

Introduction to Database Systems CSE 414

Lecture 9: Datalog

Class Overview

- Unit 1: Intro
- Unit 2: Relational Data Models and Query Languages
 - Data models, SQL, Relational Algebra, **Datalog**
- Unit 3: Non-relational data
- Unit 4: RDBMS internals and query optimization
- Unit 5: Parallel query processing
- Unit 6: DBMS usability, conceptual design
- Unit 7: Transactions

What is Datalog?

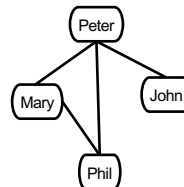
- Another query language for relational model
 - Designed in the 80's
 - Simple, concise, elegant
 - Extends relational queries with recursion
- Today is a hot topic:
 - Souffle (we will use in HW4)
 - Eve <http://witheve.com/>
 - Differential datalog <https://github.com/frankmcsherry/differential-dataflow>
 - Beyond databases in many research projects: network protocols, static program analysis



- Open-source implementation of Datalog DBMS
- Under active development
- Commercial implementations are available
 - More difficult to set up and use
- "sqlite" of Datalog
 - Set-based rather than bag-based
- Install in your VM
 - Run `sudo yum install souffle` in terminal
 - More details in upcoming HW4

Why bother with *yet* another relational query language?

Example: storing FB friends



Or

Person1	Person2	is_friend
Peter	John	1
John	Mary	0
Mary	Phil	1
Phil	Peter	1
...

As a graph

As a relation

We will learn the tradeoffs of different data models later this quarter

Compute your friends graph

p1	p2	isFriend
Peter	John	1
John	Mary	0
Mary	Phil	1
Phil	Peter	1
...

Friends(p1, p2, isFriend)

```
SELECT f.p2
FROM Friends as f
WHERE f.p1 = 'me' AND f.isFriend = 1
```

My own friends

```
SELECT f1.p2
FROM Friends as f1,
(SELECT f.p2
FROM Friends as f
WHERE f.p1 = 'me' AND
f.isFriend = 1) as f2
WHERE f1.p1 = f2.p2 AND
f1.isFriend = 1
```

My FoF

My FoFoF... My FoFoFoF...

Datalog allows us to write recursive queries easily

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year) ← Schema

Datalog: Facts and Rules

Facts = tuples in the database

Rules = queries

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

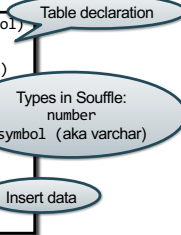
Datalog: Facts and Rules

Facts = tuples in the database

Rules = queries

```
.decl Actor(id:number, fname:symbol, lname:symbol)
.decl Casts(id:number, mid:number)
.decl Movie(id:number, name:symbol, year:number)
```

```
Actor(344759, 'Douglas', 'Fowley').
Casts(344759, 29851).
Casts(355713, 29000).
Movie(7909, 'A Night in Armour', 1910).
Movie(29000, 'Arizona', 1940).
Movie(29445, 'Ave Maria', 1940).
```



Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Datalog: Facts and Rules

Facts = tuples in the database

Rules = queries

```
Actor(344759, 'Douglas', 'Fowley').
Casts(344759, 29851).
Casts(355713, 29000).
Movie(7909, 'A Night in Armour', 1910).
Movie(29000, 'Arizona', 1940).
Movie(29445, 'Ave Maria', 1940).
```

```
Q1(y) :- Movie(x,y,z), z=1940.
```

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Datalog: Facts and Rules

Facts = tuples in the database

Rules = queries

```
Actor(344759, 'Douglas', 'Fowley').
Casts(344759, 29851).
Casts(355713, 29000).
Movie(7909, 'A Night in Armour', 1910).
Movie(29000, 'Arizona', 1940).
Movie(29445, 'Ave Maria', 1940).
```

```
Q1(y) :- Movie(x,y,z), z=1940.
```

