

Efficiently querying incomplete data

Cristina Sirangelo

LSV, ENS-Cachan

joint work with Amélie Gheerbrant and Leonid Libkin

February 3rd 2015



Plan

- ➔ Introduction
- ➔ Incomplete data model
- ➔ Querying incomplete data
 - ▶ the key to tractability: naïve evaluation
- ➔ What makes naïve evaluation work?
 - ▶ a general framework
- ➔ Applicability of the framework
- ➔ Moving forward

Incomplete data

Data incompleteness: missing/unknown data values, partially available data, ...

Reasons:

mistakes: wrong/missing entries

restrictions on data access

data heterogeneity: data exchange/integration

Incompleteness and data heterogeneity

Assume data has to be transferred from one data source to another

Emp

<i>name</i>	<i>dpt</i>
<i>Green</i>	<i>sales</i>
<i>White</i>	<i>sales</i>
<i>Brown</i>	<i>eng</i>
<i>Black</i>	<i>eng</i>

source database

data restructuring



Employee

<i>name</i>
<i>Green</i>
<i>White</i>
<i>Brown</i>
<i>Black</i>

Manager

<i>mgr</i>	<i>emp</i>
?	<i>Green</i>
?	<i>White</i>
?	<i>Brown</i>
?	<i>Black</i>

target database

Commonly only **some concepts shared** by the two data sources

- ✓ e.g. no information about the manager from the source

Querying incomplete data

Employee	Manager	
<i>name</i>	<i>mgr</i>	<i>emp</i>
<i>Green</i>	<i>?</i>	<i>Green</i>
<i>White</i>	<i>?</i>	<i>White</i>
<i>Brown</i>	<i>?</i>	<i>Brown</i>
<i>Black</i>	<i>?</i>	<i>Black</i>

incomplete database

Q: which employees are managers?



Querying incomplete data

Semantics of query answering

- ▶ How should the result of a query be defined in the presence of incompleteness?

Query evaluation

- ▶ How do we evaluate a query on an incomplete database?
- ▶ Can this be done efficiently ?

Employee	Manager	
<i>name</i>	<i>mgr</i>	<i>emp</i>
Green	?	Green
White	?	White
Brown	?	Brown
Black	?	Black

incomplete database

Q: which employees are managers?



Incompleteness in theory and practice

Incompleteness in database systems

- ▶ **Semantics of query answering**: poorly designed
- ▶ **Query evaluation**: very efficient, optimized query engines

eg:

- In **SQL**, the standard relational database query language, the following are **consistent** statements for sets X, Y

$$|X| > |Y| \text{ and } X - Y = \emptyset$$

- This may occur if Y contains incomplete information (SQL *nulls*)

Incompleteness in theory and practice

Incompleteness in database systems

- ▶ **Semantics of query answering:** poorly designed
- ▶ **Query evaluation:** very efficient, optimized query engines

Theoretical framework for incompleteness

[Imielinski-Lipski, Abiteboul-Kanellakis-Grahne, etc. 80's]

- ▶ **Semantics of query answering:** clean framework, suitable semantics
- ▶ **Query evaluation:** hard

Incompleteness in theory and practice

Incompleteness in database systems

- ▶ **Semantics of query answering:** poorly designed
- ▶ **Query evaluation:** very efficient, optimized query engines

Theoretical framework for incompleteness

[Imielinski-Lipski, Abiteboul-Kanellakis-Grahne, etc. 80's]

- ▶ **Semantics of query answering:** clean framework, suitable semantics
- ▶ **Query evaluation:** hard

Bridging the gap between theory and systems:
answer queries correctly, use classical query engines

not satisfactorily addressed even in the simplest data model

Plan

- ➔ Introduction
- ➔ **Incomplete data model**
- ➔ Querying incomplete data
 - ▶ the key to tractability: naïve evaluation
- ➔ What makes naïve evaluation work?
 - ▶ a general framework
- ➔ Applicability of the framework
- ➔ Moving forward

Incomplete relational data

Database schema (relational signature) σ : a set of relation symbols, with arities

eg $\sigma = \{ \mathbf{Employee}, \mathbf{Manager} \}$ $arity(\mathbf{Employee})=1$, $arity(\mathbf{Manager})=2$

Incomplete database instance (naive table) of schema σ [Imielinski, Lipski '84]:

associates to each relation symbol R of σ a finite subset of $(Const \cup Var)^{arity(R)}$

	Employee	Manager
I	Green	Green x_1
	x_1	x_1 Brown
	Brown	Green x_2

$Const$: a countably infinite set of constants

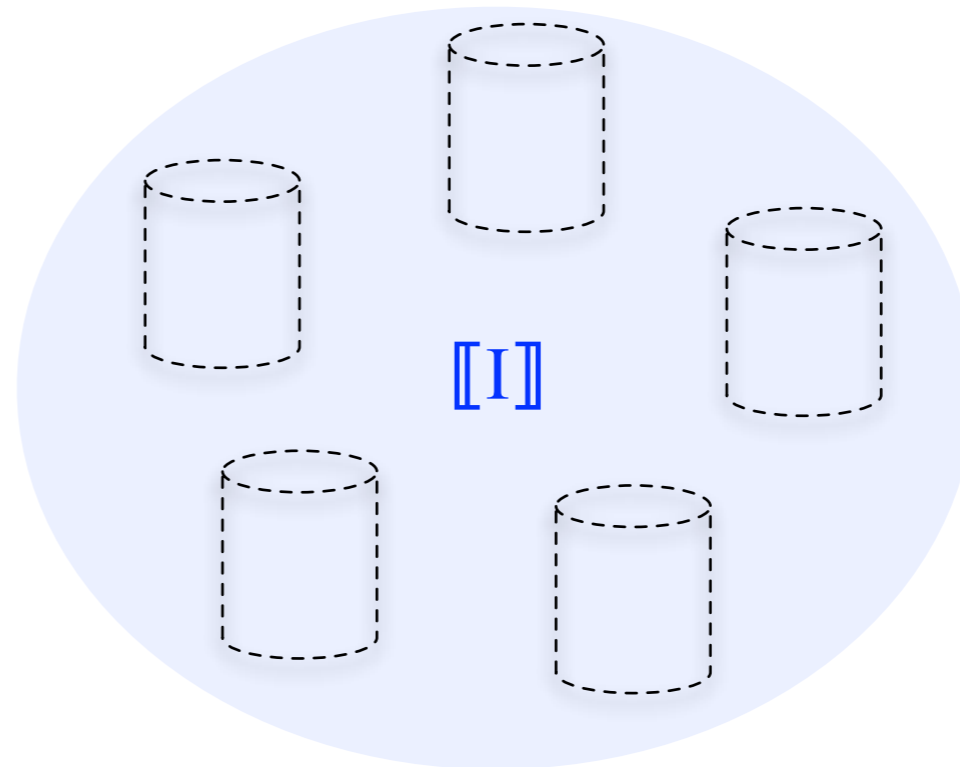
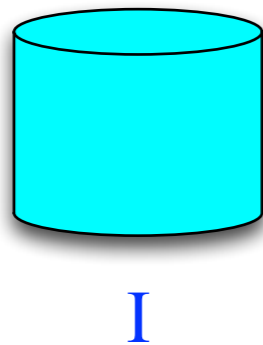
Var : a countably infinite set of variables (nulls)

Complete instance: over $Const$

$dom(I)$: the subset of $Const \cup Var$ occurring in I

Semantics of incompleteness

Any incomplete database represents a set of complete databases (*possible worlds*)



Semantics of incompleteness:

a function $[\]$ associating with each incomplete database a set of complete databases

Semantics of incompleteness

Employee

Green
x_1
Brown

Manager

Green	x_1
x_1	Brown
Green	x_2

I

Three well known relational semantics:

- ▶ **OWA** (*Open World Assumption*) [Imielinski-Lipski '84]
- ▶ **CWA** (*Closed World Assumption*) [Reiter '77, Imielinski-Lipski '84]
- ▶ **WCWA** (*Weak Closed World Assumption*) [Reiter '77]

Semantics of incompleteness

Employee

Green
x_1
Brown

Manager

Green	x_1
x_1	Brown
Green	x_2

I

$x_1=White$ $x_2=Black$

$x_1=x_2=Brown$

$x_1=White$ $x_2=Black$

Green
White
Brown

Green	White
White	Brown
Green	Black

Green
Brown

Green	Brown
Brown	Brown
Black	Brown

Green
White
Brown
Black
Smith

Green	White
White	Brown
Green	Black

⋮

$\llbracket I \rrbracket_{OWA}$

OWA:

$$\llbracket I \rrbracket_{OWA} = \{ D \text{ over } Const \mid D \supseteq v(I) \text{ for some } v: Var \rightarrow Const \}$$

v : valuation

Interpretation of incompleteness:

- missing data values, missing tuples

Semantics of incompleteness

Employee

Green
x_1
Brown

Manager

Green	x_1
x_1	Brown
Green	x_2

I

$x_1=White$ $x_2=Black$

$x_1=x_2=Brown$

$x_1=White$ $x_2=Black$

Green
White
Brown

Green	White
White	Brown
Green	Black

Green
Brown

Green	Brown
Brown	Brown

Green
White
Brown

Green	White
White	Brown
Green	Black

⋮

$\llbracket I \rrbracket_{CWA}$

CWA:

$$\llbracket I \rrbracket_{CWA} = \{ D \text{ over } Const \mid D = v(I) \text{ for some } v: Var \rightarrow Const \}$$

Interpretation of incompleteness:

- missing data values
- no missing tuples

Semantics of incompleteness

Employee

Green
x_1
Brown

Manager

Green	x_1
x_1	Brown
Green	x_2

I

$x_1=White$ $x_2=Black$

$x_1=x_2=Brown$

$x_1=White$ $x_2=Black$

Green
White
Brown

Green	White
White	Brown
Green	Black

Green
Brown

Green	Brown
Brown	Brown
Green	Green

Green
White
Brown
Black

Green	White
White	Brown
Green	Black
Green	Brown

⋮

$\llbracket I \rrbracket_{WCWA}$

WCWA:

$$\llbracket I \rrbracket_{WCWA} = \{D \text{ over } Const \mid D \supseteq v(I), \text{ dom}(D) = \text{dom}(v(I)) \text{ for some } v: Var \rightarrow Const\}$$

Interpretation of incompleteness:

- missing data values, missing tuples
- no missing domain elements

Plan

- ➔ Introduction
- ➔ Incomplete data model
- ➔ Querying incomplete data
 - ▶ the key to tractability: naïve evaluation
- ➔ What makes naïve evaluation work?
 - ▶ a general framework
- ➔ Applicability of the framework
- ➔ Moving forward

Queries

Query over σ :

a mapping Q associating to each complete instance I of σ a relation over $\text{dom}(I)$

Q: which employees are managers?

Employee

Green
White
Brown

Manager

Green	White
White	Brown
Green	Black



Green
White

I

$Q(I)$

usually expressed in fragments of *First Order logic (FO)*

$$\varphi_Q(x) = \text{Employee}(x) \wedge \exists y \text{ Manager}(x, y)$$

Queries

Boolean query over σ :

a mapping Q associating to each complete instance I of σ a value in $\{true, false\}$

$Q: \exists x, y, z (\text{Manager}(x, y) \wedge \text{Manager}(z, x))$ (there is a manager who has a manager)

Employee

Green
White
Brown

Manager

Green	White
White	Brown
Green	Black



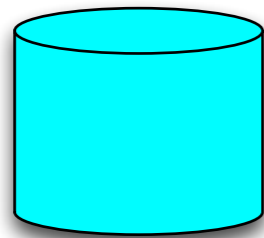
true

I

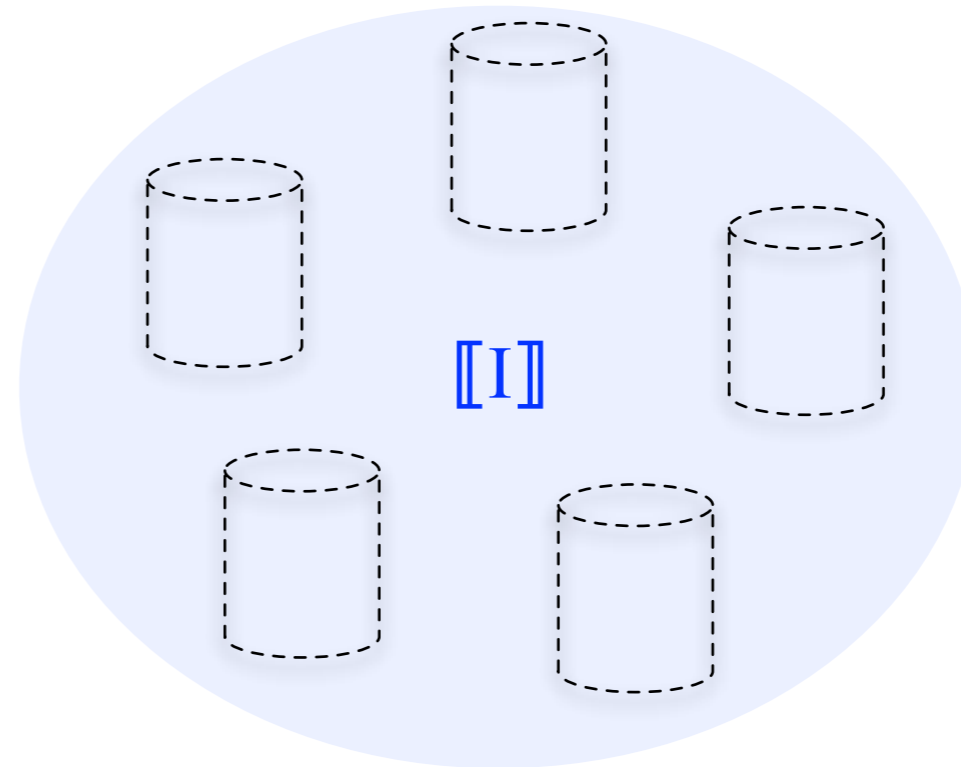
$Q(I)$

Querying incomplete databases

Semantics of query answering: *certain answers*



I



Q (Boolean)

$$\text{cert}_Q(I) = \bigwedge_{D \in \llbracket I \rrbracket} Q(D)$$

Certain answers

Example Q : “There is a manager who has a manager”

