

A Query Language Based on the Ambient Logic

Luca Cardelli

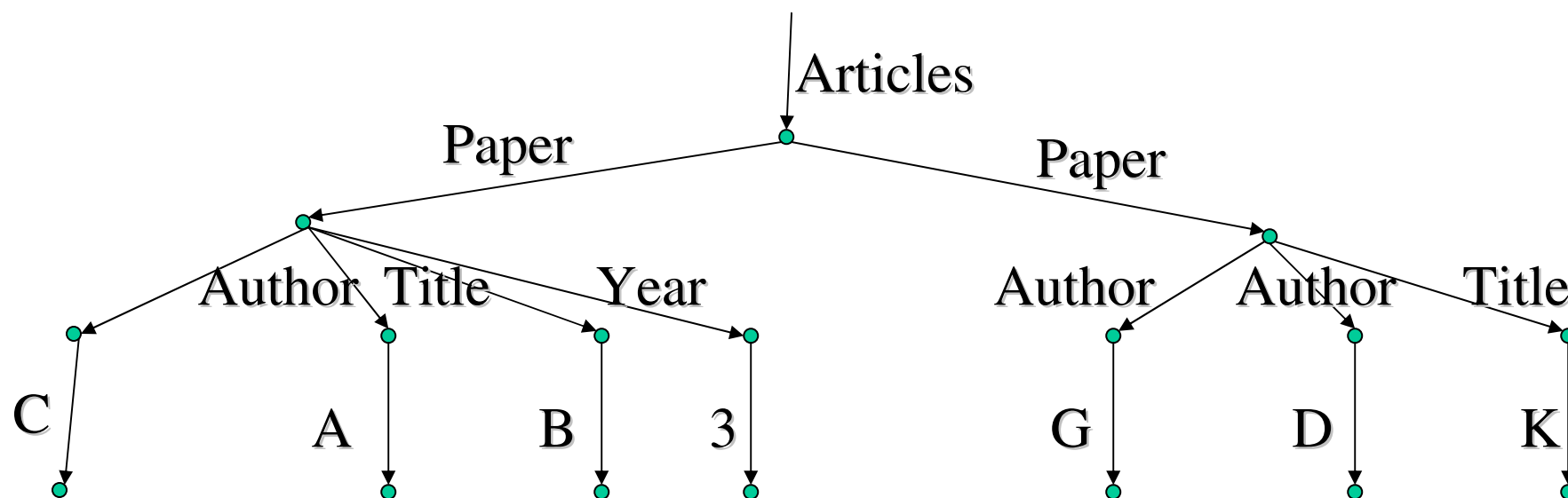
Microsoft Research

Giorgio Ghelli

University of Pisa

2001-07-06 Concoord - Lipari

Semistructured Data



- A tree (or graph), unordered (or ordered).
- Invented for “flexible” data representation (just like S-expressions...) for semi-irregular data like address books and bibliographies.
- Adopted by the DB community for the purpose of merging databases from different uncoordinated sources (without a common schema). Typically, web data that belongs to different institutions.

Unusual Data

- Not really arrays/lists:
 - Many children with the same label, instead of indexed children.
 - Mixture of repeated and non repeated labels under a node.
- Not really records:
 - Many children with the same label.
 - Missing/additional fields with no tagging information.
- Not really variants:
 - Labeled but untagged unions.
- Unusual data.
 - Yet, it aims to be the new universal standard for interoperability of programming languages, databases, e-commerce...

Unusual Data Manipulation

- New “flexible” type theories are required.
 - Based on the “effects” of processes over trees (Ambient Types).
 - Based on tree automata (Xduce).
- New processing languages required.
 - Xduce.
 - Various web scripting abominations.
- New query languages required. Various approaches.
 - From simple: Query existence of paths through the tree.
 - To fuzzy: Query whether a tree is kind of similar to another one.
 - To fancy: Query whether a tree is produced by a tree grammar.
 - To popular: “SQL for trees/graphs”.