Find Movies made in 1940

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Datalog: Facts and Rules

Facts = tuples in the database

Rules = queries

```
Actor(344759, 'Douglas', 'Fowley').
Casts(344759, 29851).
Casts(355713, 29000).
Movie(7909, 'A Night in Armour', 1910).
Movie(29000, 'Arizona', 1940).
Movie(29445, 'Ave Maria', 1940).
```

```
Q1(y) :- Movie(x,y,z), z=1940.
```

SQL

```
SELECT name
FROM Movie
WHERE year = 1940
```

Find Movies made in 1940

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Datalog: Facts and Rules

Facts = tuples in the database **Rules** = queries

```

Actor(344759, 'Douglas', 'Fowley').
Casts(344759, 29851).
Casts(355713, 29000).
Movie(7909, 'A Night in Armour', 1910).
Movie(29000, 'Arizona', 1940).
Movie(29445, 'Ave Maria', 1940).
  
```

Q1(y) :- Movie(x,y,z), z=1940.

id → name → year

Order of variable matters!

Find Movies made in 1940

CSE 414 - Autumn 2018 13

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Datalog: Facts and Rules

Facts = tuples in the database **Rules** = queries

```

Actor(344759, 'Douglas', 'Fowley').
Casts(344759, 29851).
Casts(355713, 29000).
Movie(7909, 'A Night in Armour', 1910).
Movie(29000, 'Arizona', 1940).
Movie(29445, 'Ave Maria', 1940).
  
```

Q1(y) :- Movie(idontCare,y,z), z=1940.

Find Movies made in 1940

CSE 414 - Autumn 2018 14

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Datalog: Facts and Rules

Facts = tuples in the database **Rules** = queries

```

Actor(344759, 'Douglas', 'Fowley').
Casts(344759, 29851).
Casts(355713, 29000).
Movie(7909, 'A Night in Armour', 1910).
Movie(29000, 'Arizona', 1940).
Movie(29445, 'Ave Maria', 1940).
  
```

Q1(y) :- Movie(_,y,z), z=1940.

= "don't care" variables

Find Movies made in 1940

CSE 414 - Autumn 2018 15

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Datalog: Facts and Rules

Facts = tuples in the database **Rules** = queries

```

Actor(344759, 'Douglas', 'Fowley').
Casts(344759, 29851).
Casts(355713, 29000).
Movie(7909, 'A Night in Armour', 1910).
Movie(29000, 'Arizona', 1940).
Movie(29445, 'Ave Maria', 1940).
  
```

Q1(y) :- Movie(x,y,z), z=1940.

Q2(f,l) :- Actor(z,f,l), Casts(z,l).

Find Actors who acted in Movies made in 1940

CSE 414 - Autumn 2018 16

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Datalog: Facts and Rules

Facts = tuples in the database **Rules** = queries

```

Actor(344759, 'Douglas', 'Fowley').
Casts(344759, 29851).
Casts(355713, 29000).
Movie(7909, 'A Night in Armour', 1910).
Movie(29000, 'Arizona', 1940).
Movie(29445, 'Ave Maria', 1940).
  
```

Q1(y) :- Movie(x,y,z), z=1940.

Q2(f,l) :- Actor(z,f,l), Casts(z,x), Movie(x,y,1940).

Find Actors who acted in Movies made in 1940

CSE 414 - Autumn 2018 17

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Datalog: Facts and Rules

Facts = tuples in the database **Rules** = queries

```

Actor(344759, 'Douglas', 'Fowley').
Casts(344759, 29851).
Casts(355713, 29000).
Movie(7909, 'A Night in Armour', 1910).
Movie(29000, 'Arizona', 1940).
Movie(29445, 'Ave Maria', 1940).
  
```

Q1(y) :- Movie(x,y,z), z=1940.

Q2(f,l) :- Actor(z,f,l), Casts(z,x), Movie(x,y,1940).

Q3(f,l) :- Actor(z,f,l), Casts(z,x1), Movie(x1,y1,1910), Casts(z,x2), Movie(x2,y2,1940).

