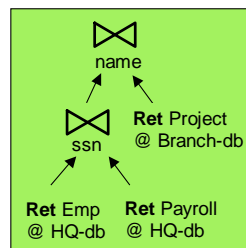
Branch-db
Project(name proj)



Join Swap



Remote Join Eval



Craig Knoblock

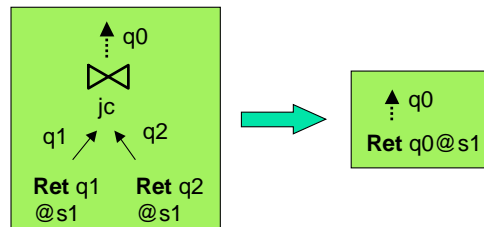
University of Southern California

40

Rewriting Rules: Distributed Environment remote-join-eval

(define-rule **remote-join-eval**

```
:if (:operators ((?n1 (retrieve ?source ?query1))
                  (?n2 (retrieve ?source ?query2)
                  (?n3 (join ?join-conds ?query0 ?query1 ?query2)))
:constraints (capability ?source join))
:replace (:operators (?n1 ?n2 ?n3))
:with (:operators ((?n4 (retrieve ?source ?query0))))))
```



Craig Knoblock

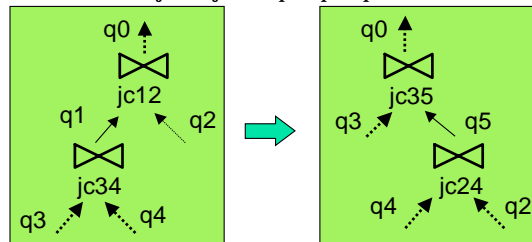
University of Southern California

41

Rewriting Rules: Relational Algebra join-associativity

(define-rule :name **join-associativity**

```
:if (:operators ((?n1 (join ?jc34 ?q1 ?q3 ?q4)
                  (?n2 (join ?jc12 ?q0 ?q1 ?q2)))
:constraints (join-swappable ?jc34 ?q1 ?q3 ?q4 ?jc12 ?q0 ?q2 ;;in
                  ?jc24 ?jc35 ?q5) ;; out
:replace (:operators (?n1 ?n2))
:with (:operators ((?n3 (join ?jc24 ?q5 ?q4 ?q2))
                  (?n4 (join ?jc35 ?q0 ?q3 ?q5))))
```



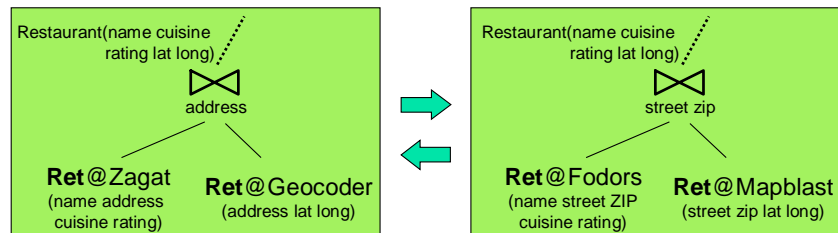
Craig Knoblock

University of Southern California

42


Rewriting Rules: Integration Axioms

- Rules computed from integration axioms relevant to query:
 $\text{Restaurant}(\text{name cuisine rating lat long}) =$
 - $\text{Zagat}(\text{name address cuisine rating}) \wedge \text{Geocoder}(\text{address lat long})$
 - $\text{Fodors}(\text{name street zip cuisine rating}) \wedge \text{Mapblast}(\text{street zip lat long})$



PbR in Query Planning: Summary

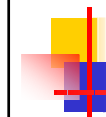
- Operators: output, retrieve, assign, select, join, union
- Plan rewriting rules:
 - Relational algebra: join-swap, selection-push-in, selection-push-out, assignment-push-in, assignment-push-out, selection-assignment-swap, push-join-thru-union, and push-union-thru-join.
 - Distributed environment: source-swap, remote-join-eval, remote-selection-eval, and remote-assignment-eval.
 - Integration axioms: computed automatically from the relevant integration axioms for classes in the query
- Search: gradient descent + random restart
 - first-improvement
 - steepest descent



Outline

- Information Gathering
- Planning for Information Gathering
 - View Integration
 - Query Reformulation
 - Source Capabilities
- Optimizing Information Gathering Plans
 - Removing Redundant Sources
 - Optimizing Sources and Queries
- Interleaving Planning and Sensing
 - Sensing to Handle Incomplete Information
 - Sensing to Optimize Plans
- Contingent Planning for Information Gathering
- Planning to Compose Web Sources
- Discussion