Analogies

- An accidental(?) similarity between two areas:
- Semistructured Data is the way it is because:
 - “*You cannot rely on uniform structure*” of data.
Must abandon schemas based on records and tables.
 - Adopt “self-describing” data structures:
Use edge-labeled trees (or graphs).
- Mobile Computation is the way it is because:
 - “*You cannot rely on static structure*” of networks.
Must abandon type systems based on records and disjoint unions.
 - Adopt “self-describing” network structures:
Use edge-labeled trees (or graphs) of locations and agents.
- Both arose out of the Internet, because things there are just too dynamic for traditional notions of data and computation.

Relevance to Concurrency and Coordination

- Immediate implication: a new, uniform, model of data and computation on the Web, with opportunities for cross-fertilization:
 - Programming technology can be used to typecheck, navigate, and transform both dynamic network structures and the semistructured data they contain. Uniformly.
 - Database technology can be used to search through both dynamic network structures (“resource discovery”), and the semistructured data they contain. Uniformly.
- This convergence is still a dream, but it did motivate us to apply a particular technology developed for mobile computation to semistructured data:
 - Specification Logic → Query Logic

Concepts

- *Information trees* $I \in \mathcal{FT}$ (semistructured data)
- *Information terms* F (denoting information trees)
- *Formulas* \mathcal{A} (denoting sets of information trees)
- A semantics of terms $\llbracket F \rrbracket \in \mathcal{FT}$
- A semantics of formulas $\llbracket \mathcal{A} \rrbracket \subseteq \mathcal{FT}$
- A satisfaction (i.e. matching) relation $F \models \mathcal{A}$ (i.e. $\llbracket F \rrbracket \in \llbracket \mathcal{A} \rrbracket$)
- A *query language* Q (including *from* $F \models \mathcal{A}$ *select* Q ')
- A (naïve) query semantics $\llbracket Q \rrbracket \in \mathcal{FT}$
- A *table algebra* for matching evaluation (i.e. for $F \models \mathcal{A}$)
- A (refined) query semantics / query evaluation procedure for Q , based on the table algebra. Correct w.r.t. $\llbracket Q \rrbracket$.

Semantics: Information Trees

- Our semantic model for semistructured data: unordered edge-labeled finite-depth trees.
- Just to make it precise:
 - Λ is a countable collection of *labels*; m, n, \dots
 - \mathcal{FT} is the collection of *information trees* I :
 - The empty multiset, $\{\}^+$ is in \mathcal{FT} .
A root.
 - If m is in Λ and I is in \mathcal{FT} ,
then the singleton multiset $\{\langle m, I \rangle\}^+$ is in \mathcal{FT} .
An edge labeled m , leading to I .
 - \mathcal{FT} is closed under multiset union $\bigcup_{j \in J}^+ I_j$,
where J is a (possibly infinite) index set.
The root-merge of all the I_j .

Syntax: Information Terms

$F ::=$

$\mathbf{0}$ denoting the empty multiset

$m[F]$ denoting a singleton multiset

$F \mid F$ denoting binary multiset union

$\llbracket \mathbf{0} \rrbracket \triangleq \mathbf{0}$ where $\mathbf{0} \triangleq \{\}^+$

$\llbracket m[F] \rrbracket \triangleq m[\llbracket F \rrbracket]$ where $m[I] \triangleq \{\langle m, I \rangle\}^+$

$\llbracket F' \mid F'' \rrbracket \triangleq \llbracket F' \rrbracket \mid \llbracket F'' \rrbracket$ where $I' \mid I'' \triangleq I' \cup^+ I''$

- Often, $m[\mathbf{0}]$ is written $m[]$, or simply m .
- Define an equivalence \equiv such that $F \equiv F'$ iff $\llbracket F \rrbracket = \llbracket F' \rrbracket$.

Equivalence of Information Term

$$F \equiv F$$

$$F \equiv F' \Rightarrow F' \equiv F$$

$$F \equiv F', F' \equiv F'' \Rightarrow F \equiv F''$$

$$F \equiv F' \Rightarrow m[F] \equiv m[F']$$

$$F \equiv F' \Rightarrow F \mid F'' \equiv F' \mid F''$$

$$F \mid \mathbf{0} \equiv F$$

$$F \mid F' \equiv F' \mid F$$

$$(F \mid F') \mid F'' \equiv F \mid (F' \mid F'')$$

N.B.: $F \mid F' \neq F$ (these are multisets)

$$F \equiv F' \quad \text{iff} \quad \llbracket F \rrbracket = \llbracket F' \rrbracket$$

