

## Datalog (1)

Goal: Learn *the* advanced optimization technique central to the (logical) optimization of queries on graphs

(In RDBMS, “home” for this material: *recursive queries*)

3 lectures:

- Learn Datalog, syntax, semantics, execution model(s)
  - a query language that leverages notation of first-order logic
  - concise & accessible
  - very commonly used in research papers
- Learn the Optimizing Transform
  - Magic-sets
  - a.k.a. sideways information passing
- Importance for graph queries will become self-evident

slide thanks: mostly Ullman, also Michael Lam

10 Datalog I

Database Systems

1

## Datalog

Objectives, lecture 1:

- Introduce Datalog
  - Context, relation with/to
    - rule-based programming and logical inference
    - triggers
    - *business rules*, and supporting systems
  - recursion and querying graph-based data

slide thanks: mostly Ullman, also Michael Lam

10 Datalog I

Database Systems

2

## Logic As a Query Language

- ◆ If-then logic rules have been used in many systems.
  - ▶ Prolog
  - ▶ Triggers a.k.a. ECA rules: on **e**vent, if **c**ondition, **a**ction
  - ▶ Depending on your religion
    - AI, Knowledge-Base systems
      - Forward-chaining rules/production systems
      - Backward chaining rules, automatic proof systems
    - guarded command languages (do od)
      - Nod to Dijkstra and the formal methods community
      - Nondeterministic, fixed-point semantics

10 Datalog I

Database Systems

3

## Logic As a Query Language

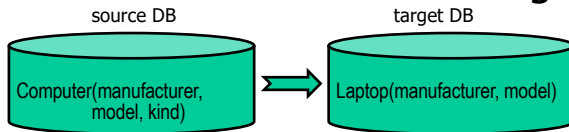
- ◆ Most important today: EII (Enterprise Information Integration)
  - ▶ Too many products to list
  - ▶ Special Mention: Logicblox (Datalog-based)
- ◆ Business logic/ workflow (BPEL, business process execution language)
  - ▶ IBM Websphere – ILOG rules
  - ▶ Redhat – DROOLS
  - ▶ ... many others

10 Datalog I

Database Systems

4

## View-based Information Integration



```
CREATE VIEW targetLaptop AS
SELECT manufacturer, model
FROM Computer
WHERE Computer.kind = 'Laptop'
```

Datalog:

```
Laptop(manufacturer, model)
<- Computer(manufacturer, model, kind) AND kind = 'Laptop'
```

10 Datalog I Database Systems 5

## Logic As a Query Language

- ◆ Business logic/ workflow (BPEL, business process execution language)

### Examples:

```
OvertimePay(person, date)
<- Role(person, title) AND OvertimeOK(title) AND Holiday(date)
```

```
VicePresidentSignatureRequired(purchase_id)
<- PurchaseOrder(purchase_id, amount) AND amount > $100,000
```

10 Datalog I

Database Systems

6

## With Datalog Came Recursive Queries

- ◆ Nonrecursive subset of Datalog is equivalent to the core relational algebra.
- ◆ Recursive rules extend relational algebra --- have been used to add recursion to SQL-99.

10 Datalog I

Database Systems

7

## Example Rule 1: a “happy drinker”

- ◆ Given
  - ◆ **Frequents(drinker,bar)** // defines drinker
  - ◆ **Likes(drinker,beer),** // what they like
  - ◆ **Sells(bar,beer,price).** // where sold, etc.
- ◆ Query who are the happy drinkers?
  - > drinkers whose bars serve their favorite beer.

10 Datalog I

Database Systems

8

## As a Datalog Rule

```
Happy(d) <- Frequents(d,bar) AND
           Likes(d,beer) AND
           Sells(bar,beer,p)
```

10 Datalog I

Database Systems

9

## Example: Interpretation

Happy(**d**) <- Frequents(**d**,**bar**) AND  
 Likes(**d**,**beer**) AND Sells(**bar**,**beer**,**p**)

Distinguished variable      Nondistinguished variables

Interpretation: drinker *d* is happy *if there exists* a bar, a beer, and a price *p* such that *d* frequents the bar, likes the beer, and the bar sells the beer at price *p*.