Craig Knoblock University of Southern California 45



Planning for the Internet Softbot

[Golden et al., 1996, Golden, 1998]

- XII and Puccini planners for the Internet Softbot
- Plans both gathering and manipulation actions
 - e.g., ls -a, chmod +r *
- Used to model Internet resources such as netfind
- Each resource modeled as an operator

Name: (netfind ?person)

Preconds:

```
(current.shell csh)
(isa netfind.server ?server)
(firstname ?person ?firstname)
(lastname ?person ?lastname)
(or
  (person.city ?person ?keyword)
  (person.institution ?person ?keyword))
```

Postconds:

```
(userid ?person !userid)
(person.machine ?person !machine)
```

Netfind Operator from XII

Craig Knoblock University of Southern California 46




Observational Effects and Knowledge Preconditions

- **Observational Effects**
 - Effect that changes the state of the world
`chmod +r foo.tex -- cause(readable(foo.tex))`
 - Effect that changes the agent's model of the world
`wc -- observe(word.count (file, !word))`
- **Knowledge Preconditions**
 - Information goal -- `find-out(length (paper.tex, l))`
 - Goals of achievement -- `satisfy(readable (f) False)`
- **Verification Links**
 - Alternative to knowledge preconditions
 - Assume secondary condition is true and then use an observational effect to determine whether it is true after execution



Sensing for Locally Complete Information


- **Reasons about incomplete information**
 - Uses LCW to reason about what it knows and what it doesn't know
 - e.g., `ls -a *` gives it locally complete information about the current directory
- **Interleaves sensing actions to gather LCW information**
 - LCW statements are a way of satisfying universally quantified goals
- **Provides fine-grained reasoning**
 - e.g., can request all recent techreports by X not already stored locally



Outline

- Information Gathering
- Planning for Information Gathering
 - View Integration
 - Query Reformulation
 - Source Capabilities
- Optimizing Information Gathering Plans
 - Removing Redundant Sources
 - Optimizing Sources and Queries
- Interleaving Planning and Sensing
 - Sensing to Handle Incomplete Information
 - Sensing to Optimize Plans
- Contingent Planning for Information Gathering
- Planning to Compose Web Sources
- Discussion

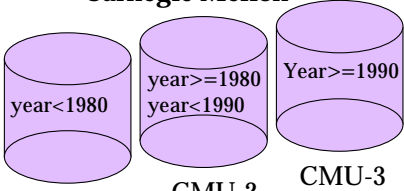
Craig Knoblock University of Southern California 49



Sensing to Determine Relevant Sources [Ashish, Knoblock, & Levy, 1997]

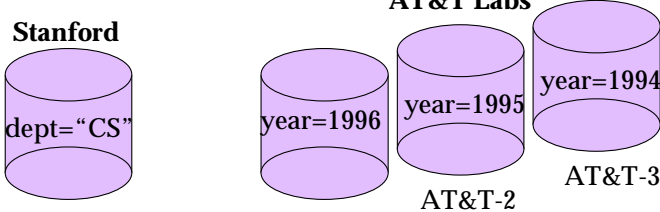
Technical Report Repositories

Carnegie Mellon



CMU-1 CMU-2 CMU-3

Stanford



AT&T-1 AT&T-2 AT&T-3

Craig Knoblock University of Southern California 50



Building a Discrimination Matrix

- Discrimination matrix specifies the relevant sources for each region of each attribute
- Approach:
 - Analyze source descriptions to build a discrimination matrix
 - Matrix partitions sources along some attribute
 - Discrimination matrix used to estimate the cost of querying with and without sensing
- Useful discriminations provided when:
 - Sources can be partitioned by some attribute
 - Exists another source that provides that attribute
- Example: Information about the year of a tech report reduces the relevant sources from 7 to 3



Discrimination Matrix

<u>Region</u>	<u>Relevant Sources</u>
< 1980	CMU-1, Stanford
[1980,1990)	CMU-2, Stanford
[1990,1994)	CMU-3, Stanford
[1994,1994]	CMU-3, Stanford, AT&T-1
[1995,1995]	CMU-3, Stanford, AT&T-2
[1996,1996]	CMU-3, Stanford, AT&T-3
> 1996	CMU-3, Stanford