Example

```
Articles[
  Paper[
    Author[Cardelli] | Author[Gordon] |
    Title [Anywhere] |
    Year[2000] |
    Conf[POPL]
  ] |
  Paper[
    Author[Ghelli] |
    Title[Recursive] |
    Proceedings[VLDB] |
    Year[1998] |
    Editor[SV]
  ]
]
```

The Query Logic

$\mathcal{A}, \mathcal{B} \in \Phi ::=$	Formulas	(η is a name n or a variable x)
\mathbf{T}	true	
$\neg \mathcal{A}$	negation	
$\mathcal{A} \wedge \mathcal{B}$	conjunction	
$\exists x. \mathcal{A}$	existential quantification over label variables	
$\eta \sim \eta'$	label comparison (e.g. equality, regexp match..)	
$\mathbf{0}$	root	
$\eta[\mathcal{A}]$	edge	
$\mathcal{A} \mathcal{B}$	composition	
X	tree variable	
$\exists X. \mathcal{A}$	existential quantification over tree variables	
ξ	recursion variable	
$\mu \xi. \mathcal{A}$	recursive formula (least fixpoint)	
		ξ may occur only <i>positively</i> in \mathcal{A}

Satisfaction (first try)

- “Match a database F to a query \mathcal{A} , collect matches in ρ .”

$$F \models_{\rho} \mathbf{T}$$

$$F \models_{\rho} \neg \mathcal{A}$$

$$F \models_{\rho} \mathcal{A} \wedge \mathcal{B}$$

$$F \models_{\rho} \exists x. \mathcal{A}$$

$$F \models_{\rho} n \sim m$$

$$F \models_{\rho} \mathbf{0}$$

$$F \models_{\rho} n[\mathcal{A}]$$

$$F \models_{\rho} \mathcal{A} \mid \mathcal{B}$$

$$F \models_{\rho} X$$

$$F \models_{\rho} \exists X. \mathcal{A}$$

$$F \models_{\rho} \mu \xi. \mathcal{A}$$

$$\Leftrightarrow \neg F \models_{\rho} \mathcal{A}$$

$$\Leftrightarrow F \models_{\rho} \mathcal{A} \wedge F \models_{\rho} \mathcal{B}$$

$$\Leftrightarrow \exists n \in \Lambda. F \models_{\rho} \mathcal{A}\{x \leftarrow n\}$$

$$\Leftrightarrow n \sim m$$

$$\Leftrightarrow F \equiv \mathbf{0}$$

$$\Leftrightarrow \exists F'. F \equiv n[F'] \wedge F' \models_{\rho} \mathcal{A}$$

$$\Leftrightarrow \exists F', F''. F \equiv F' \mid F'' \wedge F' \models_{\rho} \mathcal{A} \wedge F'' \models_{\rho} \mathcal{B}$$

$$\Leftrightarrow F \equiv \rho(X)$$

$$\Leftrightarrow \exists F'. F \models_{\rho[X \mapsto F']} \mathcal{A} \quad (\mathcal{A}\{X \leftarrow F'\} \text{ bogus})$$

$$\Leftrightarrow F \models_{\rho} \mathcal{A}\{\xi \leftarrow \mu \xi. \mathcal{A}\}$$

“success”, but no further info

a match

- However, this is not a well-formed definition of $F \models_{\rho} \mathcal{A}$, because the μ case is circular. We need a proper semantics of fixpoints.

Example: Search

- Search:
 - “Find one of my articles (ignore non-articles); bind to X all info under the *article* label”:

$$S = \exists X. \text{article}[(\text{author}[\text{Cardelli}] \mid \mathbf{T}) \wedge X] \mid \mathbf{T}$$

- Can use recursive formulas to search deeper:

$$\mu \xi. S \vee \exists x. (x[\xi] \mid \mathbf{T})$$

- Not a query language yet.
 - It searches for one instance, not all instances.
 - Some *collecting* primitive must be added. This is going to be based on the logical notion of *satisfaction*.
- N.B.: $\exists X. \mathcal{A}$ is, logically, a bit strange:
 - X denotes a *singleton* set of trees (the match), not a general set of trees like a regular formula.