10 Datalog I

Database Systems

10

## Straight to the Critical Connection

Frequents		Likes		Sells		
Person	Bar	Person	Beer	Bar	Beer	Price
Juan	Dog & Duck	Juan	Dos XX	Dog & Duck	Dos XX	\$5.00
Dan	Crown & Anchor	Dan	Shiner	Crown & Anchor	Dos XX	\$4.00

Frequents(d,bar)      AND Likes(d,beer)      AND Sells(bar,beer,p)

```
select Person
from Frequents, Likes, Sells
where Frequents.Person = Likes.Person and
      Frequents.Bar = Sells.Bar and
      Likes.Beer = Sells.Beer
```

datalog, good: concise, looks familiar.    bad: positional dependence

10 Datalog I

Database Systems

11

## A Word on Notation

Frequents		Likes		Sells		
Person	Bar	Person	Beer	Bar	Beer	Price
Juan	Dog & Duck	Juan	Dos XX	Dog & Duck	Dos XX	\$5.00
Dan	Crown & Anchor	Dan	Shiner	Crown & Anchor	Dos XX	\$4.00

Frequents(d,bar) AND Likes(d,beer) AND Sells(bar,beer,p)

bad: positional dependence, d/Person, bar/Bar  
 don't care, use underscore: Frequents(\_,bar)

### Variables/Predicates

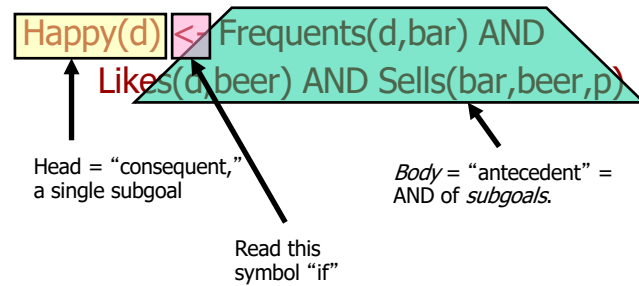
- Datalog (text): Variables, start with lower case, Predicates, upper case
- Datalog & Prolog (everywhere else):  
     Predicates, start with lower case, Variables upper case

10 Datalog I

Database Systems

12

## Anatomy of a Rule



10 Datalog I

Database Systems

13

## Subgoals Are Atoms

- ◆ An *atom* is a *predicate*, or relation name with variables or constants as arguments.
- ◆ The head is an atom; the body is the AND of one or more atoms.
- ◆ Convention: Predicates begin with a capital, variables begin with lower-case.

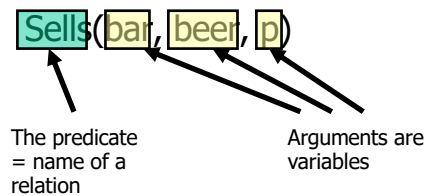
// At risk of confusion, your text repairs (changes) tradition: Prolog and Datalog, Predicates begin with lowercase letter, variables, an uppercase letter

10 Datalog I

Database Systems

14

## Example: Atom



arity of a predicate: number of arguments in a predicate

10 Datalog I

Database Systems

15

## Interpreting Rules

- ◆ A variable appearing in the head is called *distinguished*; otherwise it is *nondistinguished*.
- ◆ Rule meaning: The head is true of the distinguished variables if there exist values of the nondistinguished variables that make all subgoals of the body true.

10 Datalog I

Database Systems

16

## Arithmetic Subgoals

- ◆ In addition to relations as predicates, a predicate for a subgoal of the body can be an arithmetic comparison.
  - ▶ We write such subgoals in the usual way, e.g.:  $x < y$ .

10 Datalog I

Database Systems

17

## Example: Arithmetic

- ◆ A beer is “cheap” if there are at least two bars that sell it for under \$2.