$$\exists x, y, z (\text{Manager}(x, y) \wedge \text{Manager}(z, x))$$

Employee

Green
x_1
Brown

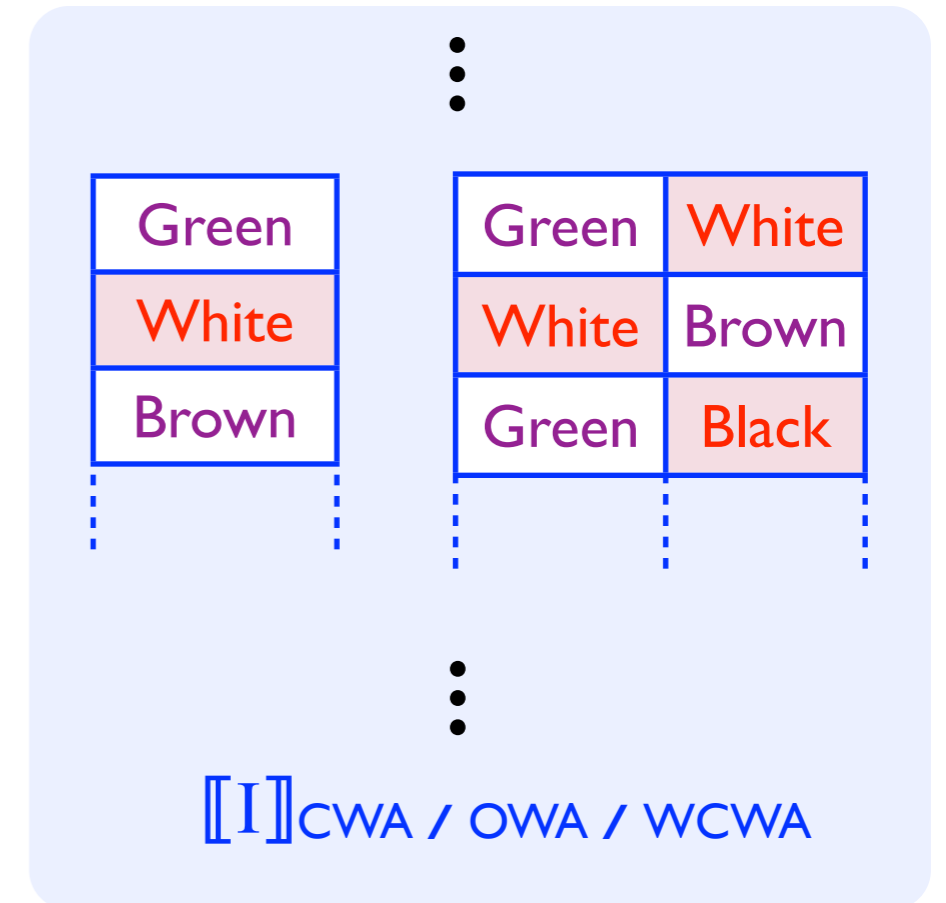
Manager

Green	x_1
x_1	Brown
Green	x_2

I

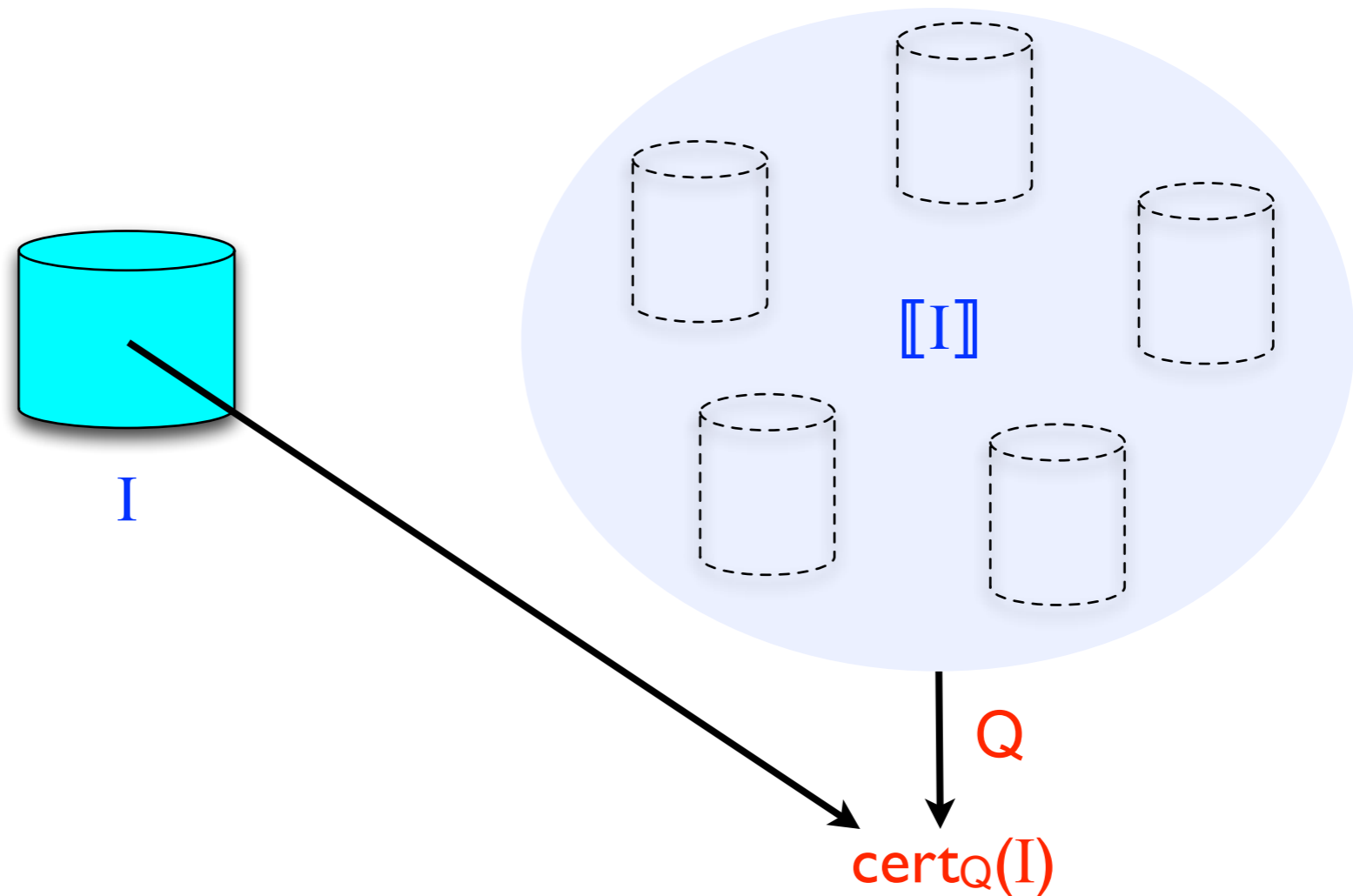
$\text{cert}_Q(I) = \text{true}$

under either CWA, OWA and WCWA



Computing certain answers

Need to use the available (incomplete) data

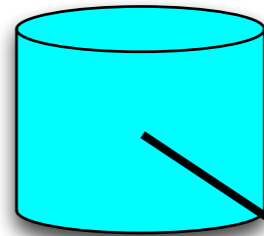


Computing certain answers on I : usually hard

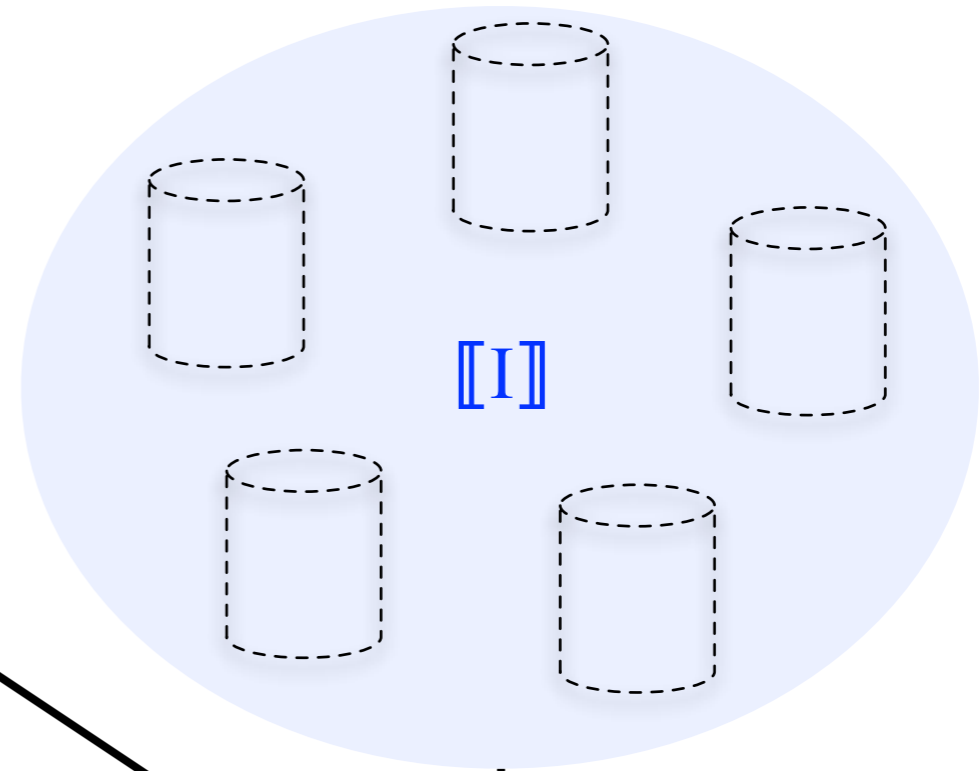
- ▶ from coNP-complete to undecidable for FO [Imielinski-Lipski '84, Abiteboul et al '91]

Naive evaluation

Naive evaluation works for Q :



I



Q

Q

$Q(I)$: Q evaluated directly on I ,
as if variables were new distinct constants

($x_i \neq x_j$ for $i \neq j$

$x_i \neq c$ for all $c \in Const$)

$Q(I) = cert_Q(I)$ for all I

Naive evaluation

Example Q : “There is a manager who has a manager”

$$\exists x, y, z (\text{Manager}(x, y) \wedge \text{Manager}(z, x))$$

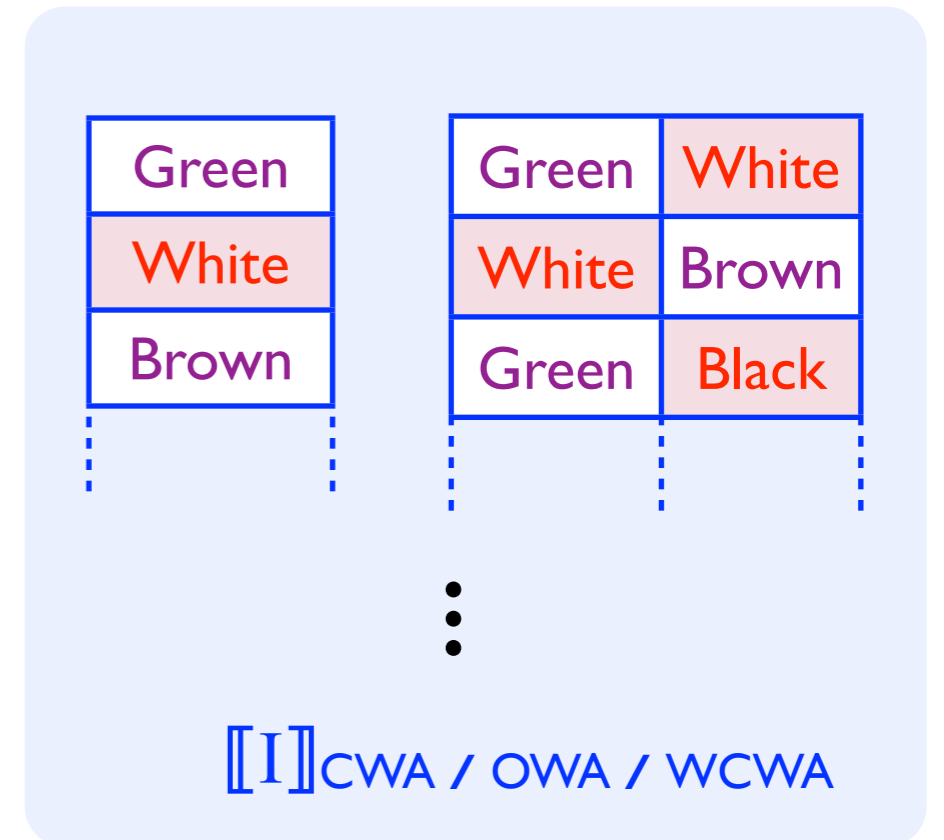
Employee

Green
x_1
Brown

Manager

Green	x_1
x_1	Brown
Green	x_2

I



$Q(I) = \text{true}$

$\text{cert}_Q(I) = \text{true}$ under CWA, OWA and WCWA

Naive evaluation

Example Q : “There is a manager who has a manager”

$$\exists x, y, z (\text{Manager}(x, y) \wedge \text{Manager}(z, x))$$

Employee

Manager

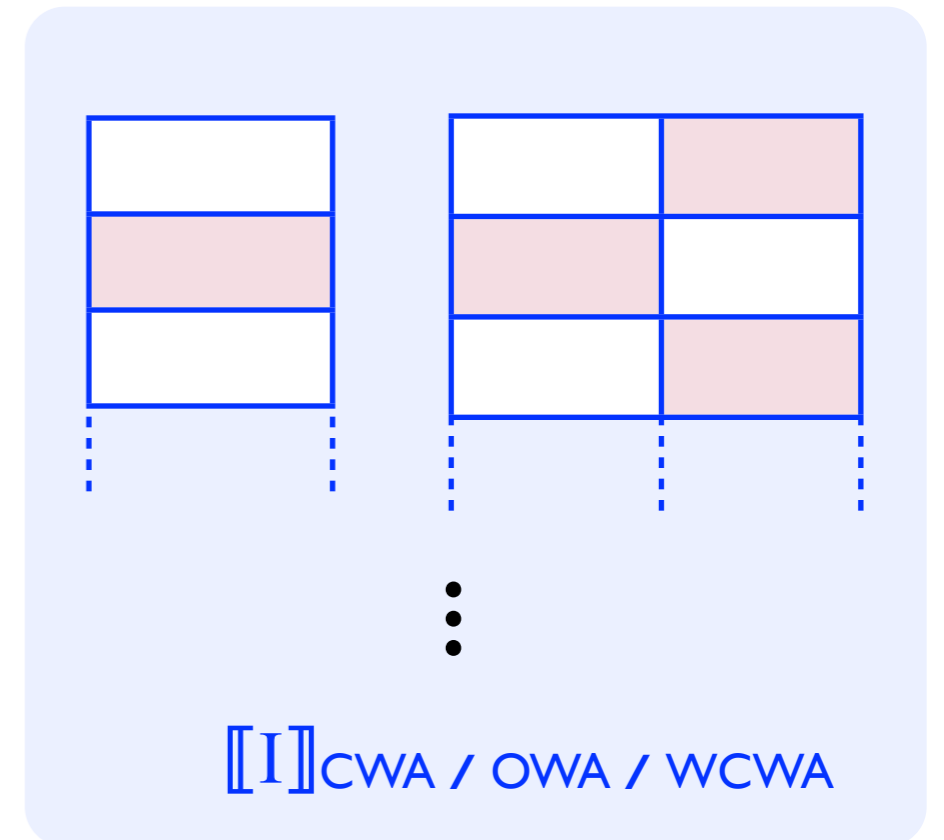
I

Generalizing:

$$Q(I) = \text{cert}_Q(I) \text{ for all } I$$

⇒ naive evaluation works for Q

under CWA, OWA and WCWA



Naïve evaluation in theory and practice

Certain answers: an entailment problem (checking that I entails Q) **HARD**

Naïve evaluation: a model-checking problem (checking $I \models Q$) **EFFICIENT**

- ▶ **PTIME** in the size of the instance for FO queries
- ▶ based on classical query evaluation algorithms of database engines
- ▶ can benefit from query optimization techniques

Naïve evaluation works

correct query answering semantics, classical query evaluation algorithms /
entailment reduces to (straightforward) model-checking

clearly not always possible ! (undecidable vs. PTIME)