Planning with Discriminating Queries

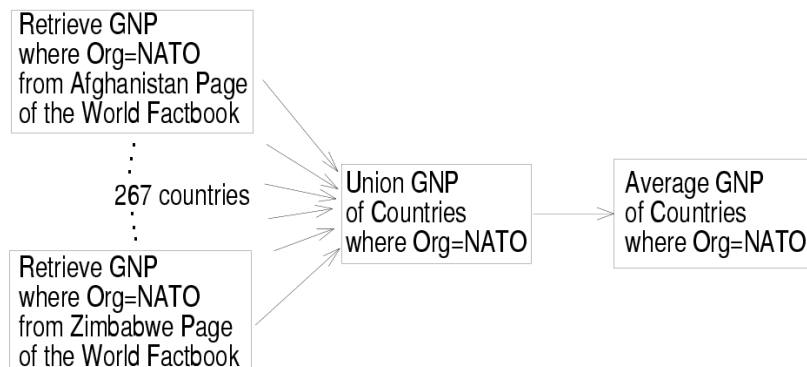
- Consider inserting a discriminating query for any subquery that:
 - Requires accessing multiple sources
 - There exists a discriminating attribute in the matrix
- Compare the cost of no discrimination to the combined cost of discriminating and querying
- Since we cannot know the results of the discrimination, use the average estimated cost
- Potentially relevant sources: $S = S_1, \dots, S_6$
- Discriminating queries: R_1, R_2
- Possible plans: $S, R_1 S', R_2 S'', R_1 R_2 S'''$
 - $R_1: \{\{S_1, S_2\}, \{S_3, S_4, S_5\}, \{S_6\}\}$
 - $R_2: \{\{S_1\}, \{S_2, S_3\}, \{S_4, S_5, S_6\}\}$
 - $R_1 R_2: \{\{S_1\}, \{S_2\}, \{S_3\}, \{S_4, S_5\}, \{S_6\}\}$

Craig Knoblock

University of Southern California

53

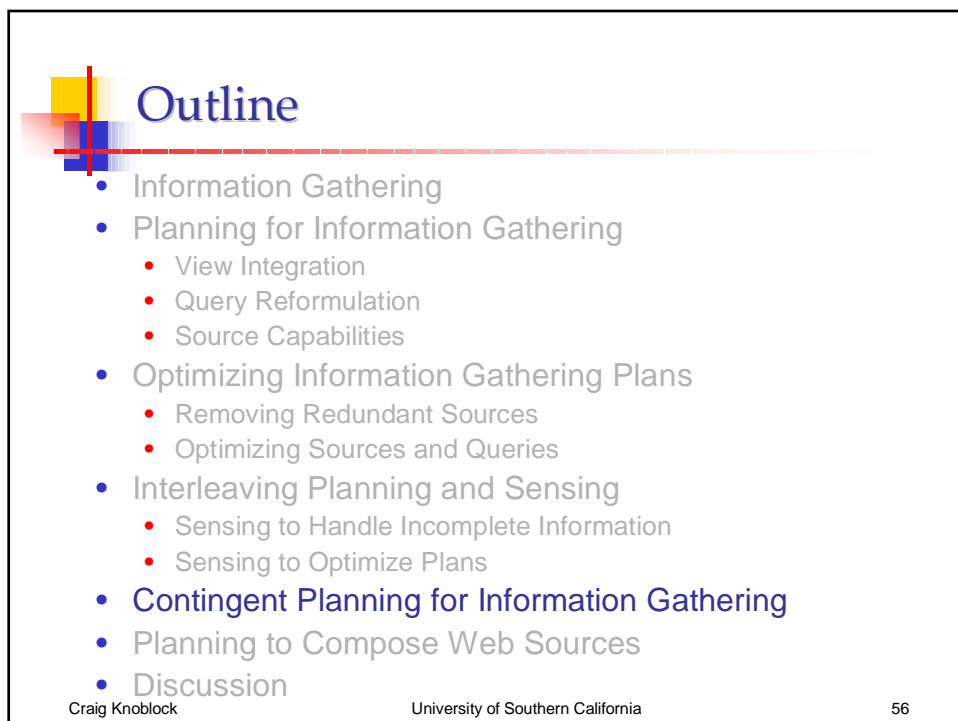
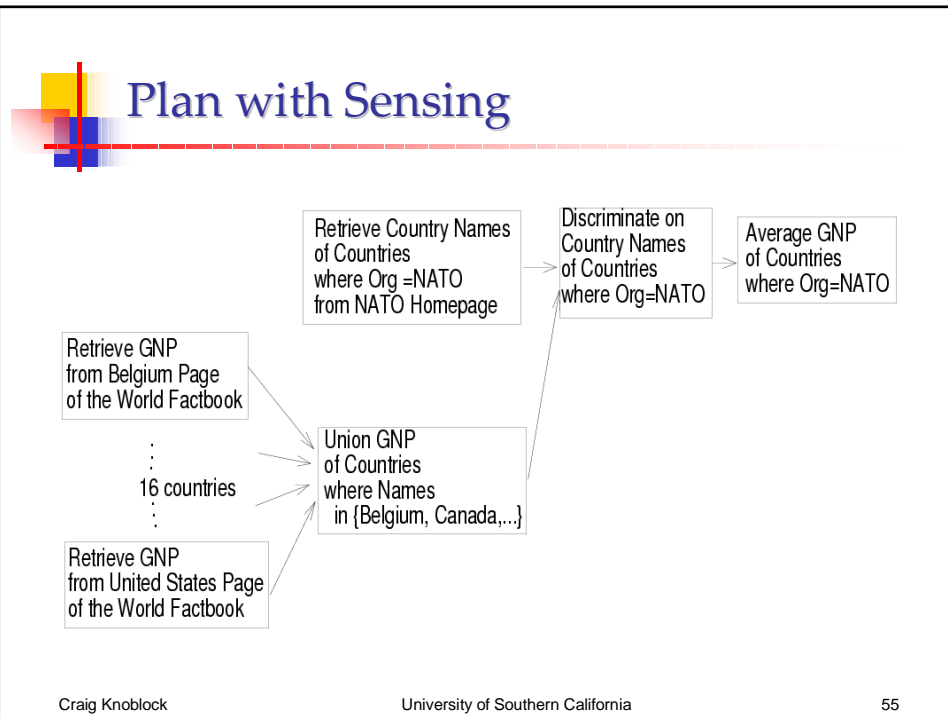
Plan without Sensing