$\text{Cheap}(\text{beer}) \leftarrow \text{Sells}(\text{bar1}, \text{beer}, p1) \text{ AND}$   
 $\text{Sells}(\text{bar2}, \text{beer}, p2) \text{ AND } p1 < 2.00$   
 $\text{AND } p2 < 2.00 \text{ AND } \text{bar1} \neq \text{bar2}$

10 Datalog I

Database Systems

18

## Negated Subgoals

- ◆ We may put NOT in front of a subgoal, to negate its meaning.
- ◆ Example:
- ◆ given additional relation  
 $\text{HappyHour}(\text{bar}, \text{hours})$
- ◆ unhappy drinker: none of his favorite bars have a happy hour

$\text{UnHappy}(d) \leftarrow \text{Frequents}(d, \text{bar}) \text{ AND}$   
 $\text{Not HappyHour}(\text{bar})$

10 Datalog I

Database Systems

19

## Relational Interpretation of Negation

- ◆ Set difference
  - ▶ A.k.a. negation as failure
  - ▶ Closed world assumption

10 Datalog I

Database Systems

20

## Negated Subgoals

- ◆ unhappy drinker: none of his favorite bars have a happy hour

$\text{UnHappy}(d) \leftarrow \text{Frequents}(d, \text{bar}) \text{ AND Not HappyHour}(\text{bar})$

$\text{Frequents}(d, \text{bar}) - \text{HappyHour}(\text{bar})$

Frequents: {(dan, 'crown anchor') (dan, 'jack allens')}

If

HappyHour: {'crown anchor'}, 'jack allens' remains, dan is happy

If

HappyHour: {'crown anchor', 'jack allens'}, Dan has no where to have happy hour, Dan is *bound* to d, Dan is unhappy

10 Datalog I

Database Systems

21

## Closed and Open Worlds

- ◆ Given

$\text{unHappy}(d) \leftarrow \text{Frequents}(d, \text{bar}) \text{ AND Not HappyHour}(\text{bar})$   
 $\text{Frequents}(d, \text{bar}) - \text{HappyHour}(\text{bar})$

Frequents: {(dan, 'crown anchor') (dan, 'jack allens')}

HappyHour: {'crown anchor', 'jack allens'},

Suppose, Dan also frequents Z'Tejas to drink. **But** it is not recorded in the database

Dan has a place to drink, he's happy, but database "doesn't know it"

Database is omniscient (knows everything): **Closed world**, set difference correctly computes NOT

Database does not know everything: **Open World**, can't logically conclude NOT

10 Datalog I

Database Systems

22

## Datalog assumes a closed world

- ◆ that plus other properties make things computable.

- ◆ next property *safety*

10 Datalog I

Database Systems

23

## Safe Rules

- ◆ A rule is *safe* if: // all the variables that appear in the  
 // head, negated subgoal, or arithmetic subgoal,  
 // appear in a nonnegated subgoal (can be bound)

1. Each distinguished variable,
2. Each variable in an arithmetic subgoal,
3. Each variable in a negated subgoal,  
 also appears in a nonnegated,  
 relational subgoal.

- ◆ Datalog allows only safe rules.

10 Datalog I

Database Systems

24

## Example: Unsafe Rules

- ◆ Each of the following is unsafe and not allowed:
  1.  $S(x) \leftarrow R(y)$
  2.  $S(x) \leftarrow R(y) \text{ AND NOT } R(x)$
  3.  $S(x) \leftarrow R(y) \text{ AND } x < y$
- ◆ In each case, an infinity of  $x$ 's can satisfy the rule, even if  $R$  is a finite relation.

10 Datalog I

Database Systems

25

## Datalog Programs

- ◆ A *Datalog program* is a collection of rules.
- ◆ In a program, predicates can be either
  1. EDB = *Extensional Database* = stored table. // base tables
  2. IDB = *Intensional Database* = relation defined by rules.
- ◆ Never both! No EDB in heads.

10 Datalog I

Database Systems

26

## Evaluating Datalog Programs

- ◆ As long as there is no recursion, we can pick an order to evaluate the IDB predicates, so that all the predicates in the body of its rules have already been evaluated.
- ◆ If an IDB predicate has more than one rule, each rule contributes tuples to its relation.