Naïve evaluation does not always work

A concrete example Q: “All employees are managers”

$$\forall x(\text{Employee}(x) \rightarrow \exists y \text{ Manager}(x, y))$$

Employee

Green
x_1
Brown

Manager

Green	x_1
x_1	Brown
Brown	x_2

I

$Q(I) = \text{true}$

Green
White
Brown
Black

Green	White
White	Brown
Brown	Black

D

⋮

$\llbracket I \rrbracket_{\text{OWA/WCWA}}$

$\text{cert}_Q(I) = \text{false}$ under OWA and WCWA

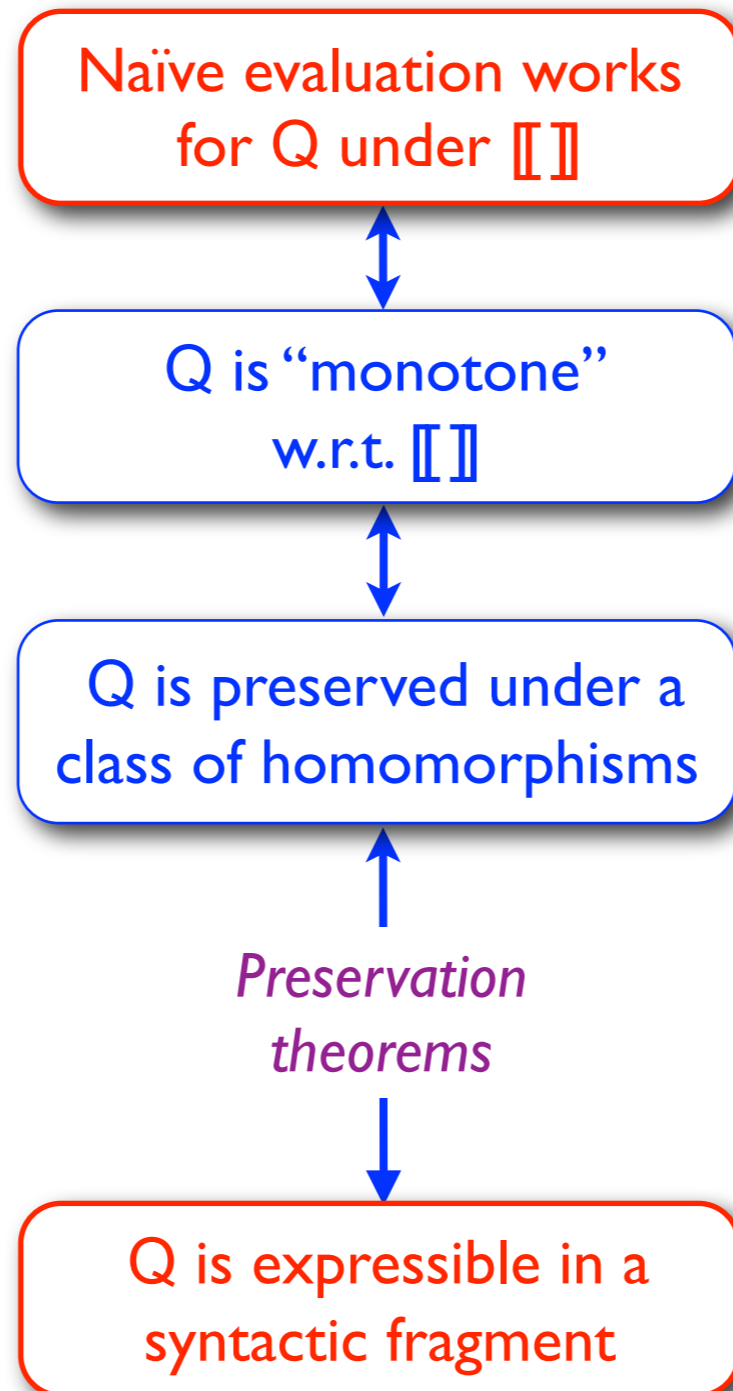
\Rightarrow naïve evaluation does not work for Q under OWA and WCWA

Plan

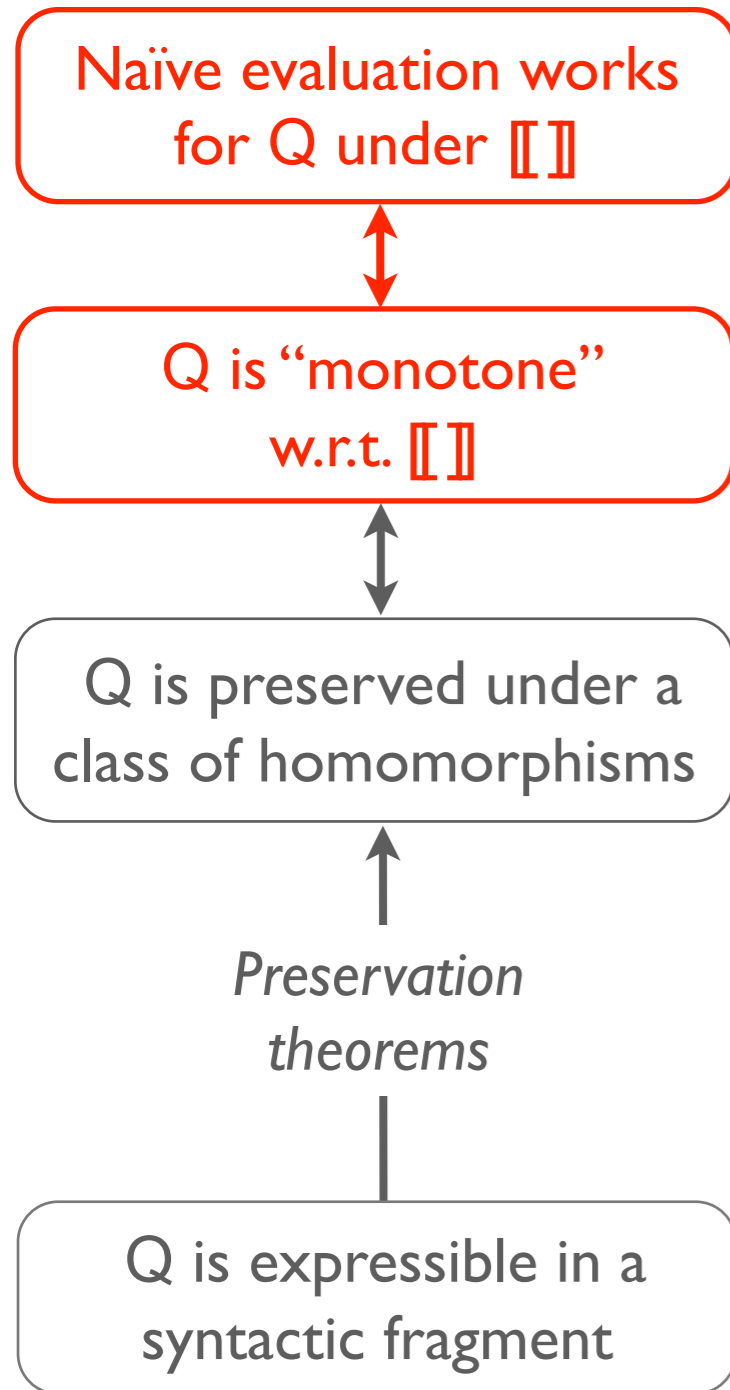
- ➔ Introduction
- ➔ Incomplete data model
- ➔ Querying incomplete data
 - ▶ the key to tractability: naïve evaluation
- ➔ What makes naïve evaluation work?
 - ▶ a general framework
- ➔ Applicability of the framework
- ➔ Moving forward

Relating naïve evaluation and syntactic fragments

A unified framework for relating naïve evaluation and syntactic fragments for several possible semantics:



Monotonicity and preservation



Shown in a very general setting subsuming every data model / semantics of incompleteness (even beyond relational databases)

Naïve evaluation and monotonicity

Database domain: a quadruple $\langle \mathcal{D}, C, \llbracket \cdot \rrbracket, \approx \rangle$

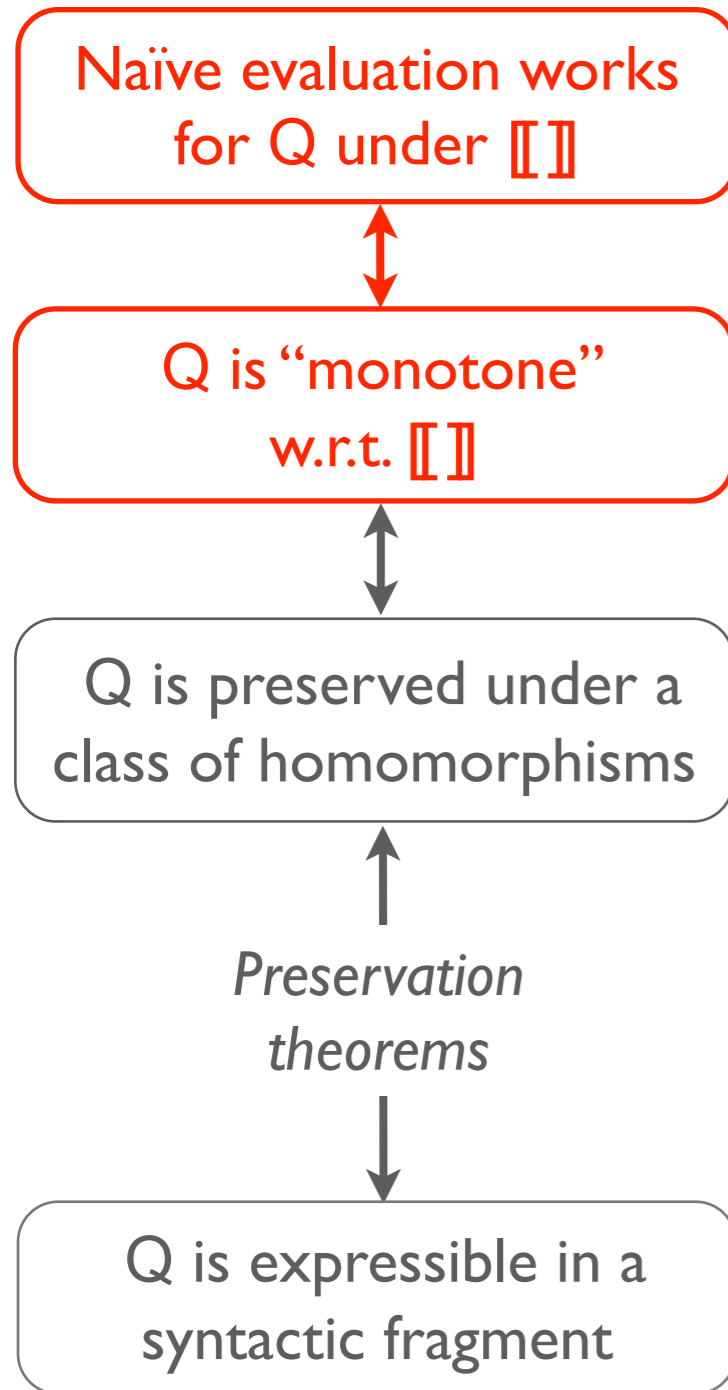
	description	example
\mathcal{D} : a set	database objects (complete and incomplete)	all naïve tables over a fixed schema σ
C : a subset of \mathcal{D}	complete database objects	all complete instances over σ
$\llbracket \cdot \rrbracket$: $\mathcal{D} \rightarrow 2^C$	semantics of incompleteness	$\llbracket \cdot \rrbracket_{\text{OWA}}, \llbracket \cdot \rrbracket_{\text{CWA}}, \text{etc.}$
\approx : an equivalence relation on \mathcal{D}	equivalence of objects (w.r.t. queries)	isomorphism of relational instances

Boolean query: $Q : \mathcal{D} \rightarrow \{\text{true}, \text{false}\}$

generic : $x \approx y$ implies $Q(x) = Q(y)$

monotone w.r.t. $\llbracket \cdot \rrbracket$: $y \in \llbracket x \rrbracket$ implies $Q(x) \Rightarrow Q(y)$

Naïve evaluation and monotonicity



Over a **saturated** database domain, if Q is a generic Boolean query:

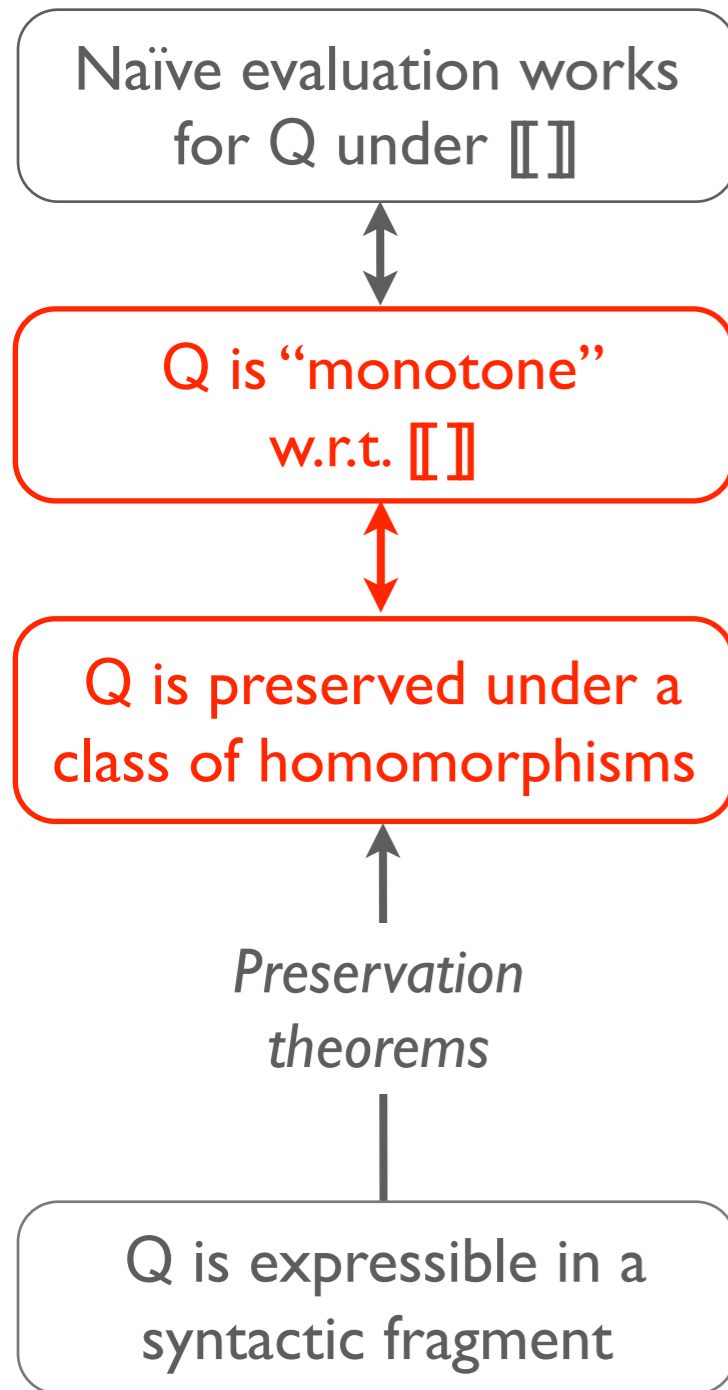
Naïve evaluation works for Q iff Q is monotone w.r.t. $\llbracket \cdot \rrbracket$

Saturation property for $\langle \mathcal{D}, \mathcal{C}, \llbracket \cdot \rrbracket, \approx \rangle$:

For all $x \in \mathcal{D}$ there exists $y \in \llbracket x \rrbracket$ $y \approx x$

holds for most common semantics

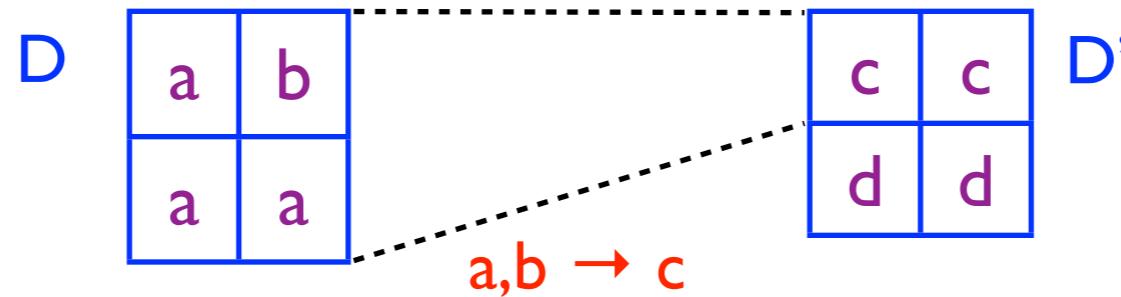
Monotonicity and preservation



- ▶ monotonicity = "preservation" under the semantics
- ▶ relational semantics: usually homomorphism-based