Example: Path Expressions

- Various kinds of path expressions can be defined:

$pq[\mathcal{A}]$	\triangleq	$p[q[\mathcal{A}]]$	path concatenation
$p^*[\mathcal{A}]$	\triangleq	$\mu\xi. \mathcal{A} \vee p[\xi]$	path iteration
$(p \vee q)[\mathcal{A}]$	\triangleq	$p[\mathcal{A}] \vee q[\mathcal{A}]$	path disjunction
$p(X)[\mathcal{A}]$	\triangleq	$p[X \wedge \mathcal{A}]$	binding the tree at the end
$.\eta[\mathcal{A}]$	\triangleq	$(\exists x. x \sim \eta \wedge x[\mathcal{A}]) \mid \mathbf{T}$	some η leads to \mathcal{A}
$.\neg\eta[\mathcal{A}]$	\triangleq	$(\exists x. \neg x \sim \eta \wedge x[\mathcal{A}]) \mid \mathbf{T}$	some non- η leads to \mathcal{A}
$!\eta[\mathcal{A}]$	\triangleq	$(\forall x. x \sim \eta \Rightarrow x[\Rightarrow \mathcal{A}]) \parallel \mathbf{F}$	all η lead to \mathcal{A}
etc.			

- “ X is an article that deals with *SSD*, or from each one can reach, through citations, an article that deals with *SSD*”:

$.article(X)(.cites.article)^*.keyword[SSD]$

Example: Schemas

- A logic is a “very rich type system”. Hence we can comfortably represent various kinds of schemas.
 - However, extensions (or unpleasant encodings) are required for ordered data: $\mathcal{A} \mid \mathcal{B}$ vs. $\mathcal{A}; \mathcal{B}$.
- Ex.: Xduce-like schemas:

$\mathbf{0}$	the empty tree
$\mathcal{A} \mid \mathcal{B}$	an \mathcal{A} next to a \mathcal{B}
$\mathcal{A} \vee \mathcal{B}$	either an \mathcal{A} or a \mathcal{B}
$n[\mathcal{A}]$	an edge n leading to an \mathcal{A}
\mathcal{A}^*	$\triangleq \mu \xi. \mathbf{0} \vee (\mathcal{A} \mid \xi)$ the merge of zero or more \mathcal{A} s
\mathcal{A}^+	$\triangleq \mathcal{A} \mid \mathcal{A}^*$ the merge of one or more \mathcal{A} s
$\mathcal{A}^?$	$\triangleq \mathbf{0} \vee \mathcal{A}$ zero or one \mathcal{A}

Semantics

$[\mathbf{T}]_{\rho,\delta}$	$\triangleq \mathcal{F}$	$F \models_{\rho,\delta} \mathcal{A} \triangleq [F] \in [[\mathcal{A}]_{\rho,\delta}]$
$[\neg \mathcal{A}]_{\rho,\delta}$	$\triangleq \mathcal{F} \setminus [[\mathcal{A}]_{\rho,\delta}]$	
$[\mathcal{A} \wedge \mathcal{B}]_{\rho,\delta}$	$\triangleq [[\mathcal{A}]_{\rho,\delta} \cap [[\mathcal{B}]_{\rho,\delta}]$	
$[\exists x. \mathcal{A}]_{\rho,\delta}$	$\triangleq \bigcup_{n \in \Lambda} [[\mathcal{A}]_{\rho[x \mapsto n], \delta}]$	
$[\eta \sim \eta']_{\rho,\delta}$	$\triangleq \text{if } \rho(\eta) \sim \rho(\eta') \text{ then } \mathcal{F} \text{ else } \{ \}$	
$[\mathbf{0}]_{\rho,\delta}$	$\triangleq \{ \mathbf{0} \}$	
$[\eta[\mathcal{A}]]_{\rho,\delta}$	$\triangleq \{ \rho(\eta)[I] \mid I \in [[\mathcal{A}]_{\rho,\delta}] \}$	
$[\mathcal{A} \mid \mathcal{B}]_{\rho,\delta}$	$\triangleq \{ I \mid I' \mid I \in [[\mathcal{A}]_{\rho,\delta}] \wedge I' \in [[\mathcal{B}]_{\rho,\delta}] \}$	
$[X]_{\rho,\delta}$	$\triangleq \{ \rho(X) \}$	
$[\exists X. \mathcal{A}]_{\rho,\delta}$	$\triangleq \bigcup_{I \in \mathcal{F}} [[\mathcal{A}]_{\rho[X \mapsto I], \delta}]$	
$[\xi]_{\rho,\delta}$	$\triangleq \delta(\xi)$	
$[\mu \xi. \mathcal{A}]_{\rho,\delta}$	$\triangleq \bigcap \{ S \subseteq \mathcal{F} \mid S \supseteq [[\mathcal{A}]_{\rho, \delta[\xi \mapsto S]}] \}$	

