

# Datalog and Recursive Queries

Universidad de Concepción, 2014

(Slides adapted from Loreto Bravo who adapted from Werner Nutt who adapted them from Thomas Eiter and Leonid Libkin)

Bases de Datos II

1

## Motivation

- Relational Calculus and Relational Algebra were considered to be “*the*” database languages for a long time
- Codd: A query language is “complete,” if it yields Relational Calculus
- However, Relational Calculus misses an important feature: *recursion*
- Example: A metro database with relation `links:line, station, nextstation`  
What stations are reachable from station “Odeon”?  
Can we go from Odeon to Tuileries?  
etc.
- It can be proved: such queries cannot be expressed in Relational Calculus
- This motivated a logic-programming extension to conjunctive queries: *datalog*

Datalog

Example: Metro Database Instance

links	line	station	nextstation
	4	St.Germain	Odeon
	4	Odeon	St.Michel
	4	St. Michel	Chatelet
	1	Chatelet	Louvres
	1	Louvres	Palais Royal
	1	Palais-Royal	Tuileries
	1	Tuileries	Concorde

Datalog program for first query:

```
reach(X,X) ← links(L,X,Y)
reach(X,X) ← links(L,Y,X)
reach(X,Y) ← links(L,X,Z), reach(Z,Y)
answer(X) ← reach('Odeon',X)
```

Note: recursive definition

Intuitively, if the part right of “←” is true, the rule “fires” and the atom left of “←” is concluded.

Datalog

Example: Ancestor-Descendant

ParentChild	parent	child
	Homer	Bart
	Homer	Lisa
	Marge	Bart
	Marge	Lisa
	Abe	Homer
	Ape	Abe

Query: Who are Bart’s ancestors?

Datalog program:

```
AncestorDescendant(X,Y) ← ParentChild(X,Y)
AncestorDescendant(X,Y) ← ParentChild(X,Z), AncestorDescendant(Z,Y)
answer(X) ← AncestorDescendant(X,'Bart')
```

Datalog

Datalog Syntax

**Definition.** A *datalog rule*  $r$  is an expression of the form

$$R_0(\vec{x}_0) \leftarrow R_1(\vec{x}_1), \dots, R_n(\vec{x}_n)$$

(1)

- where  $n \geq 0$ ,  
 $R_0, \dots, R_n$  are relations names, and  
 $\vec{x}_0, \dots, \vec{x}_n$  are vectors of variables and constants (from **dom**)
- every variable in  $\vec{x}_0$  occurs in  $\vec{x}_1, \dots, \vec{x}_n$  (“safety”)

**Remarks.**

- The *head* of  $r$ , denoted  $H(r)$ , is  $R_0(\vec{x}_0)$
- The *body* of  $r$ , denoted  $B(r)$ , is  $\{ R_1(\vec{x}_1), \dots, R_n(\vec{x}_n) \}$
- The *rule* symbol “ $\leftarrow$ ” is often also written as “ $: -$ ”

**Definition.** A *datalog program* is a finite set of datalog rules.

Datalog Programs

Let  $P$  be a datalog program.

- An *extensional relation* of  $P$  is a relation occurring only in rule bodies of  $P$
- An *intensional relation* of  $P$  is a relation occurring in the head of some rule in  $P$
- The *extensional schema* of  $P$ ,  $edb(P)$ , consists of all extensional relations of  $P$
- The *intensional schema* of  $P$ ,  $idb(P)$ , consists of all intensional relations of  $P$
- The *schema* of  $P$ ,  $sch(P)$ , is the union of  $edb(P)$  and  $idb(P)$ .

**The Metro Example /1**

Datalog program  $P$  on metro database scheme

$\mathcal{M} = \{\text{links} : \text{line}, \text{station}, \text{nextstation}\}$ :

```

reach(X, X)  ←  links(L, X, Y)
reach(X, X)  ←  links(L, Y, X)
reach(X, Y)  ←  links(L, X, Z), reach(Z, Y)
answer(X)   ←  reach('Odeon', X)

```

Here,

$$\begin{aligned}
 edb(P) &= \{\text{links}\} \quad (= \mathcal{M}), \\
 idb(P) &= \{\text{reach}, \text{answer}\}, \\
 sch(P) &= \{\text{links}, \text{reach}, \text{answer}\}
 \end{aligned}$$

Bases de Datos II

7

**Datalog Syntax (cntd)**

Datalog

- The set of constants occurring in a datalog program  $P$  is denoted as  $adom(P)$
- Given a database instance  $\mathbf{I}$ , we define the *active domain* of  $P$  with respect to  $I$  as

$$adom(P, \mathbf{I}) := adom(P) \cup adom(\mathbf{I}),$$

that is, as the set of constants occurring in  $P$  and  $\mathbf{I}$

**Definition.** Let  $\nu: \text{var}(r) \cup \mathbf{dom} \rightarrow \mathbf{dom}$  be a valuation for a rule  $r$  of form (1).