Monotonicity and preservation

homomorphism $D \rightarrow D'$: a mapping $h: \text{dom}(D) \rightarrow \text{dom}(D')$ s.t. $h(D) \subseteq D'$



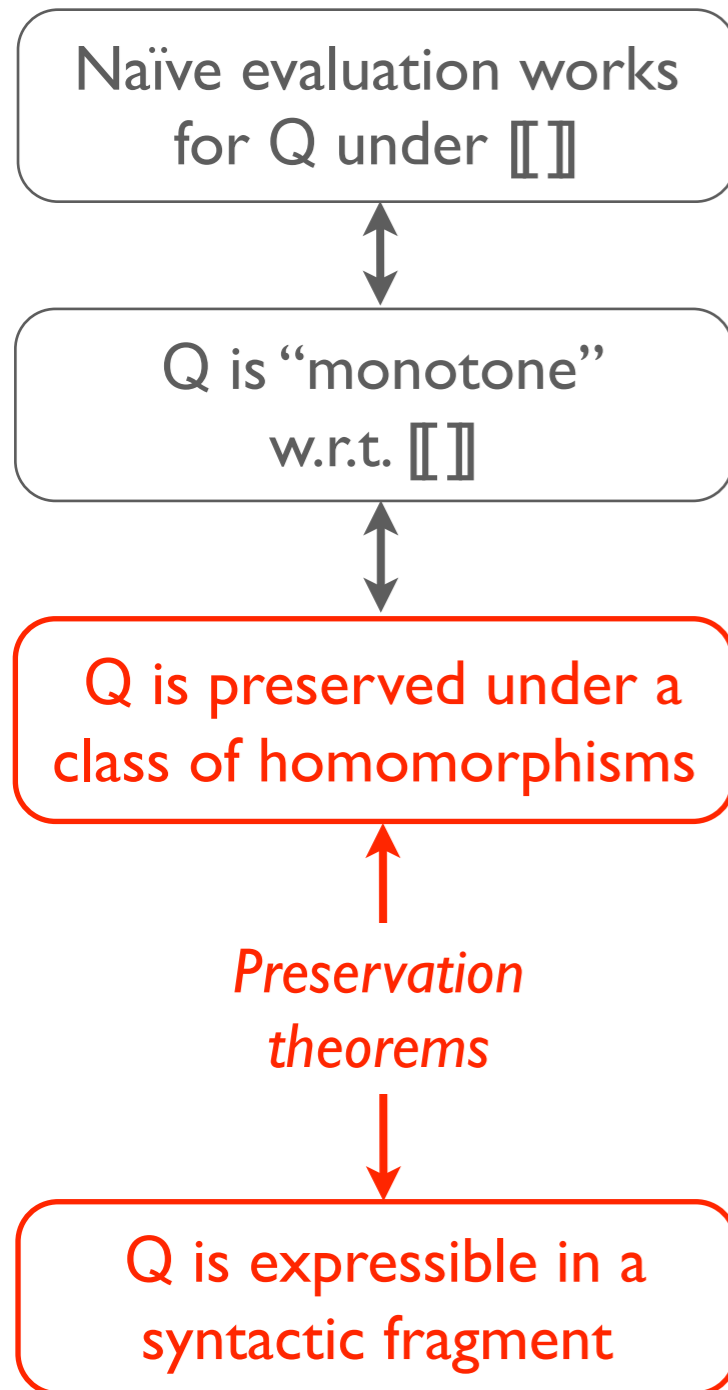
Q preserved under homomorphism:

$D \rightarrow D'$ implies $Q(D) \Rightarrow Q(D')$ for all D

Many variants: onto homomorphism, strong onto homomorphism, ...

monotonicity w.r.t different semantics \leftrightarrow
preservation under different notions of homomorphism

Preservation and syntactic fragments



Preservation theorems

- ▶ syntactic characterizations of preservation properties of queries in a given logic
- ▶ classical results in (finite) model theory

Preservation and syntactic fragments

Homomorphism Preservation Theorem: over arbitrary structures
an FO query Q is preserved under homomorphism iff Q is in $\exists\text{Pos}$

- ▶ recently proved over finite structures [Rossman '08]

$\exists\text{Pos}$: $\{\exists, \wedge, \vee\}$ -FO (Unions of Conjunctive Queries)

Preservation and syntactic fragments

Homomorphism Preservation Theorem: over arbitrary structures
an FO query Q is **preserved under homomorphism** iff Q is in $\exists\text{Pos}$

- ▶ recently proved over finite structures [Rossman '08]

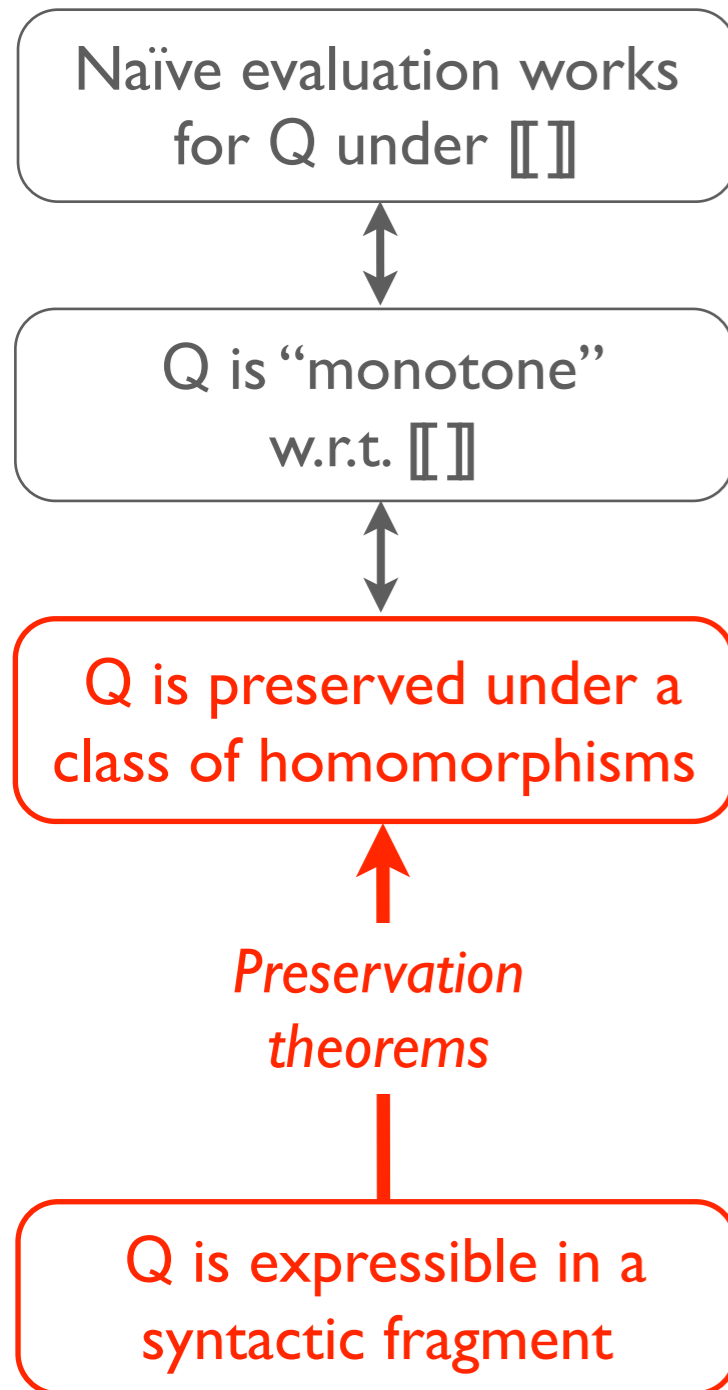
$\exists\text{Pos}$: $\{\exists, \wedge, \vee\}$ -FO (Unions of Conjunctive Queries)

Lyndon Positivity Theorem [Lyndon '59]: over arbitrary structures
an FO query Q is **preserved under onto homomorphism** iff Q is in Pos

- ▶ fails in the finite [Ajtai-Gurevich 87, Rosen '95, Stolboushkin '95]

Pos : $\{\exists, \forall, \wedge, \vee\}$ -FO

Preservation and syntactic fragments



Preservation theorems:

(Syntax \Rightarrow Preservation) holds in the finite as well

\Rightarrow classes of queries where naïve evaluation works

Naïve evaluation and syntactic fragments

Three well known semantics as instances of our framework

Naïve evaluation works under:

Naïve evaluation works for Q under $[[\]]$

Q is “monotone” w.r.t. $[[\]]$

Q is preserved under a class of homomorphisms

Preservation theorems

Q is expressible in a syntactic fragment

OWA

Preservation under homomorphism

\exists Pos

[Rossman]

WCWA

Preservation under onto homomorphism

Pos

[Lyndon]

CWA

Preservation under “strong onto” homomorphism

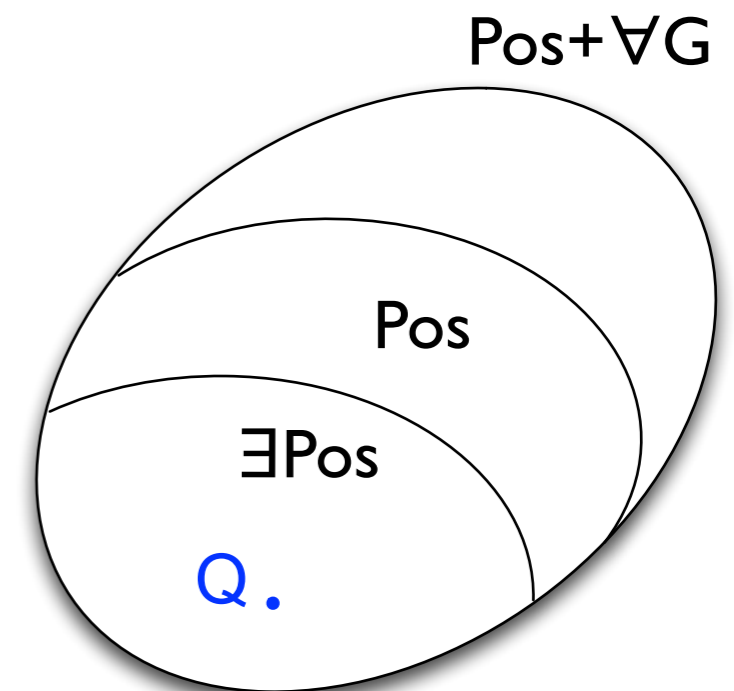
Pos+ \forall G

[Gheerbrant, Libkin, S.]

Examples revisited

Q : “There is a manager who has a manager”

$\exists x, y, z (\text{Manager}(x, y) \wedge \text{Manager}(z, x))$



\Rightarrow naïve evaluation works for Q under CWA OWA, WCWA

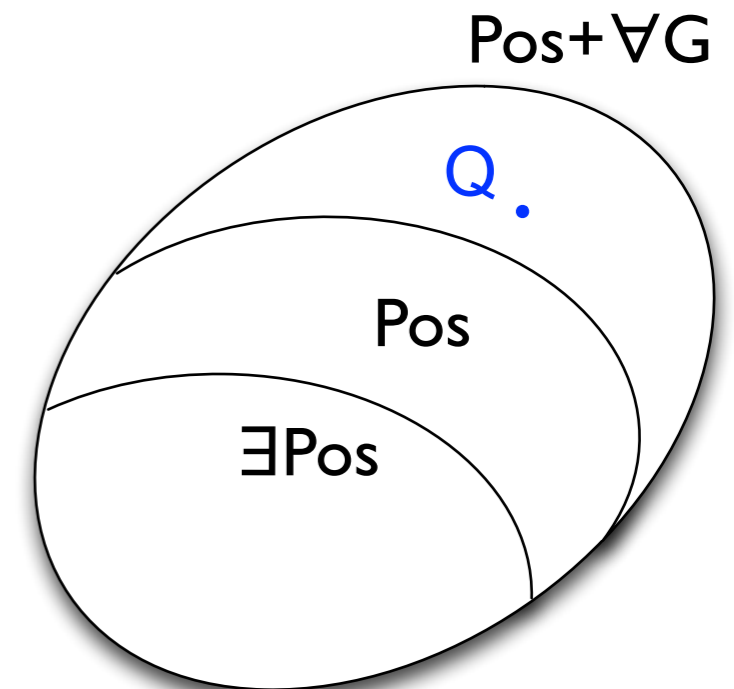
Examples revisited

Q : “All employees are managers”

$$\forall x(\text{Employee}(x) \rightarrow \exists y \text{ Manager}(x, y))$$

⇒ naïve evaluation works for Q under CWA

(recall: not true under OWA, nor under WCWA)



$\text{Pos} + \forall G$ extends Pos with a limited form of negation (universal guards)

- ▶ a very natural fragment

Naïve evaluation works well beyond $\exists \text{Pos}$ under other semantics than OWA

Plan

- ➔ Introduction
- ➔ Incomplete data model
- ➔ Querying incomplete data
 - ▶ the key to tractability: naïve evaluation
- ➔ What makes naïve evaluation work?
 - ▶ a general framework
- ➔ **Applicability of the framework**
- ➔ Moving forward

Beyond OWA, CWA and WCWA

Semantics of incompleteness have been considered in several contexts:

- ▶ programming semantics, logic programming, data exchange,...