$[[\] \in \Phi \times ((Nvar \cup \Lambda \rightarrow \Lambda) \cup (Tvar \rightarrow \mathcal{F})) \times (Rvar \rightarrow \mathcal{P}(\mathcal{F})) \rightarrow \mathcal{P}(\mathcal{F})$

Derived Operators

\mathbf{F}	$\triangleq \neg \mathbf{T}$	falsity
$\mathcal{A} \vee \mathcal{B}$	$\triangleq \neg(\neg \mathcal{A} \wedge \neg \mathcal{B})$	disjunction
$\eta \neq \eta'$	$\triangleq \neg(\eta \sim \eta')$	label diversity
$\forall x. \mathcal{A}$	$\triangleq \neg \exists x. \neg \mathcal{A}$	universal name quantification
$\forall X. \mathcal{A}$	$\triangleq \neg \exists X. \neg \mathcal{A}$	universal tree quantification
$\nu \xi. \mathcal{A}$	$\triangleq \neg \mu \xi. \neg \mathcal{A} \{ \xi \leftarrow \neg \xi \}$	maximal fixpoint

$\mathbf{1}$	$\triangleq \neg \mathbf{0}$	non-empty
$\eta[\Rightarrow \mathcal{A}]$	$\triangleq \neg(\eta[\neg \mathcal{A}])$	edge implication

if edge is η , it leads to \mathcal{A}

$\mathcal{A} \parallel \mathcal{B}$	$\triangleq \neg(\neg \mathcal{A} \mid \neg \mathcal{B})$	decomposition
-------------------------------------	---	---------------

for every partition, either \mathcal{A} holds one part, or \mathcal{B} holds on the other
 ($\mathcal{A} \parallel \mathbf{F}$: all partitions satisfy \mathcal{A})

Dualization

- We can dualize all operators.

$$\neg \mathbf{0} \rightsquigarrow \mathbf{1}$$

$$\neg \mathbf{1} \rightsquigarrow \mathbf{0}$$

$$\neg \eta[A] \rightsquigarrow \eta[\Rightarrow \neg A]$$

$$\neg \eta[\Rightarrow A] \rightsquigarrow \eta[\neg A]$$

$$\neg(A \mid B) \rightsquigarrow \neg A \parallel \neg B$$

$$\neg(A \parallel B) \rightsquigarrow \neg A \mid \neg B$$

$$\neg \nu \xi. A \rightsquigarrow \mu \xi. \neg A \{ \xi \leftarrow \neg \xi \}$$

$$\neg \mu \xi. A \rightsquigarrow \nu \xi. \neg A \{ \xi \leftarrow \neg \xi \}$$

- Plus the usual DeMorgan laws.
- This gives us an implementation strategy for $\neg A$.
 - If we take the dual operators as primitive and implement them directly, we can then “push negation to the leaves”.

Logical Equations

- We can commute/distribute/simplify many operators.
 - This gives us opportunities for query optimization.