Then the *instantiation* of  $r$  with  $\nu$ , denoted  $\nu(r)$ , is the rule

$$R_0(\nu(\vec{x}_0)) \leftarrow R_1(\nu(\vec{x}_1)), \dots, R_n(\nu(\vec{x}_n))$$

which results from replacing each variable  $x$  with  $\nu(x)$ .

Datalog

The Metro Example /2

- For the datalog program  $P$  above, we have that  $adom(P) = \{ \text{Odeon} \}$
- We consider the database instance  $\mathbf{I}$ :

links	line	station	nextstation
	4	St.Germain	Odeon
	4	Odeon	St.Michel
	4	St. Michel	Chatelet
	1	Chatelet	Louvres
	1	Louvres	Palais-Royal
	1	Palais-Royal	Tuileries
	1	Tuileries	Concorde

Then  $adom(\mathbf{I}) = \{4, 1, \text{St.Germain}, \text{Odeon}, \text{St.Michel}, \text{Chatelet}, \text{Louvres}, \text{Palais-Royal}, \text{Tuileries}, \text{Concorde}\}$

- Also  $adom(P, \mathbf{I}) = adom(\mathbf{I})$ .

The Metro Example /3

- The rule

$reach(\text{St.Germain}, \text{Odeon}) \leftarrow links(\text{Louvres}, \text{St.Germain}, \text{Concorde}), reach(\text{Concorde}, \text{Odeon})$

is an instance of the rule

$reach(X, Y) \leftarrow links(L, X, Z), reach(Z, Y)$

of  $P$ :

take  $\nu(X) = \text{St.Germain}, \nu(L) = \text{Louvres}, \nu(Y) = \text{Odeon}, \nu(Z) = \text{Concorde}$

## Datalog Semantics

- There exist several approaches to defining the semantics:

### Operational (fixpoint) approach:

Obtain query result by applying an inference procedure,  
until a fixpoint is reached

### Model-theoretic approach:

View rules as logical sentences, which state the query result

### Proof-theoretic approach:

Obtain proofs of facts in the query result, following a proof calculus  
(based on resolution)

## Fixpoint Semantics

Basic idea:

“If all facts in  $\mathbf{I}$  hold, which other facts must hold after firing the rules in  $P$ ?”

Approach:

- Define an *immediate consequence operator*  $\mathbf{T}_P(\mathbf{K})$  on db instances  $\mathbf{K}$ .
- Start with  $\mathbf{K} = \mathbf{I}$ .
- Apply  $\mathbf{T}_P$  to obtain a new instance:  $\mathbf{K}_{new} := \mathbf{T}_P(\mathbf{K}) = \mathbf{I} \cup \text{new facts}$ .
- Iterate until nothing new can be produced.
- The result yields the semantics.

**Immediate Consequence Operator**

Let  $P$  be a datalog program and  $\mathbf{K}$  be a database instance of  $sch(P)$ .

A fact  $R(\vec{t})$  is an *immediate* consequence for  $\mathbf{K}$  and  $P$ , if either

- $R \in edb(P)$  and  $R(\vec{t}) \in \mathbf{K}$ , or
- there exists a ground instance  $r$  of a rule in  $P$  such that  
 $H(r) = R(\vec{t})$  and  $B(r) \subseteq \mathbf{K}$ .