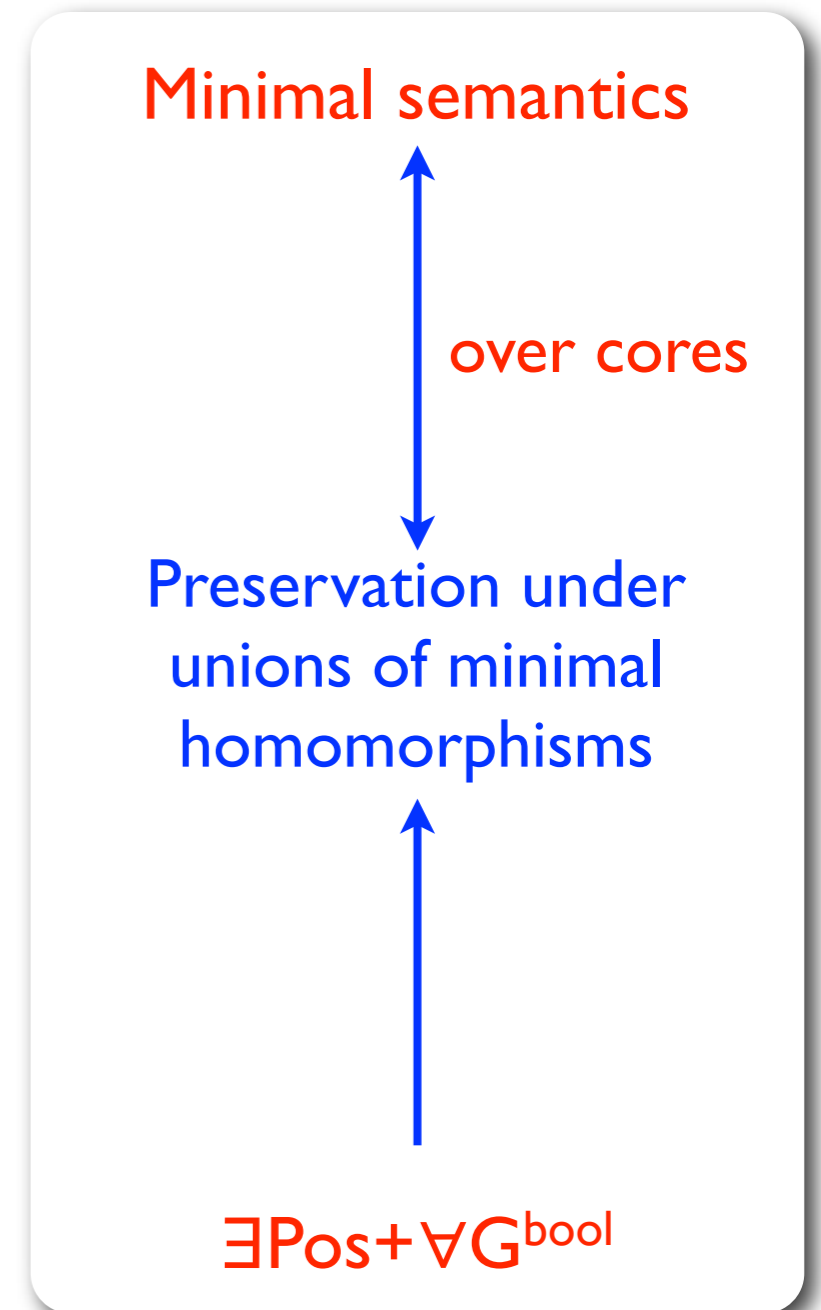
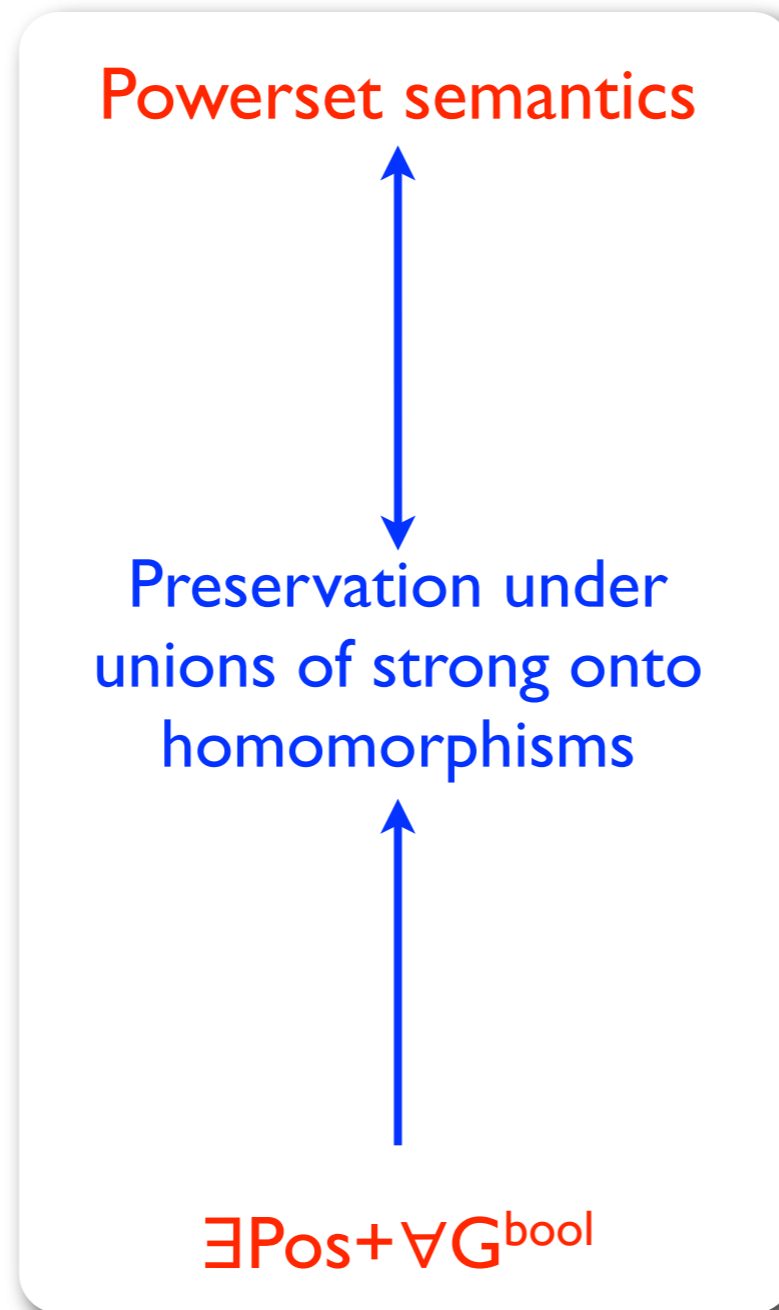
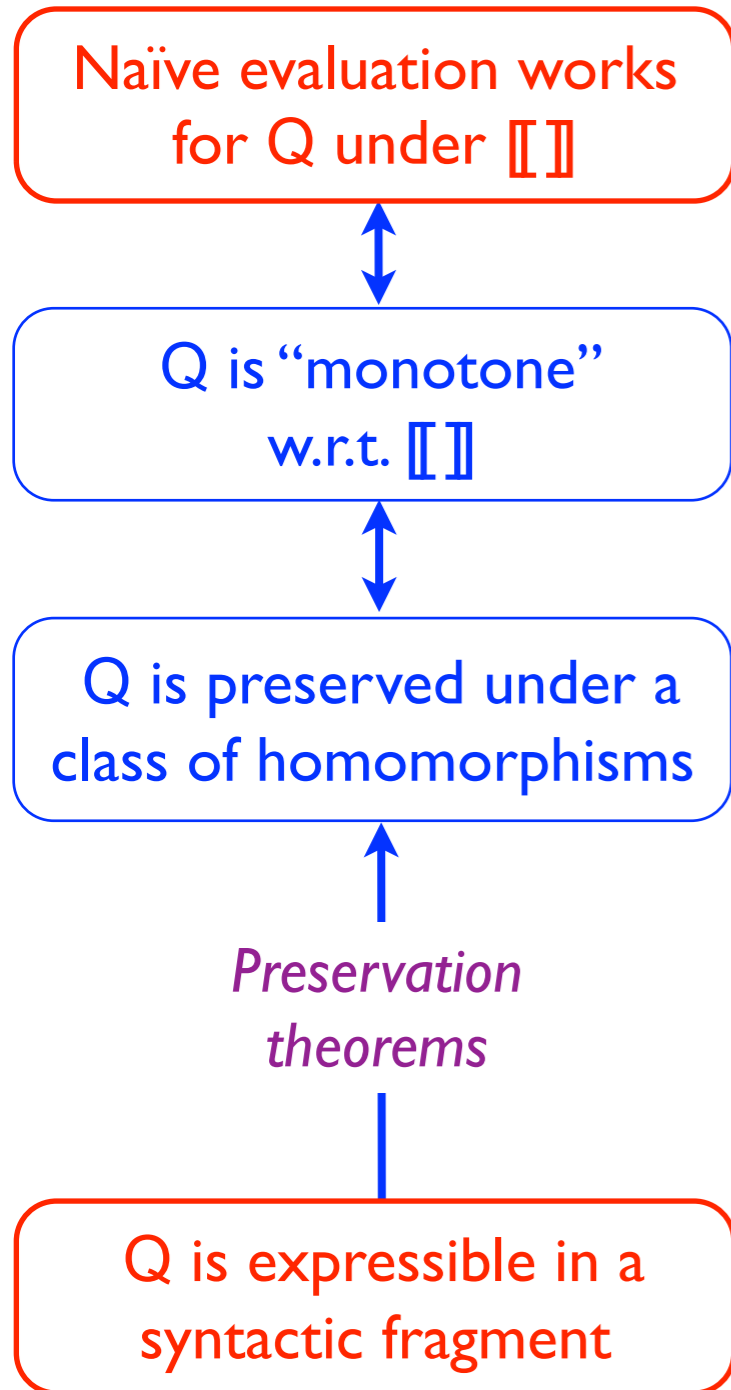
[Minker'82, Ohori'90, Rounds'91, Libkin'95, Hernich'11]

- powerset semantics

- minimal semantics

Beyond OWA, CWA and WCWA

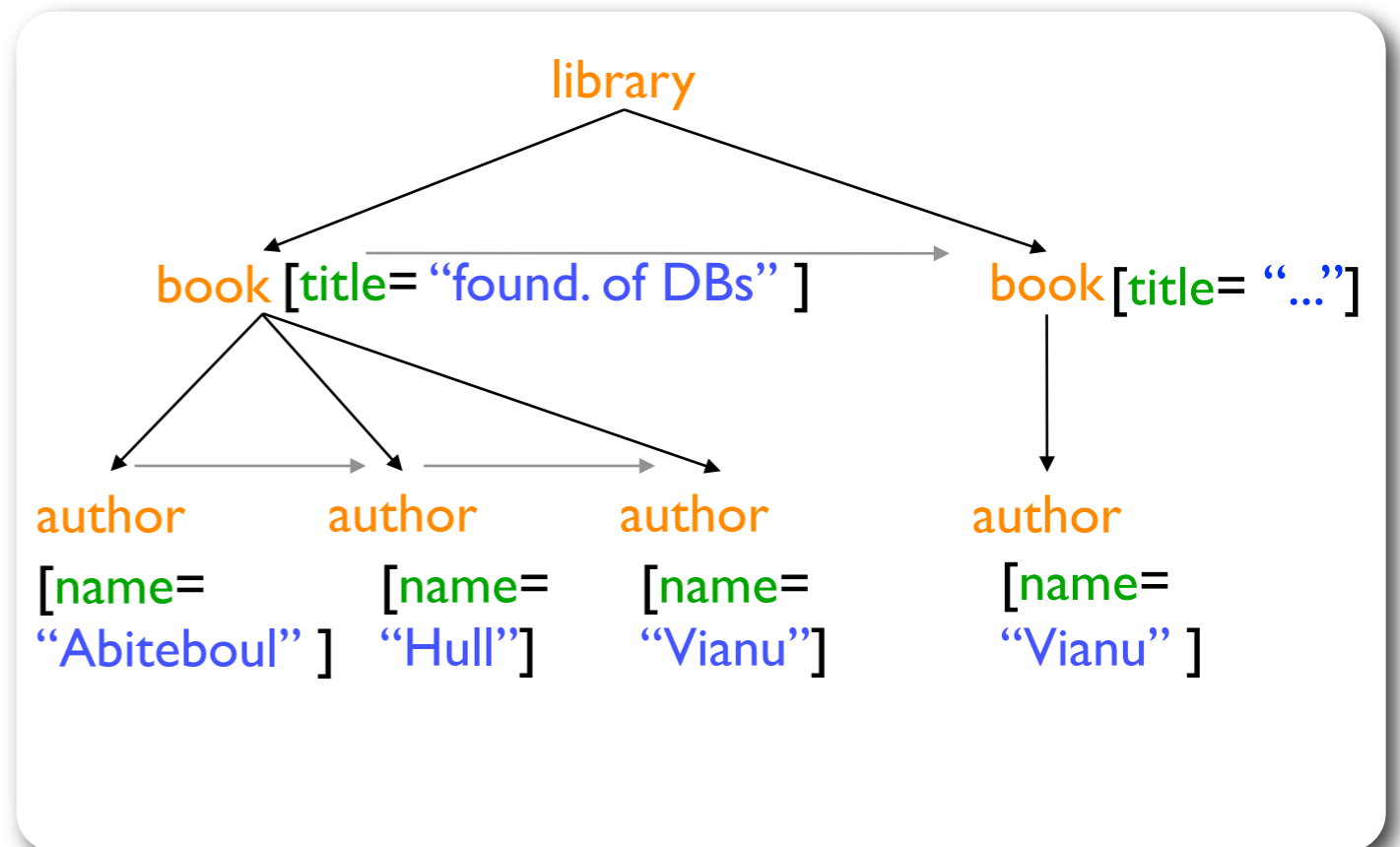
Naïve evaluation works under:



Beyond the relational data model

XML: hierarchically structured data

```
<?xml version="1.0" encoding="UTF-8"?>
<library>
  <book title="found. of DBs">
    <author name="Abiteboul" > </author>
    <author name="Hull"> </author>
    <author name="Vianu"> </author>
  </book>
  <book title="...">
    <author name="Vianu" > </author>
    ...
  </book>
</library>
```

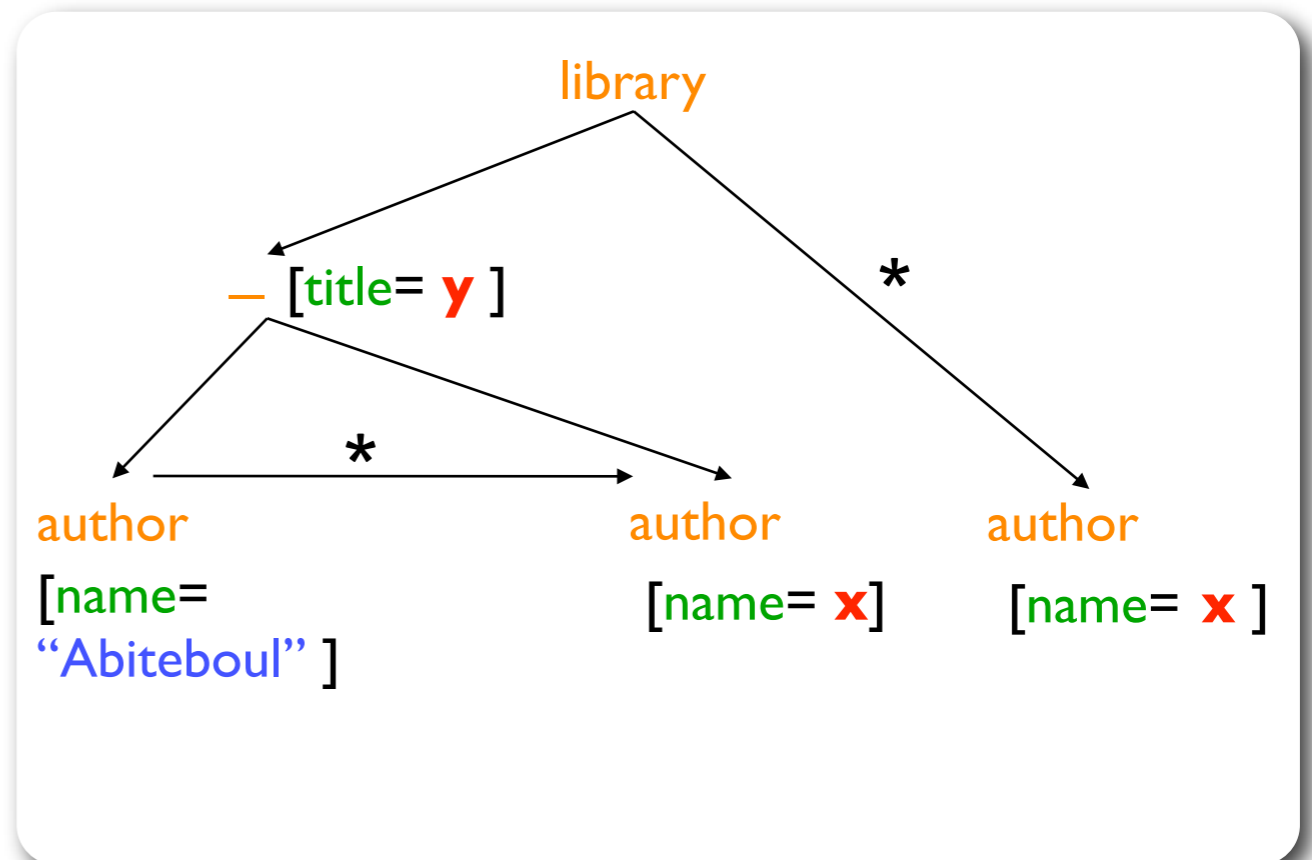


modeled as trees with data values associated to nodes

Beyond the relational data model

Incomplete XML based on a form of “tree patterns” [Barcelo-Libkin-Poggi-S. '10]

- ▶ missing data values
- ▶ missing nodes
- ▶ missing structural information
- ✓ labels
- ✓ parent-child, next-sibling relationships
- ✓ etc.