$\eta[A]$	$\Leftrightarrow \eta[\mathbf{T}] \wedge \eta[\Rightarrow A]$	$\eta[\Rightarrow A]$	$\Leftrightarrow \eta[\mathbf{T}] \Rightarrow \eta[A]$
$\eta[\mathbf{F}]$	$\Leftrightarrow \mathbf{F}$	$\eta[\Rightarrow \mathbf{T}]$	$\Leftrightarrow \mathbf{T}$
$\eta[A \wedge B]$	$\Leftrightarrow \eta[\mathbf{T}] \wedge \eta[B]$	$\eta[\Rightarrow A \vee B]$	$\Leftrightarrow \eta[\Rightarrow A] \vee \eta[\Rightarrow B]$
$\eta[A \vee B]$	$\Leftrightarrow \eta[A] \vee \eta[B]$	$\eta[\Rightarrow A \wedge B]$	$\Leftrightarrow \eta[\Rightarrow A] \wedge \eta[\Rightarrow B]$
$A \mathbf{F}$	$\Leftrightarrow \mathbf{F}$	$A \mathbf{T}$	$\Leftrightarrow \mathbf{T}$
$\mathbf{T} \mathbf{T}$	$\Leftrightarrow \mathbf{T}$	$\mathbf{F} \mathbf{F}$	$\Leftrightarrow \mathbf{F}$
$A \mathbf{0}$	$\Leftrightarrow A$	$A \mathbf{1}$	$\Leftrightarrow A$
$A B$	$\Leftrightarrow B A$	$A B$	$\Leftrightarrow B A$
$(A B) C$	$\Leftrightarrow A (B C)$	$(A B) C$	$\Leftrightarrow A (B C)$
$A (B \vee C)$	$\Leftrightarrow (A B) \vee (A C)$	$A (B \wedge C)$	$\Leftrightarrow (A B) \wedge (A C)$

- Dualization and other manipulations are based on logically valid equations; these have been studied extensively for the original Ambient Logic.

The Query Language

$Q ::=$

$from\ Q \models \mathcal{A}\ select\ Q'$

X

\emptyset

$\eta[Q]$

$Q \mid Q'$

$f(Q)$

Query

match and collect

matching variable

empty result

nesting of result

composition of results

tree functions (for extensibility)

- $from\ Q \models \mathcal{A}\ select\ Q'$

All the matches of Q with \mathcal{A} are computed, producing bindings for the x and X variables that are free in \mathcal{A} . The result expression Q' is evaluated for each (distinct!) such binding, and all the results are merged by \mid .

- N.B.: This general approach to building a query language Q for a logic \mathcal{A} , is fairly independent from the details of the logic.

Query Examples

- Joins

Merge info about persons from two db's:

```
from db1  $\vDash$  person[name[ $X^\lambda$ ] |  $Y^\lambda$ ] | T select  
from db2  $\vDash$  person[name[ $X$ ] |  $Z^\lambda$ ] | T select  
person[name[ $X$ ] |  $Y$  |  $Z$ ]
```

λ : binding occurrence

- Restructuring

Rearrange publications from by-article to by-year,
for each distinct year (i.e., for each distinct binding of X):

```
from db  $\vDash$  .article[.year[ $X^\lambda$ ]] select  
publications-by-year[  
  year[ $X$ ] |  
  from db  $\vDash$  .article[year[ $X$ ] |  $Z^\lambda$ ] select article[ $Z$ ]]
```

Z binds all fields except *year*; this is rather unusual in QL's

Query Examples

- Recursion

Find all email (or e-mail) addresses:

```
from db  $\models \mu\xi. .email[X^\lambda] \vee .e-mail[X^\lambda] \vee \exists x. .x[\xi]$   
select email[X]
```

- Unsafe queries

To be avoided by static or dynamic detection:

```
from db  $\models (male[X^{\lambda?}] \vee female[Y^{\lambda?}]) \mid \mathbf{T}$  select X? \mid Y?
```

```
from db  $\models \neg author[X^{\lambda?}]$  select X?
```

Reference Query Semantics

- Schemas (for rows and tables)
 - $\mathbf{V} = V_1 \dots V_n$ is a collection of distinct variables: either label variables x or tree variables X .
- Rows (or Valuations)
 - $\rho^{\mathbf{V}}$ is a row with schema \mathbf{V} : it maps each variable in \mathbf{V} to a value of the appropriate kind (label or tree).
 - A row can be seen as the result of a match, or as an environment in which to evaluate further matches.
- Tables (or Relations)
 - A set of rows with a common schema \mathbf{V} is called a table. That is, \mathbf{V} names the columns of the table.
- Query Semantics $\llbracket Q \rrbracket_{\rho^{\mathbf{V}}}$
 - Produces a tree: the result of the query Q , where $\rho^{\mathbf{V}}$ is used as an environment for the free variables of Q (included in \mathbf{V}).