10 Datalog I

Database Systems

27

## Example: Datalog Program

- ◆ Using EDB *Sells(bar, beer, price)* and *Beers(name, manf)*, find the manufacturers of beers Joe doesn't sell.

```
JoeSells(b) <- Sells(' Joe' 's Bar' , b, p)
Answer(m) <- Beers(b,m)
               AND NOT JoeSells(b)
```

10 Datalog I

Database Systems

28

## Definition of Recursion

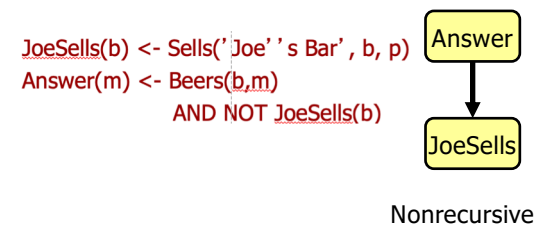
- ◆ Form a *dependency graph* whose nodes = IDB predicates.
- ◆ Arc  $X \rightarrow Y$  if and only if there is a rule with  $X$  in the head and  $Y$  in the body.
- ◆ Cycle = recursion; no cycle = no recursion.

10 Datalog I

Database Systems

29

## Example: Dependency Graphs



10 Datalog I

Database Systems

30

## Recursive Example

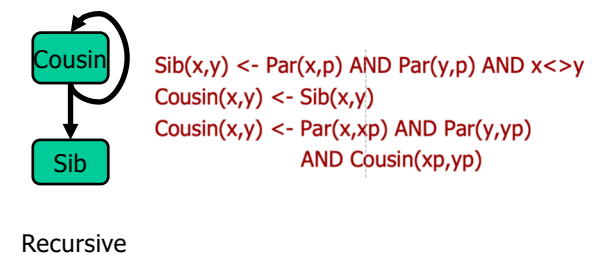
- ◆ EDB:  $\text{Par}(c, p) = p$  is a parent of  $c$ .
  - ◆ Generalized cousins: people with common ancestors one or more generations back:
- $\text{Sib}(x, y) \leftarrow \text{Par}(x, p) \text{ AND } \text{Par}(y, p) \text{ AND } x \neq y$   
 $\text{Cousin}(x, y) \leftarrow \text{Sib}(x, y)$   
 $\text{Cousin}(x, y) \leftarrow \text{Par}(x, xp) \text{ AND } \text{Par}(y, yp) \text{ AND } \text{Cousin}(xp, yp)$

10 Datalog I

Database Systems

31

## Example: Dependency Graphs



10 Datalog I

Database Systems

32



## Expressive Power of Datalog

- ◆ Without recursion, Datalog can express all and only the queries of core relational algebra.
  - ▶ The same as SQL select-from-where, without aggregation and grouping.
- ◆ But with recursion, Datalog can express more than these languages.
- ◆ Yet still not Turing-complete.

10 Datalog I

Database Systems

33

## Evaluation of Non-Recursive Datalog

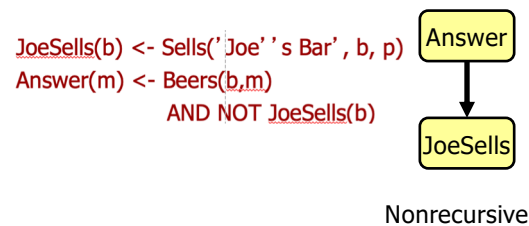
- ◆ Topologically sort the dependency graph.
- ◆ Evaluate rules in sorted order

10 Datalog I

Database Systems

34

## Example: Dependency Graphs



10 Datalog I

Database Systems

35

## Consider Subset of Cousin Program

```

(1) sibling(X,Y) :- parent(X,Z) & parent(Y,Z) & X≠Y.
(2) cousin(X,Y) :- parent(X,Xp) & parent(Y,Yp) &
                    sibling(Xp,Yp).
  