Tree-pattern-queries: the analog of \exists Pos on trees

Beyond the relational data model

The analog of naïve evaluation works for tree-pattern-queries under OWA on *rigid tree patterns* [Barcelo-Libkin-Poggi-S. '10]

- ▶ rigidity: essentially avoids structural incompleteness

Our framework explains this result:

- ▶ database domain:
 - ✓ the set of complete/incomplete trees,
 - ✓ OWA semantics: homomorphism-based
- ▶ tree pattern queries are preserved under homomorphisms of trees
- ▶ *rigidity ensures the saturation property*

Moving forward

Naïve evaluation on combinations of data models/semantics, e.g

- ➔ XML/CWA
- ➔ graph-structured data

Query languages beyond FO

- ➔ fixed-point logics, fragments of SO, etc.

Naïve evaluation over restricted instances

- ➔ Applications: data integration/exchange

Beyond naïve-evaluation

- ➔ rewriting of the query/instance
(classical in ontology-based query answering)

Thank you!

Real life paradoxes

- SQL adopts a three-valued logic
 - ▶ essentially any comparison involving null values evaluates to *unknown*
- An SQL condition checking $X - Y \neq \emptyset$

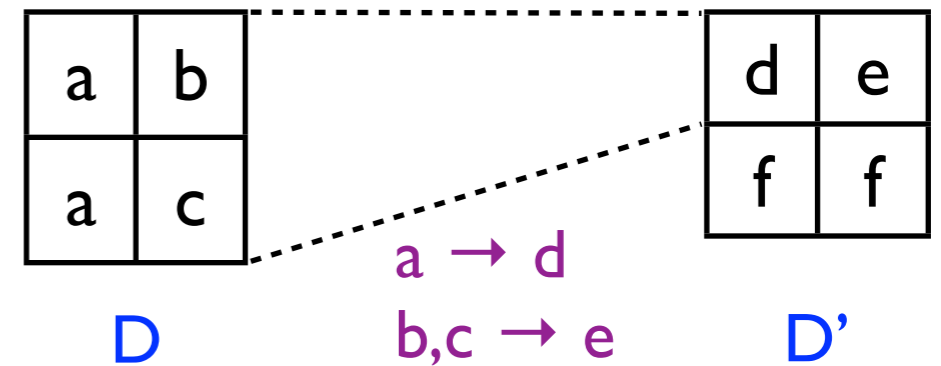
```
EXISTS ( SELECT X.A FROM X
WHERE X.A NOT IN ( SELECT Y.A FROM Y ) )
```
- $X.A = \{1, 2, 3, \dots, N\}$ and $Y.A = \{\text{null}\}$, then $X - Y = \emptyset$ no matter what N is!
- That's how SQL programs work: this is part of the SQL 1999 ANSI Standard

Homomorphisms

Homomorphism $D \rightarrow D'$:

a mapping $h: \text{dom}(D) \rightarrow \text{dom}(D')$ s.t.

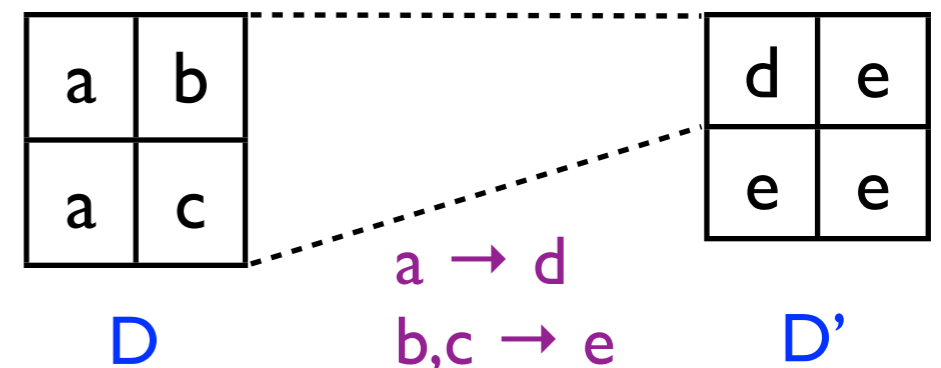
$$h(D) \subseteq D'$$



Onto homomorphism $D \rightarrow D'$:

a homomorphism $h: D \rightarrow D'$ s.t.

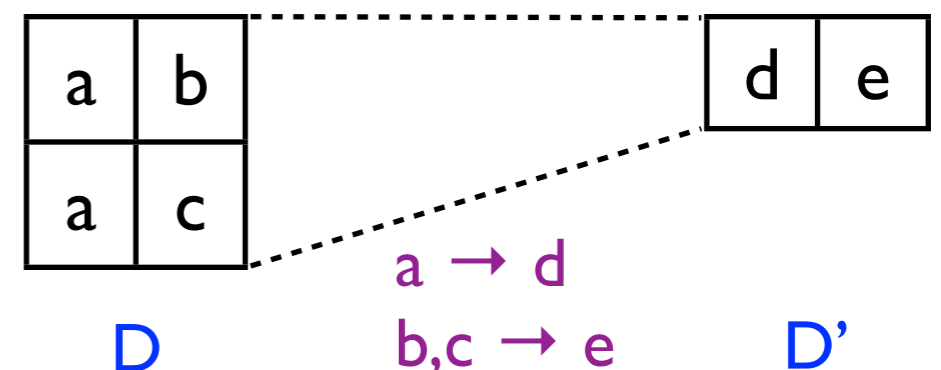
$$h(\text{dom}(D)) = \text{dom}(D')$$



Strong onto homomorphism $D \rightarrow D'$:

a homomorphism $h: D \rightarrow D'$ s.t.

$$h(D) = D'$$



Homomorphisms

- ▶ **Union of strong onto homomorphisms** $D \rightarrow D' : \bigcup_i h_i(D) = D'$
- ▶ **D-minimal homomorphism** h on D :
there exists no h' , preserving all constants preserved by h , s.t. $h'(D) \subsetneq h(D)$
- ▶ **Union of minimal homomorphisms** $D \rightarrow D' : \bigcup_i h_i(D) = D'$
with $h_1 \dots h_n$ **D-minimal** and **preserving the same constants**

Homomorphism-based relational semantics



Homomorphism-based relational semantics



- Essentially based on **two steps**: 1) **valuation of nulls** 2) **extension of the instance**
- Other well-known semantics follow the same paradigm:

Homomorphism-based relational semantics

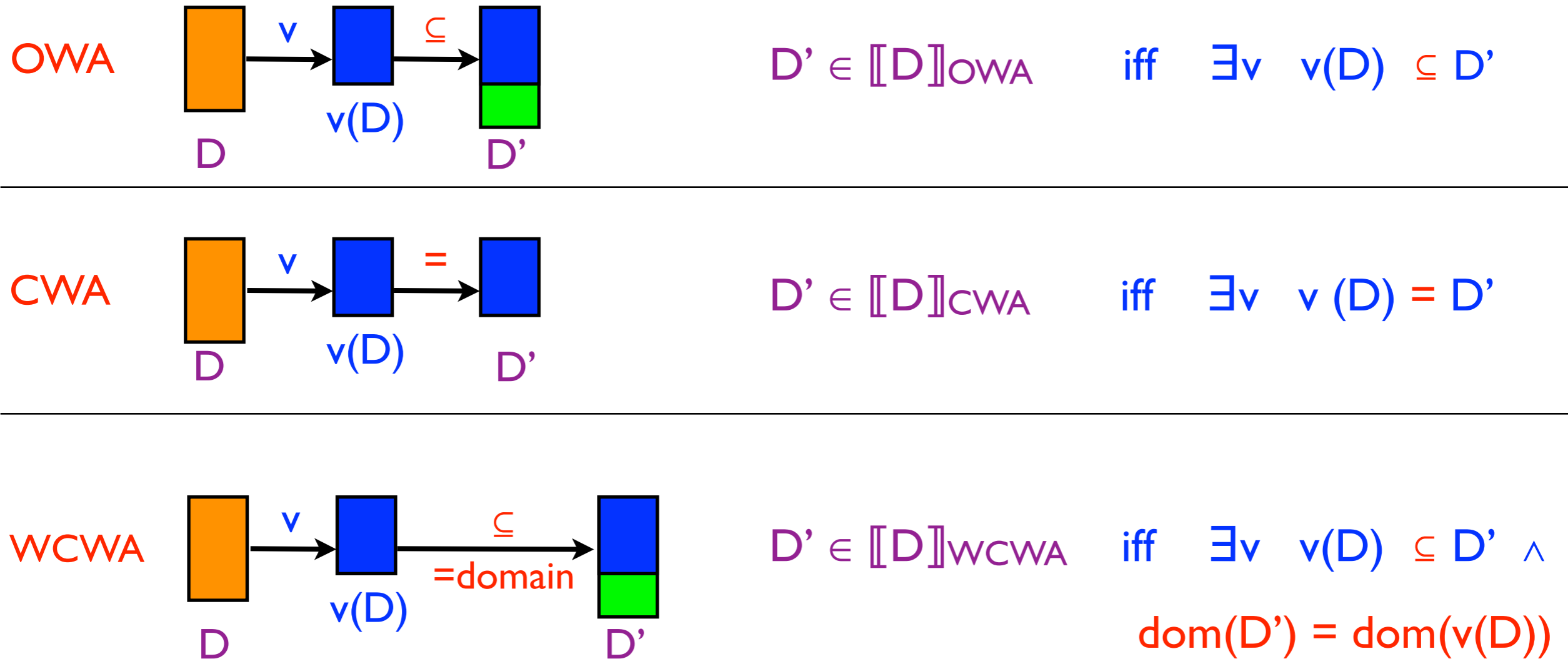


- Essentially based on **two steps**: 1) **valuation of nulls** 2) **extension of the instance**
- Other well-known semantics follow the same paradigm:

Weak Closed World Assumption [Reiter 77]



Homomorphism-based relational semantics



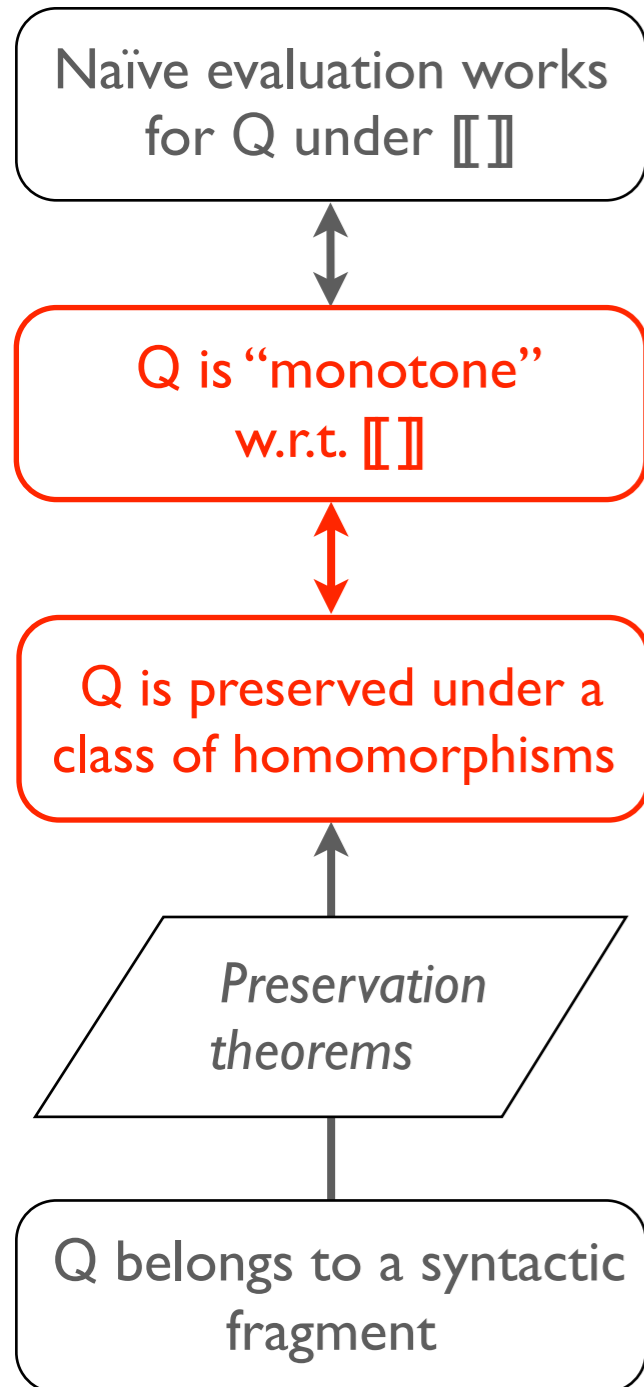
Can be generalized to arbitrary semantic relations...