**Definition.** The *immediate consequence operator* of a datalog program  $P$  is the mapping

$$\mathbf{T}_P: inst(sch(P)) \rightarrow inst(sch(P))$$

where

Bases de Datos II

13

$$\mathbf{T}_P(\mathbf{K}) = \{ A \mid A \text{ is an immediate consequence for } \mathbf{K} \text{ and } P \}.$$

Datalog

**Example**

Consider

$$P = \{ \text{reachable}(\mathbf{a}) \\ \text{reachable}(\mathbf{Y}) \leftarrow \text{arc}(\mathbf{X}, \mathbf{Y}), \text{reachable}(\mathbf{X}) \}$$

where  $edb(P) = \{\text{arc}\}$  and  $idb(P) = \{\text{reachable}\}$ .

$$\begin{aligned} \mathbf{I} = \mathbf{K}_1 &= \{ \text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}) \} \\ \mathbf{K}_2 &= \{ \text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}), \text{reachable}(\mathbf{a}) \} \\ \mathbf{K}_3 &= \{ \text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}), \text{reachable}(\mathbf{a}), \text{reachable}(\mathbf{b}) \} \\ \mathbf{K}_4 &= \{ \text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}), \text{reachable}(\mathbf{a}), \text{reachable}(\mathbf{b}), \text{reachable}(\mathbf{c}) \} \end{aligned}$$

Datalog

**Example (cntd)**

Then,

$$\begin{aligned}
 \mathbf{T}_P(\mathbf{K}_1) &= \{\text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}), \text{reachable}(\mathbf{a})\} = \mathbf{K}_2 \\
 \mathbf{T}_P(\mathbf{K}_2) &= \{\text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}), \text{reachable}(\mathbf{a}), \text{reachable}(\mathbf{b})\} = \mathbf{K}_3 \\
 \mathbf{T}_P(\mathbf{K}_3) &= \{\text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}), \text{reachable}(\mathbf{a}), \text{reachable}(\mathbf{b}), \text{reachable}(\mathbf{c})\} = \mathbf{K}_4 \\
 \mathbf{T}_P(\mathbf{K}_4) &= \{\text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}), \text{reachable}(\mathbf{a}), \text{reachable}(\mathbf{b}), \text{reachable}(\mathbf{c})\} = \mathbf{K}_4
 \end{aligned}$$

Thus,  $\mathbf{K}_4$  is a *fixpoint* of  $\mathbf{T}_P$ .

**Definition.**  $\mathbf{K}$  is a *fixpoint* of operator  $\mathbf{T}_P$  if  $\mathbf{T}_P(\mathbf{K}) = \mathbf{K}$ .

**Properties**

**Proposition.** For every datalog program  $P$  we have:

1. The operator  $\mathbf{T}_P$  is monotonic, that is,  $\mathbf{K} \subseteq \mathbf{K}'$  implies  $\mathbf{T}_P(\mathbf{K}) \subseteq \mathbf{T}_P(\mathbf{K}')$ ;
2. For any  $\mathbf{K} \in \text{inst}(\text{sch}(P))$  we have:

$\mathbf{K}$  is a model of  $\Sigma_P$  if and only if  $\mathbf{T}_P(\mathbf{K}) \subseteq \mathbf{K}$ ;

3. If  $\mathbf{T}_P(\mathbf{K}) = \mathbf{K}$  (i.e.,  $\mathbf{K}$  is a fixpoint), then  $\mathbf{K}$  is a model of  $\Sigma_P$ .

Note: The converse of 3. does not hold in general.



## Fixpoint Iteration

For a datalog program  $P$  and database instance  $\mathbf{I}$ , define the sequence  $(\mathbf{I}_i)_{i \geq 0}$  by

$$\begin{aligned}\mathbf{I}_0 &= \mathbf{I} \\ \mathbf{I}_i &= \mathbf{T}_P(\mathbf{I}_{i-1}) \quad \text{for } i > 0.\end{aligned}$$

- By monotonicity of  $\mathbf{T}_P$ , we have  $\mathbf{I}_0 \subseteq \mathbf{I}_1 \subseteq \mathbf{I}_2 \subseteq \dots \subseteq \mathbf{I}_i \subseteq \mathbf{I}_{i+1} \subseteq \dots$
- For every  $i \geq 0$ , we have  $\mathbf{I}_i \subseteq \mathbf{B}(P, \mathbf{I})$ , with  $\mathbf{B}(P, \mathbf{I})$  the instance of  $sch(P)$
- Hence, for some integer  $n \leq |\mathbf{B}(P, \mathbf{I})|$ , we have  $\mathbf{I}_{n+1} = \mathbf{I}_n$  ( $=: \mathbf{T}_P^\omega(\mathbf{I})$ )
- It holds that  $\mathbf{T}_P^\omega(\mathbf{I}) = lfp(P, \mathbf{I}) = P(\mathbf{I})$ .