Reference Query Semantics

$$\begin{aligned}
 \llbracket X \rrbracket_{\rho, \mathbf{v}} &\triangleq \rho^{\mathbf{v}}(X) \\
 \llbracket 0 \rrbracket_{\rho, \mathbf{v}} &\triangleq \mathbf{0} \\
 \llbracket \eta[Q] \rrbracket_{\rho, \mathbf{v}} &\triangleq \rho^{\mathbf{v}}(\eta) \llbracket [Q] \rrbracket_{\rho, \mathbf{v}} && (\rho^{\mathbf{v}}(n) \triangleq n) \\
 \llbracket Q \mid Q' \rrbracket_{\rho, \mathbf{v}} &\triangleq \llbracket [Q] \rrbracket_{\rho, \mathbf{v}} \mid \llbracket [Q'] \rrbracket_{\rho, \mathbf{v}} \\
 \llbracket f(Q) \rrbracket_{\rho, \mathbf{v}} &\triangleq f(\llbracket [Q] \rrbracket_{\rho, \mathbf{v}})
 \end{aligned}$$

$$\begin{aligned}
 \llbracket \text{from } Q \models \mathcal{A} \text{ select } Q' \rrbracket_{\rho, \mathbf{v}} &\triangleq \\
 &\bigcup^+ \text{for } \rho', \mathbf{v}' \supseteq \rho, \mathbf{v} \text{ s.t. } \mathbf{V}' = \mathbf{V} \cup FV(\mathcal{A}) \wedge \llbracket [Q] \rrbracket_{\rho, \mathbf{v}} \in \llbracket [\mathcal{A}] \rrbracket_{\rho', \mathbf{v}', \emptyset} \\
 &\text{of } \llbracket [Q'] \rrbracket_{\rho', \mathbf{v}'}
 \end{aligned}$$

– $\llbracket \text{from } Q \models \mathcal{A} \text{ select } Q' \rrbracket_{\rho, \mathbf{v}}$

Consider all the valuations ρ' that extend the current valuation ρ with the (still-free) variables of \mathcal{A} , and such that Q matches \mathcal{A} under the appropriate valuations.

For all such (distinct) valuations, compute Q' and put all the results in parallel.

Just a “reference” semantics

- The reference semantics is clean, but is hopelessly inefficient.
 - Literally, it requires computing beforehand a potentially infinite set of $\rho'v \supseteq \rho^v$, which are then filtered by *checking* that $\llbracket Q \rrbracket_{\rho} v \in \llbracket \mathcal{A} \rrbracket_{\rho, v, \emptyset}$. The $\rho'v$ that survive (hopefully a finite set) are used to build the result.
 - It should be much better to compute this (hopefully) finite set of useful $\rho'v$ on the fly, by *matching* Q to \mathcal{A} .
- This idea leads to a *table algebra*, used for building the relevant valuations while matching Q to \mathcal{A} , and to a refined evaluation procedure.

The Table Algebra

- A relational-style algebra for *relations over labels & trees*:

1^V		the largest table with schema V		
$R^V \cup^V R'^V$	\triangleq	$R^V \cup R'^V$	\subseteq	1^V
$Co^V(R^V)$	\triangleq	$1^V \setminus R^V$	\subseteq	1^V
$R^V \times^{V,V'} R'^V$	\triangleq	$\{\rho; \rho' \mid \rho \in R^V, \rho' \in R'^V\}$	\subseteq	$1^{V \cup V'}$ ($V \cap V' = \emptyset$)
$\Pi_{V'}^V R^V$	\triangleq	$\{\rho' \in 1^{V'} \mid \exists \rho \in R^V. \rho \supseteq \rho'\}$	\subseteq	$1^{V'}$ ($V' \subseteq V$)
$\sigma_{\eta \sim \eta'}^V R^V$	\triangleq	$\{\rho \in R^V \mid \rho^V(\eta) \sim \rho^V(\eta')\}$	\subseteq	1^V ($FV(\eta, \eta') \subseteq V$)

- Derived operators

$Ext_{V'}^V(R^V)$	\triangleq	$R^V \times^{V,V' \setminus V} 1^{V' \setminus V}$	\subseteq	$1^{V'}$ ($V \subseteq V'$)
$R^V \cap^V R'^V$	\triangleq	$Co^V(Co^V(R^V) \cup^V Co^V(R'^V))$	\subseteq	1^V
$R^V \bowtie^{V \cup V'} R'^V$	\triangleq	$Ext_{V \cup V'}^V(R^V) \cap^{V \cup V'} Ext_{V \cup V'}^{V'}(R'^V)$	\subseteq	$1^{V \cup V'}$
...				