Homomorphism-based relational semantics



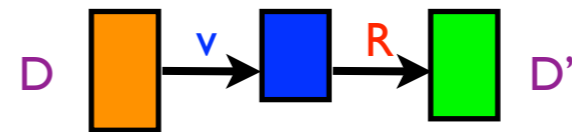
R: reflexive binary relation between complete instances

Monotonicity and preservation



- **R-homomorphism** $D \rightarrow D'$ (D and D' complete):
a mapping h over $\text{dom}(D)$ s.t. $h(D) R D'$

- R-homomorphisms "mimic" the semantic mapping:



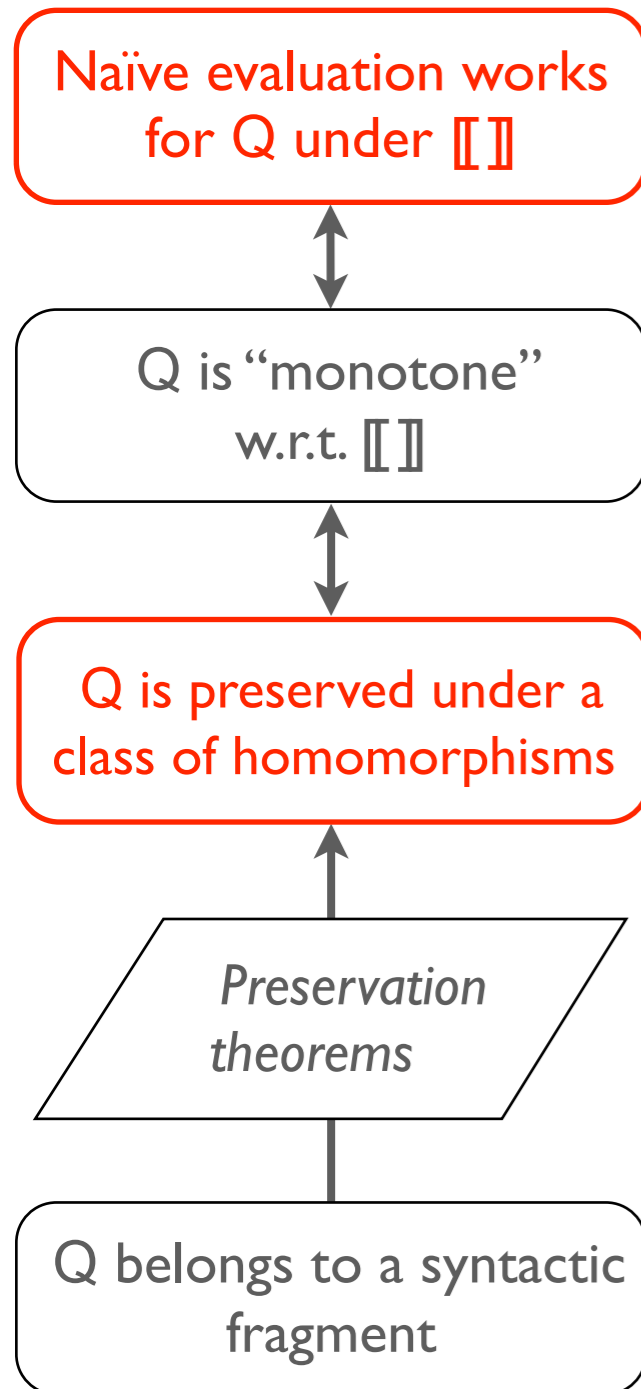
- except that valuations distinguish constants from nulls
- However, using query genericity:

If a relational semantics $\llbracket \cdot \rrbracket$ is given by R
and Q is a generic Boolean query

Q is monotone w.r.t. $\llbracket \cdot \rrbracket$ iff

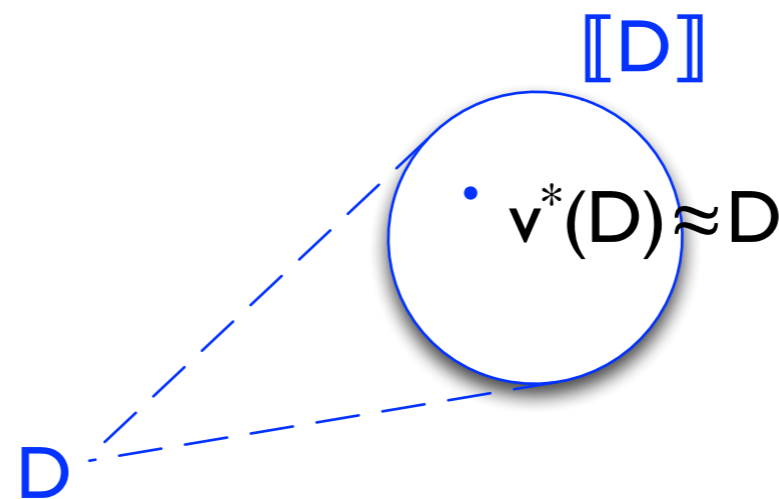
Q is preserved under R-homomorphisms

Naïve evaluation and preservation



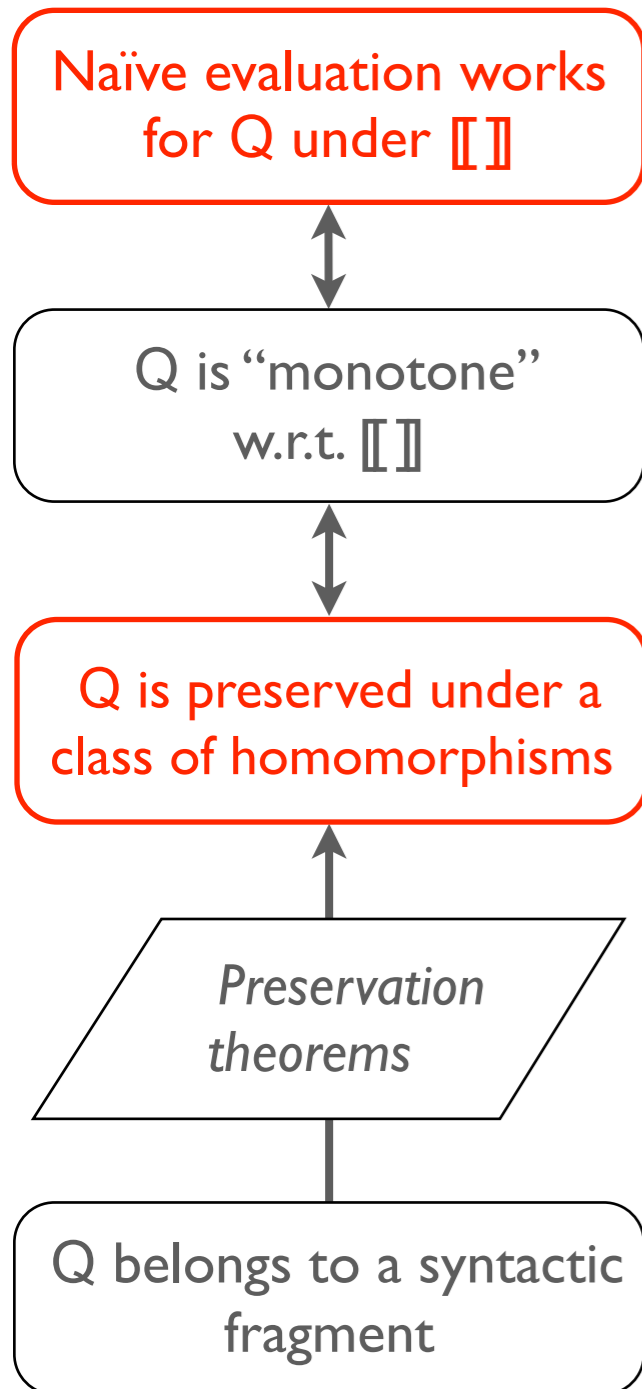
Combining the two steps:

If a relational semantics is given by R , the saturation property holds:



v^* : distinct nulls to distinct constants not occurring in D

Naïve evaluation and preservation



Theorem
 If a relational semantics is given by R and Q is a generic Boolean query
 Naive evaluation works for Q iff
 Q is preserved under R-homomorphisms

	R	R-homomorphism
OWA	\subseteq	homomorphisms
CWA	=	strong onto homomorphisms (i.e. homomorphisms $D \rightarrow h(D)$)
WCWA	\subseteq =domain	onto homomorphisms

Preservation and syntactic fragments of FO

- What about **strong onto homomorphisms** ?
 - ▶ There is a preservation result in the infinite [Keisler '65]
 - ▶ complex syntactic restrictions, one binary relation only, problematic to extend...
- A new sufficient condition for preservation, with a good syntax:

Positive fragment with Universal Guards (Pos+ $\forall G$)

$\varphi := \top \mid \perp \mid R(\bar{x}) \mid x = y \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \exists x \varphi \mid \forall x \varphi \mid$

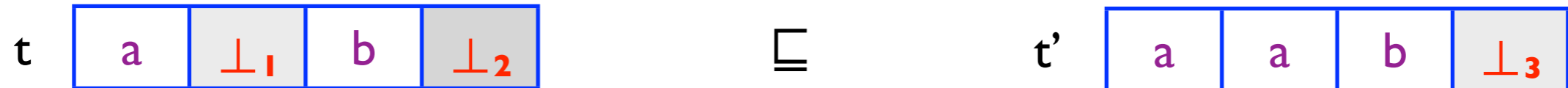
$\left(\forall \bar{x} (G(\bar{x}) \rightarrow \varphi) \right)$ with
G : a relation or equality symbol
 \bar{x} : a tuple of distinct variables

Pos+ $\forall G$ formulas are preserved under strong onto homomorphisms

Semantics arising from orderings

- Information ordering of Codd tuples :

$t \sqsubseteq t'$ if $t[i]$ constant $\Rightarrow t'[i] = t[i]$ t' is “more informative” than t



- Lifting to sets of tuples [Hoare, Plotkin 70s]

- $D \sqsubseteq^H D' : \forall t \in D \exists t' \in D' \quad t \sqsubseteq t'$ (Hoare ordering, open-world)
- $D \sqsubseteq^P D' : D \sqsubseteq^H D'$ and $\forall t' \in D' \exists t \in D \quad t \sqsubseteq t'$ (Plotkin ordering)

Semantics arising from orderings

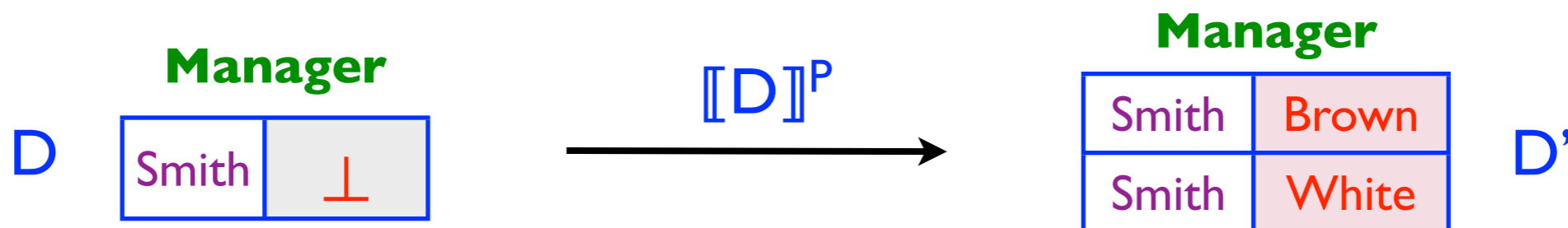
$D \sqsubseteq^H D' : \forall t \in D \ \exists t' \in D' \ t \sqsubseteq t'$ (Hoare ordering, open-world)

$D \sqsubseteq^P D' : D \sqsubseteq^H D' \text{ and } \forall t \in D' \ \exists t \in D \ t \sqsubseteq t'$ (Plotkin ordering)

- Orderings give rise to **semantics** of Codd databases:

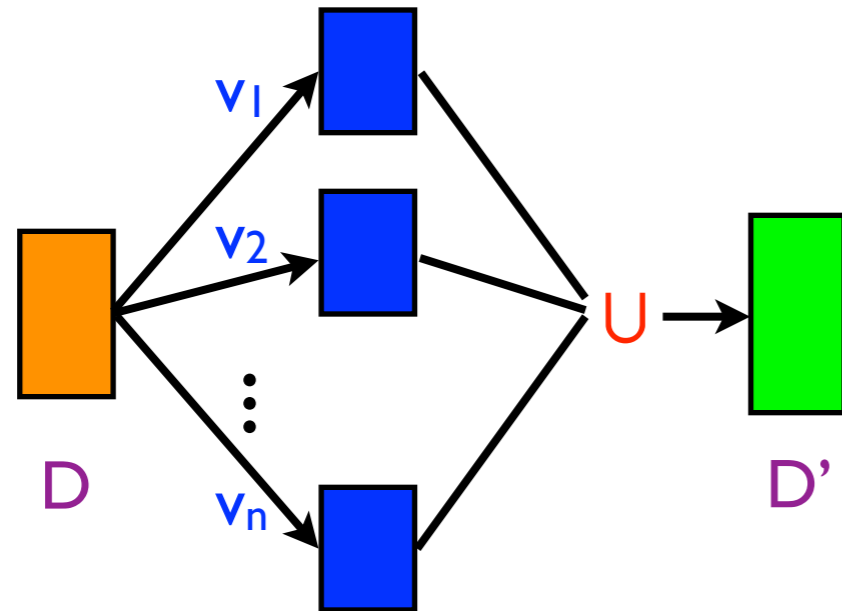
$$\llbracket D \rrbracket = \{ \text{complete } D' \mid D \sqsubseteq D' \}$$

- Observe that Plotkin semantics is more “open” than CWA :



Powerset semantics

- Extend (and generalize) Plotkin semantics to naïve databases



Powerset CWA

$D' \in (D)_{CWA}$ iff

\exists valuations $v_1, \dots, v_n \quad D' = \bigcup_i v_i(D)$

- When restricting to Codd databases Powerset CWA coincides with Plotkin
- Gives rise to a whole new class of semantics: U is replaced by any suitable relation

$$R \subseteq 2^C \times C$$

Naïve evaluation for powerset semantics

- Naïve evaluation \leftrightarrow Monotonicity \leftrightarrow Preservation continues to hold
- Under the powerset CWA the needed notion is preservation under *unions of strong onto homomorphisms* (i.e. homomorphisms $D \rightarrow \bigcup_{i=1}^n h_i(D)$)
- We have similar results for powerset semantics based on arbitrary **R**
- An FO fragment preserved under unions of strong onto homomorphisms:

$\exists\text{Pos} + \forall G^{\text{bool}}$:

$\exists\text{Pos}$ extended with *universal guarded sentences* $\forall \bar{x} (G(\bar{x}) \rightarrow \varphi(\bar{x}))$

Corollary

If Q is a Boolean query from $\exists\text{Pos} + \forall G^{\text{bool}}$

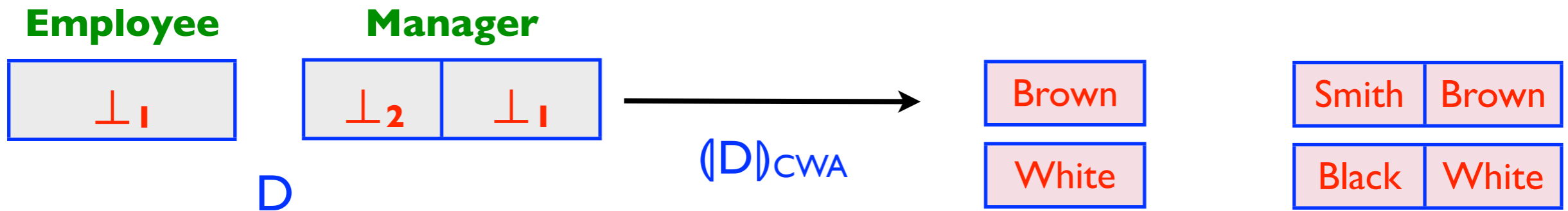
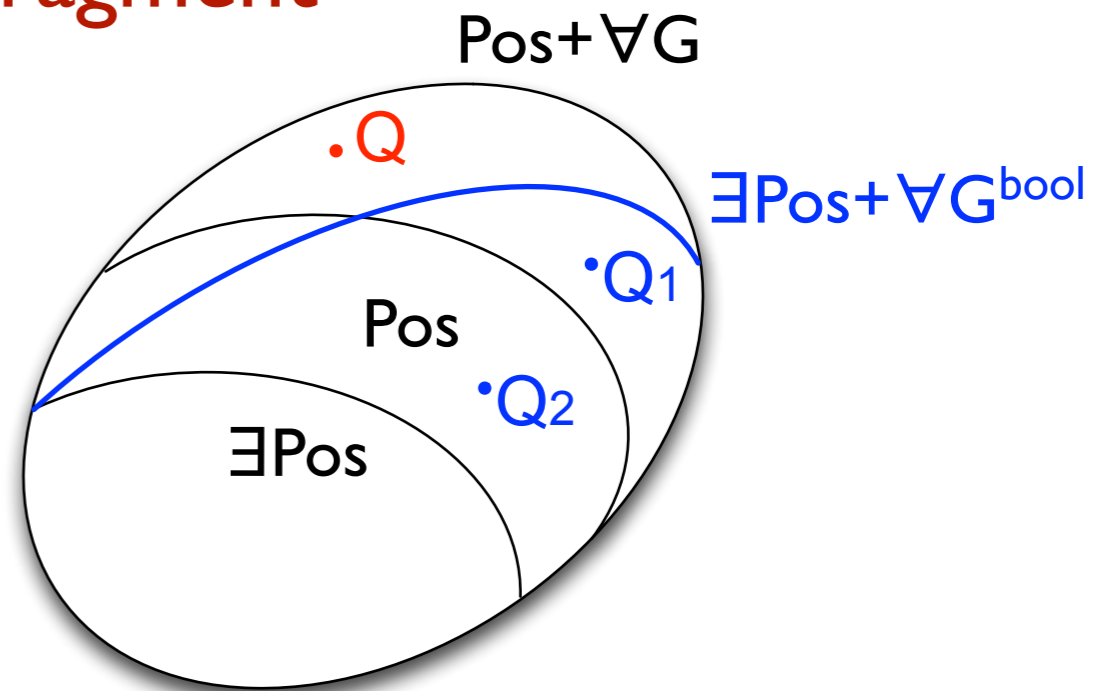
Naïve evaluation works for Q under $(\cdot)_{\text{CWA}}$