Bases de Datos II

17

This can be readily implemented by an algorithm.

Datalog

## Example

$$\begin{aligned}P &= \{ \text{reachable}(\mathbf{a}) \\ &\quad \text{reachable}(\mathbf{Y}) \leftarrow \text{arc}(\mathbf{X}, \mathbf{Y}), \text{reachable}(\mathbf{X}) \} \\ \mathbf{I} &= \{ \text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}) \}\end{aligned}$$

Then,

$$\begin{aligned}\mathbf{I}_0 &= \{ \text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}) \} \\ \mathbf{I}_1 = \mathbf{T}_P^1(\mathbf{I}) &= \{ \text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}), \text{reachable}(\mathbf{a}) \} \\ \mathbf{I}_2 = \mathbf{T}_P^2(\mathbf{I}) &= \{ \text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}), \text{reachable}(\mathbf{a}), \text{reachable}(\mathbf{b}) \} \\ \mathbf{I}_3 = \mathbf{T}_P^3(\mathbf{I}) &= \{ \text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}), \text{reachable}(\mathbf{a}), \text{reachable}(\mathbf{b}), \text{reachable}(\mathbf{c}) \} \\ \mathbf{I}_4 = \mathbf{T}_P^4(\mathbf{I}) &= \{ \text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}), \text{reachable}(\mathbf{a}), \text{reachable}(\mathbf{b}), \text{reachable}(\mathbf{c}) \} \\ &= \mathbf{T}_P^3(\mathbf{I})\end{aligned}$$

Thus,  $\mathbf{T}_P^\omega(\mathbf{I}) = lfp(P, \mathbf{I}) = \mathbf{I}_4$ .

Datalog

**Datalog: Model-Theoretic Semantics**
**General Idea:**

- We view a program as a set of first-order sentences
- Given an instance  $\mathbf{I}$  of  $edb(P)$ , the result of  $P$  is a database instance of  $sch(P)$  that extends  $\mathbf{I}$  and satisfies the sentences (or, is a *model* of the sentences)
- There can be many models
- The *intended answer* is specified by particular models
- These particular models are selected by “external” conditions

Bases de Datos II

19

Datalog

**Logical Theory  $\Sigma_P$** 

- To every datalog rule  $r$  of the form  $R_0(\vec{x}_0) \leftarrow R_1(\vec{x}_1), \dots, R_n(\vec{x}_n)$ , with variables  $x_1, \dots, x_m$ , we associate the logical sentence  $\sigma(r)$ :

$$\forall x_1, \dots, \forall x_m (R_1(\vec{x}_1) \wedge \dots \wedge R_n(\vec{x}_n) \rightarrow R_0(\vec{x}_0))$$

- To a program  $P$ , we associate the set of sentences  $\Sigma_P = \{\sigma(r) \mid r \in P\}$ .

**Definition.** Let  $P$  be a datalog program and  $\mathbf{I}$  an instance of  $edb(P)$ . Then,

- A *model* of  $P$  is an instance of  $sch(P)$  that satisfies  $\Sigma_P$
- The *semantics* of  $P$  on input  $\mathbf{I}$ , denoted  $P(\mathbf{I})$ , is the *least model* of  $P$  containing  $\mathbf{I}$ , if it exists.

Datalog

Example

For program  $P$  and instance  $\mathbf{I}$  of the Metro Example, the least model is:

links	line	station	nextstation	reach		
	4	St.Germain	Odeon		St.Germain	St.Germain
	4	Odeon	St.Michel		Odeon	Odeon
	4	St. Michel	Chatelet		...	
	1	Chatelet	Louvres		Concorde	Concorde
	1	Louvres	Palais-Royal		St.Germain	Odeon
	1	Palais-Royal	Tuileries		St.Germain	St.Michel
	1	Tuileries	Concorde		St.Germain	Chatelet
					St.Germain	Louvres
					...	

answer	
	Odeon
	St.Michel
	Chatelet
	Louvres
	Palais-Royal
	Tuileries
	Concorde

Questions

Datalog

- Is the semantics  $P(\mathbf{I})$  well-defined for every input instance  $\mathbf{I}$ ?
- How can one compute  $P(\mathbf{I})$ ?