Query Evaluation

- We define a refined query semantics:
 - A procedure $Q(Q)_{\rho v}$ that evaluate queries, uses a procedure $B(I, \mathcal{A})_{\rho v}$ to evaluate binders $Q \models \mathcal{A}$ in the *from-select* case.
 - $B(I, \mathcal{A})_{\rho v}$ produces a table, and is computed using the table algebra operators.
 - Natural join, \bowtie , takes the role of “unification” for match-variables in binders.
- Some cases (simplified for non-recursive formulas):

$$Q(\text{from } Q \models \mathcal{A} \text{ select } Q')_{\rho v} \triangleq \\ \text{let } I = Q(Q)_{\rho v} \text{ and } R^{FV(\mathcal{A}) \setminus V} = B(I, \mathcal{A})_{\rho v} \text{ in } \bigcup_{\rho' \in R} Q(Q')_{(\rho v; \rho')}$$

$$B(I, \neg \mathcal{A})_{\rho v} \triangleq C_{O^{FV(\mathcal{A}) \setminus V}}(B(I, \mathcal{A})_{\rho v})$$

$$B(I, \mathcal{A} \wedge \mathcal{B})_{\rho v} \triangleq B(I, \mathcal{A})_{\rho v} \bowtie^{FV(\mathcal{A}) \setminus V, FV(\mathcal{B}) \setminus V} B(I, \mathcal{B})_{\rho v}$$

$$B(I, \exists X. \mathcal{A})_{\rho v} \triangleq \Pi_{FV(\mathcal{A}) \setminus V \setminus \{X\}}^{FV(\mathcal{A}) \setminus V} B(I, \mathcal{A})_{\rho v}$$

Correctness Theorem

- I.e., instead of computing the set $\llbracket \mathcal{A} \rrbracket$, as in the reference semantics, and checking at the end whether I is in $\llbracket \mathcal{A} \rrbracket$, we process I and \mathcal{A} together in \mathbf{B} , step by step.
- The table algebra is still specified rather abstractly. Any particular implementation of the table algebra yields an implementation of \mathbf{B} , and hence of the query procedure Q .
- Correctness: the query evaluation procedure conforms to the reference semantics, so it is correct:

$$\forall Q. \forall V \supseteq FV(Q). \forall \rho^V. \quad Q(Q)_{\rho^V} = \llbracket Q \rrbracket_{\rho^V}$$

By a simple induction, with a non-trivial lemma in the recursive formulas case.

What's left out

- Effective implementations of the table algebra:
 - In general, the tables manipulated by **B** can get infinite, e.g. for unsafe queries, or for negation (may be ok if final result is finite).
 - To supply a real implementation, one must provide a *concrete (sub-)algebra* that provides a particular, efficient, representation of tables, and of the operators over them. Some techniques:
 - Eliminate unsafe queries by static or dynamic checks.
 - Implement important derived operators directly (typically \bowtie , which behaves finitely over finite tables).
 - Push negation to the leaves (define **B** over dualized logical operators, map them to the algebra).
 - Represent certain infinite tables by finite means (e.g. by constraints) and define the operators to work over those.
 - These techniques are being investigated in an implementation (**TQL**) that Giorgio Ghelli is carrying out in Pisa.

Conclusions

- There are many proposals for SSD query languages. Given the power of recursive formulas, we think we can capture many of them (certainly not all), and at least identify a natural spot in design space.
- We have investigated the notion of *SSD query algebra*, which has been missing for too long. (Other proposals are now emerging.)
- We have provided a query language for SSD, a set of logical optimization/rewrite rules, a reference semantics, a query algebra, a specification of algebra-based implementations, and a correctness theorem for the specification w.r.t. the reference semantics.
 - L.Cardelli, G.Ghelli: *A query language based on the ambient logic*. Proc. ESOP'01 (invited paper). <http://www.luca.demon.co.uk>
 - G.Ghelli: *Evaluation of TQL queries*. To appear. <http://www.di.unipi.it/~ghelli/papers.html>