Craig Knoblock

University of Southern California

54



Contingent Planning for Information Gathering [Friedman & Weld '97]

- Use subsumption relationships to make a plan more resource conscious
 - Determined based on LCW statements
- Execution policies:
 - Brute force – ignore subsumption and execute everything greedily
 - Aggressive – execute multiple alternatives and cancel others once a subsumed source is successful
 - Frugal – execute the most general source first and only execute others if it fails

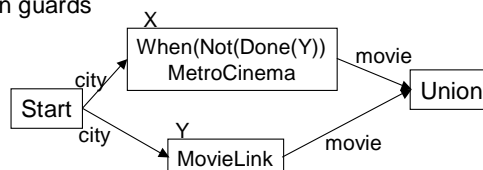
Craig Knoblock

University of Southern California

57

Augmenting the Plans


- Contingent plans
 - Operator can fire when its guard is true
 - Status variable for each operator
 - Sleeping, running, failed, and done
 - Approach:
 - Nodes initialized to running
 - Running nodes fired when input is available
 - Update status based on guards
 - Guards
 - Aggressive policy:
 - Frugal policy:
 - Failed(Y)



Craig Knoblock

University of Southern California

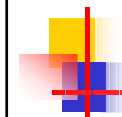
58



Outline

- Information Gathering
- Planning for Information Gathering
 - View Integration
 - Query Reformulation
 - Source Capabilities
- Optimizing Information Gathering Plans
 - Removing Redundant Sources
 - Optimizing Sources and Queries
- Interleaving Planning and Sensing
 - Sensing to Handle Incomplete Information
 - Sensing to Optimize Plans
- Contingent Planning for Information Gathering
- Planning to Compose Web Sources
- Discussion

Craig Knoblock University of Southern California 59



Composing Web Services

- Information sources only have inputs and outputs
 - Possibly with some additional constraints on those
- Services have:
 - Inputs and outputs
 - Preconditions and effects
- Could be cast as a traditional planning problem with preconditions and effects
- Example:
 - To purchase a book on Amazon has a precondition of having the money and has the effects of having the book and less money
- Services can be composed into compound services [McIlraith & Fadel, 2002]
- Stored and reused similar to Macrops [Eikes, 1972]