Observation: For any  $\mathbf{I}$ , there is a model of  $P$  containing  $\mathbf{I}$

- Let  $\mathbf{B}(P, \mathbf{I})$  be the instance of  $sch(P)$  such that

$$\mathbf{B}(P, \mathbf{I})(R) = \begin{cases} \mathbf{I}(R) & \text{for each } R \in edb(P) \\ adom(P, \mathbf{I})^{arity(R)} & \text{for each } R \in idb(P) \end{cases}$$

- Then:  $\mathbf{B}(P, \mathbf{I})$  is a model of  $P$  containing  $\mathbf{I}$   
 $\Rightarrow P(\mathbf{I})$  is a subset of  $\mathbf{B}(P, \mathbf{I})$  (if it exists)
- Naive algorithm: explore all subsets of  $\mathbf{B}(P, \mathbf{I})$

Datalog

### Elementary Properties of $P(\mathbf{I})$

Let  $P$  be a datalog program,  $\mathbf{I}$  an instance of  $edb(P)$ , and  $\mathcal{M}(\mathbf{I})$  the set of all models of  $P$  containing  $\mathbf{I}$ .

**Theorem.** The intersection  $\bigcap_{M \in \mathcal{M}(\mathbf{I})} M$  is a model of  $P$ .

**Corollary.**

1.  $P(\mathbf{I}) = \bigcap_{M \in \mathcal{M}(\mathbf{I})} M$
2.  $adom(P(\mathbf{I})) \subseteq adom(P, \mathbf{I})$ , that is, no new values appear
3.  $P(\mathbf{I})(R) = \mathbf{I}(R)$ , for each  $R \in edb(P)$ .

**Consequences:**

- $P(\mathbf{I})$  is well-defined for every  $\mathbf{I}$
- If  $P$  and  $\mathbf{I}$  are finite, the  $P(\mathbf{I})$  is finite

Bases de Datos II

23

Datalog

### Why Choose the Least Model?

There are two reasons to choose the least model containing  $\mathbf{I}$ :

1. The *Closed World Assumption*:
  - If a fact  $R(\vec{c})$  is not true in all models of a database  $\mathbf{I}$ , then infer that  $R(\vec{c})$  is false
  - This amounts to considering  $\mathbf{I}$  as complete
  - ... which is customary in database practice
2. The relationship to Logic Programming:
  - Datalog should desirably match Logic Programming (seamless integration)
  - Logic Programming builds on the minimal model semantics

Datalog

### Relating Datalog to Logic Programming

- A logic program makes no distinction between *edb* and *idb*
- A datalog program  $P$  and an instance  $\mathbf{I}$  of  $edb(P)$  can be mapped to the logic program

$$\mathcal{P}(P, \mathbf{I}) = P \cup \mathbf{I}$$

(where  $\mathbf{I}$  is viewed as a set of atoms in the Logic Programming perspective)

- Correspondingly, we define the logical theory

$$\Sigma_{P, \mathbf{I}} = \Sigma_P \cup \mathbf{I}$$

- The semantics of the logic program  $\mathcal{P} = \mathcal{P}(P, \mathbf{I})$  is defined in terms of *Herbrand interpretations* of the language induced by  $\mathcal{P}$ :

The domain of discourse is formed by the constants occurring in  $\mathcal{P}$

- Each constant occurring in  $\mathcal{P}$  is interpreted by itself

Datalog

### Herbrand Interpretations of Logic Programs

Given a rule  $r$ , we denote by  $Const(r)$  the set of all constants in  $r$

**Definition.** For a (function-free) logic program  $\mathcal{P}$ , we define

- the *Herbrand universe* of  $\mathcal{P}$ , by

$$\mathbf{HU}(\mathcal{P}) = \bigcup_{r \in \mathcal{P}} Const(r)$$

- the *Herbrand base* of  $\mathcal{P}$ , by

$$\mathbf{HB}(\mathcal{P}) = \{R(c_1, \dots, c_n) \mid R \text{ is a relation in } \mathcal{P}, \\ c_1, \dots, c_n \in \mathbf{HU}(\mathcal{P}), \text{ and } \text{arity}(R) = n\}$$