```

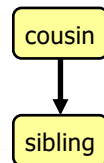
10 Datalog I

Database Systems

36

## Consider Subset of Cousin Program

```
(1) sibling(X,Y) :- parent(X,Z) & parent(Y,Z) & X≠Y.
(2) cousin(X,Y) :- parent(X,Xp) & parent(Y,Yp) &
    sibling(Xp,Yp).
```



10 Datalog I

Database Systems

37

## Evaluating Recursive Rules

◆ The following works when there is no negation:

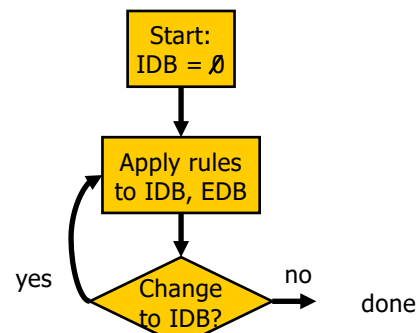
1. Start by assuming all IDB relations are empty.
2. Repeatedly evaluate the rules using the EDB and the previous IDB, to get a new IDB.
3. End when no change to IDB. (fixed-point)

10 Datalog I

Database Systems

38

## The “Naïve” Evaluation Algorithm



10 Datalog I

Database Systems

39

## Example: Evaluation of Cousin

- ◆ We'll proceed in rounds to infer Sib facts (red) and Cousin facts (green).
- ◆ Remember the rules:

```
Sib(x,y) <- Par(x,p) AND Par(y,p) AND x<>y
Cousin(x,y) <- Sib(x,y)
Cousin(x,y) <- Par(x,xp) AND Par(y,yp)
    AND Cousin(xp,yp)
```

10 Datalog I

Database Systems

40

## Seminaive Evaluation

- ◆ Since the EDB never changes, on each round we only get new IDB tuples if we use at least one IDB tuple that was obtained on the previous round.
- ◆ Saves work; lets us avoid rediscovering *most* known facts.
  - ▶ A fact could still be derived in a second way.

10 Datalog I

Database Systems

41

## Incremental Maintenance of a View

$$(R \cup S) \mid X \mid T = (R \mid X \mid T) \cup (S \mid X \mid T)$$

So, if View,  $V = (S \mid X \mid T)$

- ◆ Insert Into S, rows R, i.e.  $S' = (R \cup S)$
- ◆ New view contents,  $V'$ ,  $V' = S' \mid X \mid T$

$$\begin{aligned} V' &= (R \cup S) \mid X \mid T \\ &= (R \mid X \mid T) \cup (S \mid X \mid T) \\ &= (R \mid X \mid T) \cup V \end{aligned}$$

7:Semantics/Triggers/V  
iews

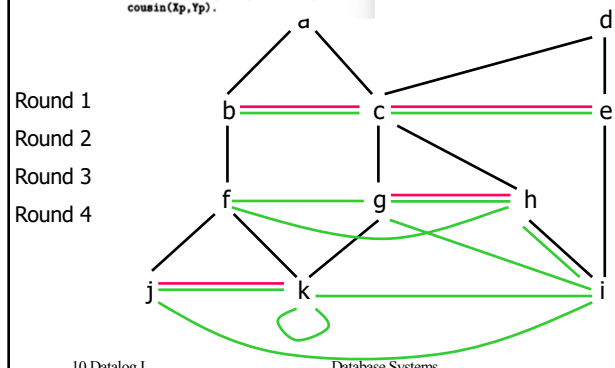
Database Systems

42

## Par Data: Parent Above Child

- (1)  $\text{ sibling}(X,Y) :- \text{ parent}(X,Z) \ \& \ \text{ parent}(Y,Z) \ \& \ X \neq Y.$
- (2)  $\text{ cousin}(X,Y) :- \text{ parent}(X,Xp) \ \& \ \text{ parent}(Y,Yp) \ \& \ \text{ sibling}(Xp,Yp).$
- (3)  $\text{ cousin}(X,Y) :- \text{ parent}(X,Xp) \ \& \ \text{ parent}(Y,Yp) \ \& \ \text{ cousin}(Xp,Yp).$

Red: sibling rows  
Green: cousin rows



10 Datalog I

Database Systems

43