Craig Knoblock University of Southern California 60



Outline

- Information Gathering
- Planning for Information Gathering
 - View Integration
 - Query Reformulation
 - Source Capabilities
- Optimizing Information Gathering Plans
 - Removing Redundant Sources
 - Optimizing Sources and Queries
- Interleaving Planning and Sensing
 - Sensing to Handle Incomplete Information
 - Sensing to Optimize Plans
- Contingent Planning for Information Gathering
- Planning to Compose Web Sources
- Discussion

Craig Knoblock

University of Southern California

61



Discussion

- Is this planning?
 - Not in the sense of composing sequences of actions with interacting effects
 - Certainly in the broader sense of formulating a scheme or program for the accomplishment or attainment of some goal
- Good ideas can be shared across fields
 - Planning by rewriting based on traditional approaches to query planning
- Lots of interesting problems with real world applications
 - Optimizing the plans (e.g., interleaving sensing actions)
 - Interleaving source selection and plan optimization
 - Efficient execution of the plans (next class)

Craig Knoblock

University of Southern California

62



Bibliography

- Planning for Information Gathering
 - View Integration
 - Levy, Alon Y. (2000). Logic-based Techniques in Data Integration. *Logic Based Artificial Intelligence*, Edited by Jack Minker, Kluwer Publishers.
 - Halevy, Alon Y. (2001). Answering Queries Using Views: A Survey. *VLDB Journal*.
 - Duschka, Oliver M. (1997). Query Planning and Optimization in Information Integration. Ph.D. Thesis, Stanford University, Computer Science Technical Report STAN-CS-TR-97-1598.
 - Duschka, Oliver M. and Alon Y. Levy (1997). Recursive Plans for Information Gathering. *Proceedings of IJCAI-97*



Bibliography

- Planning for Information Gathering
 - Traditional Planning Approaches
 - Lambrecht, Eric and Subbarao Kambhampati (1997). Planning for Information Gathering: A Tutorial Survey. ASU CSE Technical Report 96-017.
 - Knoblock, Craig A. (1995). Planning, Executing, Sensing, and Replanning for Information Gathering. In *Proceedings of IJCAI-95*.
 - Knoblock, Craig A. (1996) Building a planner for information gathering: A report from the trenches. In *Proceedings of AIPS-96*.
 - Kwok, Chung T. and Daniel S. Weld (1996). Planning to Gather Information. In *Proceedings of AAAI-96*.



Bibliography

- Optimizing Information Gathering Plans
 - Removing Redundant Sources
 - Duschka, Oliver M. (1997). Query Optimization Using Local Completeness. In *Proceedings of AAAI-97*.
 - Lambrecht, Eric and Subbarao Kambhampati, and Senthil Gnanaprakasam. (1999) Optimizing Recursive Information Gathering Plans. In *Proceedings of IJCAI-99*.
 - Optimizing Sources and Queries
 - Ambite, Jose Luis and Craig A. Knoblock (2000). Flexible and Scalable Cost-based Query Planning in Mediators: A Transformational Approach. *Artificial Intelligence Journal* , 118(1-2):115—161.
 - Jose Luis Ambite and Craig A. Knoblock (1997). Planning by Rewriting: Efficient Generating High-Quality Plans. In *Proceedings of AAAI-1997*.



Bibliography

- Interleaving Planning and Sensing
 - Sensing to Handle Incomplete Information
 - Golden, Keith, Oren Etzioni, and Daniel S. Weld (1996). Planning with Execution and Incomplete Information. University of Washington, Department of Computer Science, Technical Report UW-CSE-96-01-09.
 - Golden, Keith (1998) Leap Before you Look: Information Gathering in the PUCCINI Planner. In *Proceedings of AIPS-98*.
 - Sensing to Optimize Plans
 - Ashish, Naveen, Craig A. Knoblock, and Alon Y. Levy (1997). Information Gathering Plans with Sensing Actions, *Recent Advances in AI Planning: 4th European Conference on Planning, ECP'97* . Springer-Verlag, New York, 1997.



Bibliography

- Contingent Planning for Information Gathering
 - Friedman, Marc and Daniel S. Weld. (1997). Efficiently Executing Information-Gathering Plans. In *Proceedings of IJCAI-97*, Nagoya, Japan, August 1997.
- Planning to Compose Web Sources
 - McIlraith, Shiela and Ronald Fadel (2002). Planning with Complex Actions. In *Proceedings of AIPS-2002 Workshop: Is There Life Beyond Operator Sequencing? - Exploring Real-World Planning*.