Datalog

**Example**

$$\mathcal{P} = \{ \text{arc}(\text{a}, \text{b}). \\ \text{arc}(\text{b}, \text{c}). \\ \text{reachable}(\text{a}). \\ \text{reachable}(\text{Y}) \leftarrow \text{arc}(\text{X}, \text{Y}), \text{reachable}(\text{X}). \}$$

$$\mathbf{HU}(\mathcal{P}) = \{\text{a}, \text{b}, \text{c}\}$$

$$\mathbf{HB}(\mathcal{P}) = \{\text{arc}(\text{a}, \text{a}), \text{arc}(\text{a}, \text{b}), \text{arc}(\text{a}, \text{c}), \\ \text{arc}(\text{b}, \text{a}), \text{arc}(\text{b}, \text{b}), \text{arc}(\text{b}, \text{c}), \\ \text{arc}(\text{c}, \text{a}), \text{arc}(\text{c}, \text{b}), \text{arc}(\text{c}, \text{c}), \\ \text{reachable}(\text{a}), \text{reachable}(\text{b}), \text{reachable}(\text{c})\}$$

Bases de Datos II

27

Datalog

**Grounding**

- A rule  $r'$  is a *ground instance* of a rule  $r$  with respect to  $\mathbf{HU}(\mathcal{P})$ , if  $r' = \nu(r)$  for a valuation  $\nu$  such that  $\nu(x) \in \mathbf{HU}(\mathcal{P})$  for each  $x \in \text{var}(r)$ .
- The *grounding* of a rule  $r$  with respect to  $\mathbf{HU}(\mathcal{P})$ , denoted  $\text{Ground}_{\mathcal{P}}(r)$ , is the set of all ground instances of  $r$  wrt  $\mathbf{HU}(\mathcal{P})$
- The *grounding* of a logic program  $\mathcal{P}$  is

$$\text{Ground}(\mathcal{P}) = \bigcup_{r \in \mathcal{P}} \text{Ground}_{\mathcal{P}}(r)$$

Datalog

### Example

$Ground(\mathcal{P}) = \{ \text{arc}(a, b). \text{ arc}(b, c). \text{ reachable}(a). \\
\text{reachable}(a) \leftarrow \text{arc}(a, a), \text{ reachable}(a). \\
\text{reachable}(b) \leftarrow \text{arc}(a, b), \text{ reachable}(a). \\
\text{reachable}(c) \leftarrow \text{arc}(a, c), \text{ reachable}(a). \\
\text{reachable}(a) \leftarrow \text{arc}(b, a), \text{ reachable}(b). \\
\text{reachable}(b) \leftarrow \text{arc}(b, b), \text{ reachable}(b). \\
\text{reachable}(c) \leftarrow \text{arc}(b, c), \text{ reachable}(b). \\
\text{reachable}(a) \leftarrow \text{arc}(c, a), \text{ reachable}(c). \\
\text{reachable}(b) \leftarrow \text{arc}(c, b), \text{ reachable}(c). \\
\text{reachable}(c) \leftarrow \text{arc}(c, c), \text{ reachable}(c). \}$

### Herbrand Models

- A *Herbrand-interpretation*  $I$  of  $\mathcal{P}$  is any subset  $I \subseteq \mathbf{HB}(\mathcal{P})$
- A *Herbrand-model* of  $\mathcal{P}$  is a Herbrand-interpretation that satisfies all sentences in  $\Sigma_{P, \mathbf{I}}$

Equivalently,  $M \subseteq \mathbf{HB}(\mathcal{P})$  is a Herbrand model if

- for all  $r \in Ground(\mathcal{P})$  such that  $B(r) \subseteq M$  we have that  $H(r) \subseteq M$

**Example**

The Herbrand models of program  $\mathcal{P}$  above are exactly the following:

- $M_1 = \{ \text{arc}(\mathbf{a}, \mathbf{b}), \text{arc}(\mathbf{b}, \mathbf{c}), \\ \text{reachable}(\mathbf{a}), \text{reachable}(\mathbf{b}), \text{reachable}(\mathbf{c}) \}$
- $M_2 = \mathbf{HB}(\mathcal{P})$
- every interpretation  $M$  such that  $M_1 \subseteq M \subseteq M_2$

and no others.