Find Actors who acted in Movies made in 1940

CSE 414 - Autumn 2018 18

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Datalog: Facts and Rules

Facts = tuples in the database **Rules** = queries

```

Actor(344759, 'Douglas', 'Fowley').
Casts(344759, 29851).
Casts(355713, 29000).
Movie(7909, 'A Night in Armour', 1910).
Movie(29000, 'Arizona', 1940).
Movie(29445, 'Ave Maria', 1940).
  
```

Q1(y) :- Movie(x,y,z), z=1940.

Q2(f,l) :- Actor(z,f,l), Casts(z,x),
Movie(x,y,1940).

Q3(f,l) :- Actor(z,f,l), Casts(z,x1), Movie(x1,y1,1910),
Casts(z,x2), Movie(x2,y2,1940).

Find Actors who acted in a Movie in 1940 and in one in 1910

CSE 414 - Autumn 2018 19

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Datalog: Facts and Rules

Facts = tuples in the database **Rules** = queries

```

Actor(344759, 'Douglas', 'Fowley').
Casts(344759, 29851).
Casts(355713, 29000).
Movie(7909, 'A Night in Armour', 1910).
Movie(29000, 'Arizona', 1940).
Movie(29445, 'Ave Maria', 1940).
  
```

Q1(y) :- Movie(x,y,z), z=1940.

Q2(f,l) :- Actor(z,f,l), Casts(z,x),
Movie(x,y,1940).

Q3(f,l) :- Actor(z,f,l), Casts(z,x1), Movie(x1,y1,1910),
Casts(z,x2), Movie(x2,y2,1940).

Extensional Database Predicates = EDB = Actor, Casts, Movie
Intensional Database Predicates = IDB = Q1, Q2, Q3

CSE 414 - Autumn 2018 20

Datalog: Terminology

head body

atom atom atom (aka subgoal)

Q2(f, l) :- Actor(z,f,l), Casts(z,x), Movie(x,y,1940).

f, l = head variables
x,y,z = existential variables

CSE 414 - Autumn 2018 21

More Datalog Terminology

Q(args) :- R1(args), R2(args), ...

- $R_i(\text{args}_i)$ called an *atom*, or a *relational predicate*
- $R_i(\text{args}_i)$ evaluates to true when relation R_i contains the tuple described by args.
 - Example: Actor(344759, 'Douglas', 'Fowley') is true
- In addition we can also have arithmetic predicates
 - Example: $z > 1940$.
- Book uses AND instead of , **Q**(args) :- R1(args) AND R2(args) ...

CSE 414 - Autumn 2018 22

Datalog program

- A Datalog program consists of several rules
- Importantly, rules may be recursive!
 - Recall CSE 143!
- Usually there is one distinguished predicate that's the output
- We will show an example first, then give the general semantics.

CSE 414 - Autumn 2018 23

Example

R encodes a graph
e.g., connected cities

R=

1	2
2	1
2	3
1	4
3	4
4	5

CSE 414 - Autumn 2018 24

R encodes a graph e.g., connected cities

Multiple rules for the same IDB means OR

Example

$$T(x,y) :- R(x,y).$$

$$T(x,y) :- R(x,z), T(z,y).$$

What does it compute?

R=

1	2
2	1
2	3
1	4
3	4
4	5

CSE 414 - Autumn 2018 25

R encodes a graph e.g., connected cities

Example

$$T(x,y) :- R(x,y).$$

$$T(x,y) :- R(x,z), T(z,y).$$

What does it compute?

R=

1	2
2	1
2	3
1	4
3	4
4	5

Initially: T is empty.

--	--

CSE 414 - Autumn 2018 26

R encodes a graph e.g., connected cities

Example

$$T(x,y) :- R(x,y).$$

$$T(x,y) :- R(x,z), T(z,y).$$

What does it compute?

R=

1	2
2	1
2	3
1	4
3	4
4	5

Initially: T is empty.