The $\exists\text{Pos}+\forall\text{G}^{\text{bool}}$ fragment

$$Q_1 = \forall x(\text{Employee}(x) \rightarrow \exists y \text{Manager}(x, y))$$

$$Q_2 = \forall x(\text{Employee}(x))$$

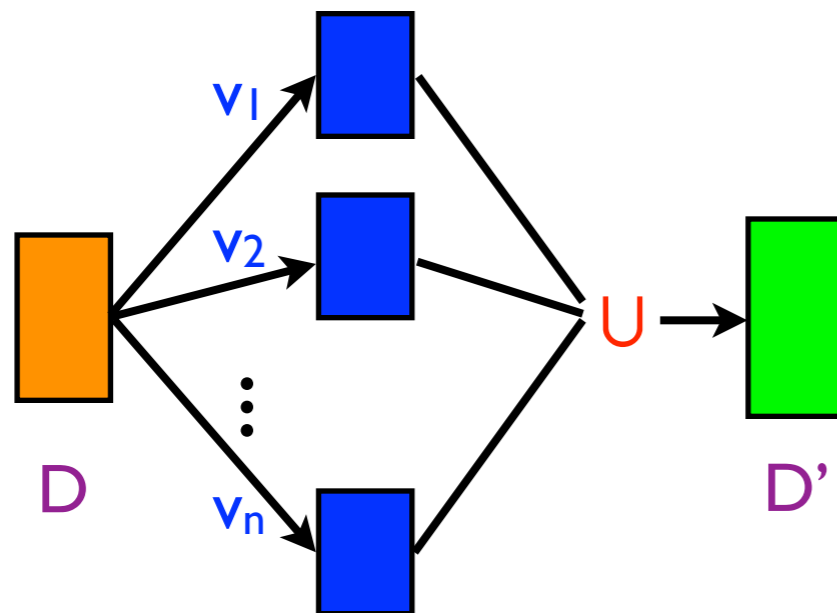
$$Q = \exists x\forall y(\text{Employee}(y) \rightarrow \text{Manager}(x, y))$$



- Naïve evaluation works for Q_1 and Q_2 under $(\cdot)_{\text{CWA}} \iff Q_1, Q_2 \in \exists\text{Pos}+\forall\text{G}^{\text{bool}}$
- Naïve evaluation does not work for Q under the $(\cdot)_{\text{CWA}} \implies Q \notin \exists\text{Pos}+\forall\text{G}^{\text{bool}}$

Minimal semantics

- A special form of powerset semantics was introduced in the field of deductive databases (GCWA [Minker '82])
- Later modified and adopted as *data exchange* semantics (GCWA* [Hernich'11])
- We define it here for arbitrary incomplete instances:



Minimal Powerset CWA

$D' \in (D)_{CWA}^{\min}$ iff

\exists *D-minimal* valuations v_1, \dots, v_n

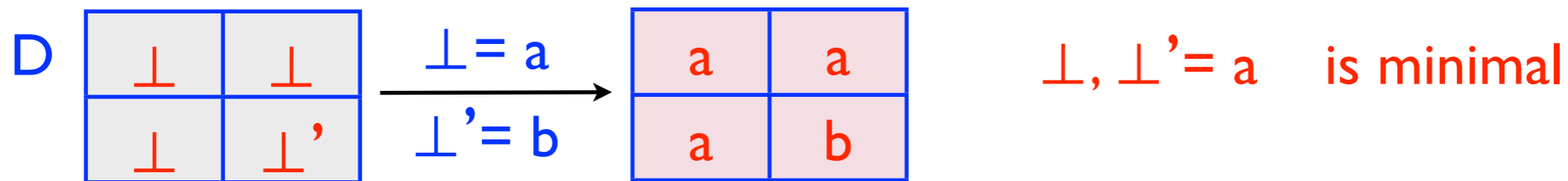
$D' = \bigcup_i v_i(D)$

A valuation v on D is *D-minimal* if there is no valuation v' s.t. $v'(D) \subsetneq v(D)$

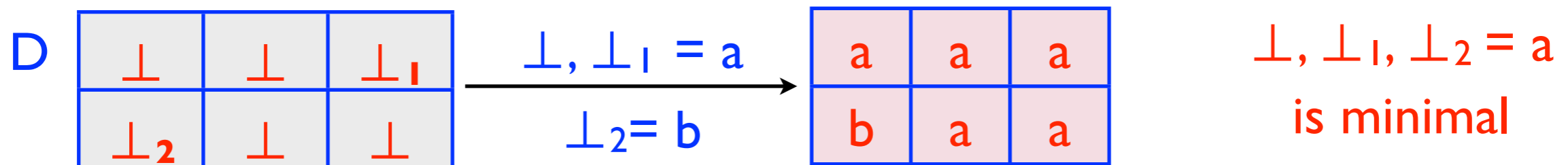
Minimal powerset semantics: U is replaced by any suitable relation $R \subseteq 2^C \times C$

Minimal semantics and the core

- Not all valuations are minimal:



- true also if D is a core



- but if v is a minimal valuation $v(D) = v(\text{core}(D))$

- There are other important connections between minimal semantics and the core (later)

Core of D

substructure D' of D such that $D \rightarrow D'$
 and there is no $D'' \subsetneq D'$ s.t. $D \rightarrow D''$
 (\rightarrow : homomorphism preserving constants)

Minimal semantics and the saturation property

Naive evaluation works
for Q under $\llbracket \cdot \rrbracket$



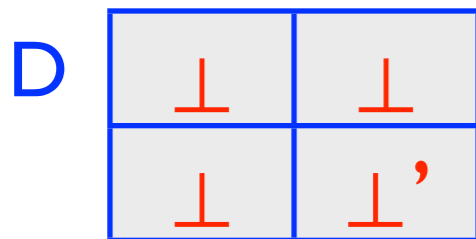
Q is “monotone”
w.r.t. $\llbracket \cdot \rrbracket$

saturated
domain

Saturation property for $\langle \mathcal{D}, C, \llbracket \cdot \rrbracket, \approx \rangle$:

For all $x \in \mathcal{D}$ there exists $y \in \llbracket x \rrbracket$ $y \approx x$

Under the minimal Powerset CWA the saturation property does not hold



- all D -minimal images are of the form

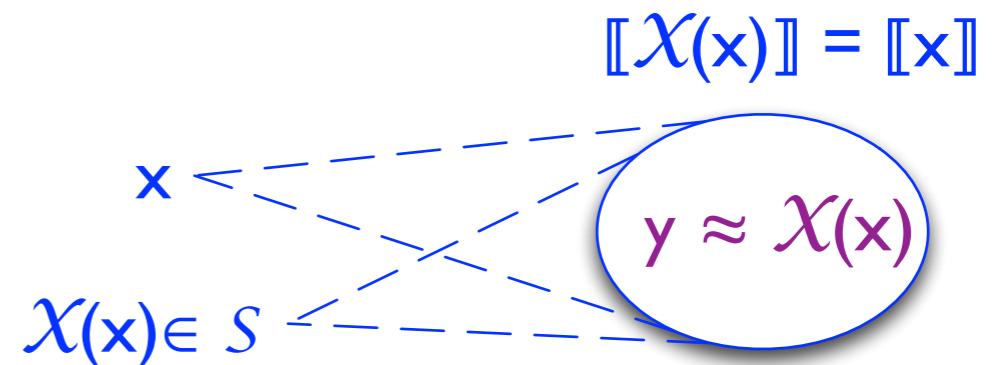
a	a
-----	-----
- No union of D -minimal images can be isomorphic to D

The saturation property revisited

$\langle \mathcal{D}, C, \llbracket \cdot \rrbracket, \approx \rangle$ has a **saturated subdomain** if $\exists S$ with $C \subseteq S \subseteq \mathcal{D}$

and a function $\mathcal{X}: \mathcal{D} \rightarrow S$ (the *representative function*) s.t.

- $\langle S, C, \llbracket \cdot \rrbracket, \approx \rangle$ is saturated
- $\llbracket \mathcal{X}(x) \rrbracket = \llbracket x \rrbracket$ for all $x \in \mathcal{D}$



Proposition

If a database domain has a saturated subdomain with representative function \mathcal{X} and Q is a generic Boolean query

Naïve evaluation works for Q iff

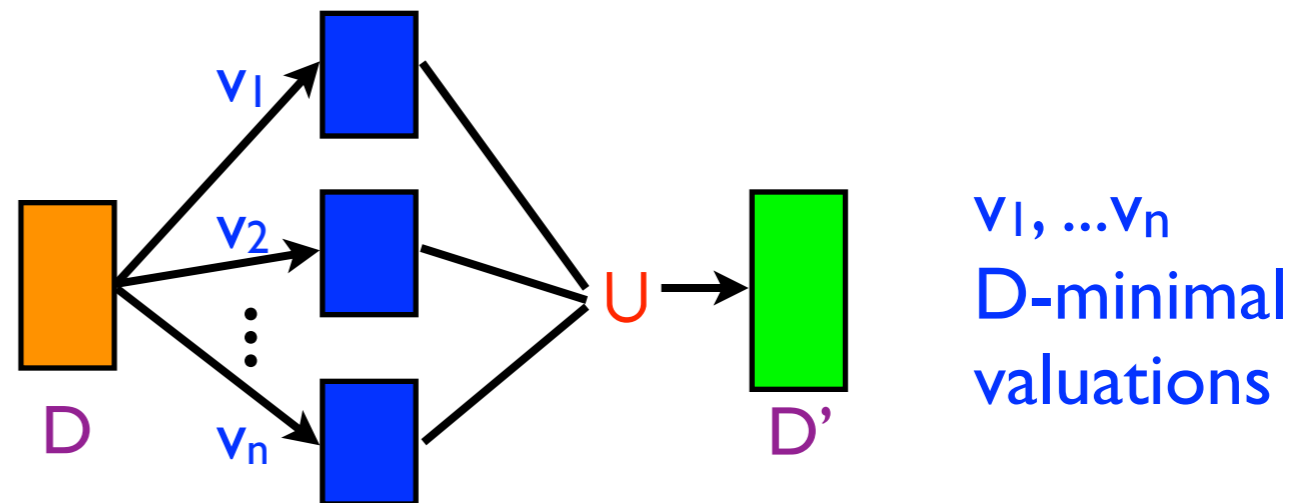
Q is **monotone w.r.t. $\llbracket \cdot \rrbracket$** and $Q(x) = Q(\mathcal{X}(x))$ for all x

Lemma: Under the minimal powerset CWA the set of **cores** is a saturated subdomain with representative function $\text{core}(\cdot)$

Monotonicity and preservation for minimal semantics

- Monotonicity under minimal powerset CWA is

preservation under the mapping



- Query genericity used with care:

- ▶ valuations are indistinguishable from homomorphisms, however
- ▶ v_1, \dots, v_n are minimal w.r.t all other valuations (not all arbitrary homomorphisms)
- ▶ v_1, \dots, v_n preserve the same elements of D

Naive evaluation and preservation for minimal semantics

The right notion of preservation:

- ▶ **D-minimal homomorphism h :**
there exists no h' , preserving all constants preserved by h , s.t. $h'(D) \subsetneq h(D)$
- ▶ **Unions of minimal homomorphisms:** homomorphisms $D \rightarrow \bigcup_{i=1}^n h_i(D)$
with $h_1 \dots h_n$ D-minimal and preserving the same constants

Theorem

If Q is a generic Boolean query

Naive evaluation works for Q under the minimal powerset CWA iff

Q is **preserved under unions of minimal homomorphisms** and

$Q(D) = Q(\text{core}(D))$ for every database D

Similar results hold for arbitrary minimal semantics

Preservation and syntactic fragments for minimal semantics

- **Preservation under unions of minimal homomorphisms** : no “tight” syntactic fragment known
- *Remark*: unions of minimal homomorphisms are also unions of strong onto homomorphisms

If Q is a Boolean query from $\exists\text{Pos}+\forall\text{G}^{\text{bool}}$, under the minimal powerset CWA:

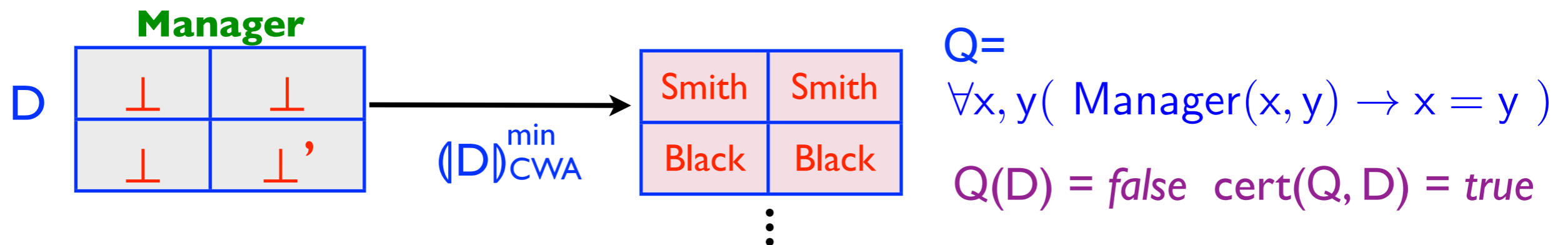
- Naïve evaluation works for Q iff $Q(D) = Q(\text{core}(D))$ for every D
- Naïve evaluation works over cores

Preservation and syntactic fragments for minimal semantics

- **Preservation under unions of minimal homomorphisms** : no “tight” syntactic fragment known
- **Remark**: unions of minimal homomorphisms are also unions of strong onto homomorphisms

If Q is a Boolean query from $\exists\text{Pos}+\forall\text{G}^{\text{bool}}$, under the minimal powerset CWA:

- Naïve evaluation works for Q iff $Q(D) = Q(\text{core}(D))$ for every D
- Naïve evaluation works over cores



$Q \in \exists\text{Pos}+\forall\text{G}^{\text{bool}}$ but naïve evaluation does not work

$Q(D) \neq Q(\text{core}(D))$

Non-Boolean queries

All results can be lifted to non-boolean relational queries. For a k -ary query Q :

- ▶ Define a new database domain whose elements are pairs (D, t)
 D : a relational database t : a k -tuple of constants
- ▶ Define a boolean query Q^* s.t. $Q^*(D, t) = \text{true}$ iff $t \in Q(D)$
- ▶ Apply previous results to Q^* and the new database domain \Rightarrow derive results for Q over the original relational database domain

For k -ary FO queries, $k \geq 0$

Semantics

Naïve evaluation works for

OWA

\exists Pos

WCWA

Pos

CWA

Pos + $\forall G$

Powerset CWA

\exists Pos + $\forall G^{\text{bool}}$

Min Powerset CWA

\exists Pos + $\forall G^{\text{bool}}$ iff $Q(D) = Q(\text{core}(D))$