**Logic Programming Semantics**

- **Proposition.**  $\mathbf{HB}(\mathcal{P})$  is always a model of  $\mathcal{P}$
- **Theorem.** For every logic program there exists a least Herbrand model (wrt “ $\subseteq$ ”).  
For a program  $\mathcal{P}$ , this model is denoted  $MM(\mathcal{P})$  (for “minimal model”).  
The model  $MM(\mathcal{P})$  is the semantics of  $\mathcal{P}$ .
- **Theorem (Datalog  $\leftrightarrow$  Logic Programming).** Let  $P$  be a datalog program and  $\mathbf{I}$  be an instance of  $edb(P)$ . Then,

$$P(\mathbf{I}) = MM(\mathcal{P}(P, \mathbf{I}))$$



## Consequences

Results and techniques for Logic Programming can be exploited for datalog.

For example,

- proof procedures for Logic Programming (e.g., SLD resolution) can be applied to datalog (with some caveats, regarding for instance termination)
- datalog can be reduced by “grounding” to propositional logic programs

## Datalog Semantics via Least Fixpoint

The semantics of  $P$  on database instance  $\mathbf{I}$  of  $edb(P)$  is a special fixpoint:

**Theorem.** Let  $P$  be a datalog program and  $\mathbf{I}$  be a database instance. Then

1.  $\mathbf{T}_P$  has a least (wrt “ $\subseteq$ ”) fixpoint containing  $\mathbf{I}$ , denoted  $lfp(P, \mathbf{I})$ .
2. Moreover,  $lfp(P, \mathbf{I}) = MM(\mathcal{P}(P, \mathbf{I})) = P(\mathbf{I})$ .

Advantage: Constructive definition of  $P(\mathbf{I})$  by *fixpoint iteration*

### Proof-Theoretic Approach

Basic idea: The answer of a datalog program  $P$  on  $\mathbf{I}$  is given by the set of facts which can be *proved* from  $P$  and  $\mathbf{I}$ .

**Definition.** A *proof tree* for a fact  $A$  from  $\mathbf{I}$  and  $P$  is a labeled finite tree  $T$  such that

- each vertex of  $T$  is labeled by a fact
- the root of  $T$  is labeled by  $A$
- each leaf of  $T$  is labeled by a fact in  $\mathbf{I}$
- if a non-leaf of  $T$  is labeled with  $A_1$  and its children are labeled with  $A_2, \dots, A_n$ , then there exists a ground instance  $r$  of a rule in  $P$  such that

$$A_1 \leftarrow B(r) \text{ and } B(r) = \{A_2, \dots, A_n\}$$

35

Datalog

### Example (Same Generation)

$$P = \{ \begin{array}{l} r_1 : \text{sgc}(X, X) \leftarrow \text{person}(X) \\ r_2 : \text{sgc}(X, Y) \leftarrow \text{par}(X, X1), \text{sgc}(X1, Y1), \text{par}(Y, Y1) \end{array} \}$$

where  $edb(P) = \{\text{person}, \text{par}\}$  and  $idb(P) = \{\text{sgc}\}$

Consider  $\mathbf{I}$  as follows:

$$\begin{aligned} \mathbf{I}(\text{person}) &= \{ \langle \text{ann} \rangle, \langle \text{bertrand} \rangle, \langle \text{charles} \rangle, \langle \text{dorothy} \rangle, \\ &\quad \langle \text{evelyn} \rangle, \langle \text{fred} \rangle, \langle \text{george} \rangle, \langle \text{hilary} \rangle \} \\ \mathbf{I}(\text{par}) &= \{ \langle \text{dorothy}, \text{george} \rangle, \langle \text{evelyn}, \text{george} \rangle, \langle \text{bertrand}, \text{dorothy} \rangle, \\ &\quad \langle \text{ann}, \text{dorothy} \rangle, \langle \text{hilary}, \text{ann} \rangle, \langle \text{charles}, \text{evelyn} \rangle \}. \end{aligned}$$

Datalog

Example (Same Generation)/2

Proof tree for  $A = \text{sgc}(\text{ann}, \text{charles})$  from  $\mathbf{I}$  and  $P$ :

$\text{sgc}(\text{ann}, \text{charles})$

$r_2 :$        $\text{par}(\text{ann}, \text{dorothy})$        $\text{sgc}(\text{dorothy}, \text{evelyn})$        $\text{par}(\text{charles}, \text{evelyn})$

$r_2 :$        $\text{par}(\text{dorothy}, \text{george})$        $\text{sgc}(\text{george}, \text{george})$        $\text{par}(\text{evelyn}, \text{george})$

$r_1 :$        $\text{person}(\text{george})$

Proof Tree Construction

Different ways to construct a proof tree for  $A$  from  $P$  and  $\mathbf{I}$  exist

- **Bottom Up construction:** From leaves to root  
Intimately related to fixpoint approach
  - Define  $S \vdash_P B$  to prove fact  $B$  from facts  $S$  if  $B \in S$  or by a rule in  $P$
  - Give  $S = \mathbf{I}$  for granted
- **Top Down construction:** From root to leaves  
In Logic Programming view, consider program  $\mathcal{P}(P, \mathbf{I})$ .
  - This amounts to a set of logical sentences  $H_{\mathcal{P}(P, \mathbf{I})}$  of the form
$$\forall x_1 \cdots \forall x_m (R_1(\vec{x}_1) \vee \neg R_2(\vec{x}_2) \vee \neg R_3(\vec{x}_3) \vee \cdots \vee \neg R_n(\vec{x}_n))$$
  - Prove  $A = R(\vec{t})$  via resolution refutation, that is, that  $H_{\mathcal{P}(P, \mathbf{I})} \cup \{\neg A\}$  is unsatisfiable.

## Datalog and SLD Resolution

- Logic Programming uses SLD resolution
- SLD: Selection Rule Driven Linear Resolution for Definite Clauses
- For datalog programs  $P$  on  $\mathbf{I}$ , resp.  $\mathcal{P}(P, \mathbf{I})$ , things are simpler than for general logic programs (no function symbols, unification is easy)
- Also non-ground atoms can be handled (e.g.,  $\text{sgc}(\text{ann}, X)$ )

Let  $SLD(\mathcal{P})$  be the set of ground atoms provable with SLD Resolution from  $\mathcal{P}$ .

**Theorem.** For any datalog program  $P$  and database instance  $\mathbf{I}$ ,

$$SLD(\mathcal{P}(P, \mathbf{I})) = P(\mathbf{I}) = \mathbf{T}_{\mathcal{P}(P, \mathbf{I})}^{\infty} = \text{lfp}(\mathbf{T}_{\mathcal{P}(P, \mathbf{I})}) = MM(\mathcal{P}(P, \mathbf{I}))$$

Bases de Datos II

39

## SLD Resolution – Termination

Datalog

- Notice: Selection rule for next rule / atom to be considered for resolution might affect termination
- Prolog's strategy (leftmost atom / first rule) is problematic

Example:

```

child_of(karl, franz).
child_of(franz, frieda).
child_of(frieda, pia).
descendent_of(X, Y) ← child_of(X, Y).
descendent_of(X, Y) ← child_of(X, Z), descendent_of(Z, Y).
← descendent_of(karl, X).

```

Datalog

**SLD Resolution – Termination /2**

```
child_of(karl, franz).
child_of(franz, frieda).
child_of(frieda, pia).
descendent_of(X, Y) ← child_of(X, Y).
descendent_of(X, Y) ← descendent_of(X, Z), child_of(Z, Y).
← descendent_of(karl, X).
```

**SLD Resolution – Termination /3**

```
child_of(karl, franz).
child_of(franz, frieda).
child_of(frieda, pia).
descendent_of(X, Y) ← child_of(X, Y).
descendent_of(X, Y) ← descendent_of(X, Z),
                        descendent_of(Z, Y).
← descendent_of(karl, X).
```

### SQL3

- SQL2 does not support recursive queries:
  - Need to write PL/SQL or embedded SQL
- SQL3 supports recursive queries:
  - WITH statement to define recursive relation
  - The recursive relation is then used in a traditional SELECT FROM WHERE query

### Example: Bart's ancestors in SQL3

Datalog

Datalog Query:

```

AncestorDescendant(X, Y) ← ParentChild(X, Y)
AncestorDescendant(X, Y) ← ParentChild(X, Z), AncestorDescendant(Z, Y)
answer(X) ← AncestorDescendant(X, 'Bart')

```

SQL3 Query:

WITH

```

    RECURSIVE AncestorDescendent(ancestor, descendent) AS
        (SELECT * FROM ParentChild)
        UNION
        (SELECT pc.parent, ad.descendent
         FROM ParentChild pc, AncestorDescendent ad
         WHERE pc.child = ad.ancestor)

```

```

SELECT ancestor
FROM AncestorDescendent
WHERE descendent = 'Bart';

```

Datalog