--	--

First iteration: T =

1	2
2	1
2	3
1	4
3	4
4	5

First rule generates this

Second rule generates nothing (because T is empty)

CSE 414 - Autumn 2018 27

R encodes a graph e.g., connected cities

Example

$$T(x,y) :- R(x,y).$$

$$T(x,y) :- R(x,z), T(z,y).$$

What does it compute?

R=

1	2
2	1
2	3
1	4
3	4
4	5

Initially: T is empty.

--	--

First iteration: T =

1	2
2	1
2	3
1	4
3	4
4	5

Second iteration: T =

1	2
2	1
2	3
1	4
2	3
3	4
1	1
2	2
2	4
1	3
2	4
1	5
3	5

First rule generates this

Second rule generates this

New facts

CSE 414 - Autumn 2018 28

R encodes a graph e.g., connected cities

Example

$$T(x,y) :- R(x,y).$$

$$T(x,y) :- R(x,z), T(z,y).$$

What does it compute?

R=

1	2
2	1
2	3
1	4
3	4
4	5

Initially: T is empty.

--	--

First iteration: T =

1	2
2	1
2	3
1	4
3	4
4	5

Second iteration: T =

1	2
2	1
2	3
1	4
3	4
4	5
1	1
2	2
1	3
2	4
1	5
3	5

Both rules

First rule

Second rule

Third iteration: T =

1	2
2	1
2	3
1	4
3	4
4	5
1	1
2	2
1	3
2	4
1	5
3	5
2	5

New fact

CSE 414 - Autumn 2018 29

R encodes a graph e.g., connected cities

Example

$$T(x,y) :- R(x,y).$$

$$T(x,y) :- R(x,z), T(z,y).$$

What does it compute?

R=

1	2
2	1
2	3
1	4
3	4
4	5

Initially: T is empty.

--	--

First iteration: T =

1	2
2	1
2	3
1	4
3	4
4	5

Second iteration: T =

1	2
2	1
2	3
1	4
2	3
3	4
1	1
2	2
1	3
2	4
1	5
3	5

Third iteration: T =

1	2
2	1
2	3
1	4
3	4
4	5
1	1
2	2
1	3
2	4
1	5
3	5
2	5

Fourth iteration: T = (same)

No new facts. DONE

CSE 414 - Autumn 2018 30

More Features

- Aggregates
- Grouping
- Negation

CSE 414 - Autumn 2018

31

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Aggregates

[aggregate name] <var> : { [relation to compute aggregate on] }

min x : { Actor(x, y, _) , y = 'John' }

Q(minId) :- minId = min x : { Actor(x, y, _) , y = 'John' }

Assign variable to the value of the aggregate

Meaning (in SQL)

```
SELECT min(id) as minId
FROM Actor as a
WHERE a.name = 'John'
```

Aggregates in Souffle:

- count
- min
- max
- sum

CSE 414 - Autumn 2018

32

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Aggregates

[aggregate name] <var> : { [relation to compute aggregate on] }

min x : { Actor(x, y, _) , y = 'John' }

head

Q(minId, y) :- minId = min x : { Actor(x, y, _) }

What does this even mean???



Can't use variable that are not aggregated in the outer /head atoms

CSE 414 - Autumn 2018

33

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Counting

Q(c) :- c = count : { Actor(_, y, _) , y = 'John' }

No variable here!

Meaning (in SQL, assuming no NULLs)

```
SELECT count(*) as c
FROM Actor as a
WHERE a.name = 'John'
```

CSE 414 - Autumn 2018

34

Actor(id, fname, lname)
Casts(pid, mid)
Movie(id, name, year)

Grouping

Q(y,c) :- Movie(_, y), c = count : { Movie(_, y) }

Meaning (in SQL)

```
SELECT m.year, count(*)
FROM Movie as m
GROUP BY m.year
```

CSE 414 - Autumn 2018

35

ParentChild(p,c)

Example

For each person, compute the total number of descendants

```
// for each person, compute his/her descendants
```

CSE 414 - Autumn 2018

36

Example

For each person, compute the total number of descendants

```
// for each person, compute his/her descendants
D(x,y) :- ParentChild(x,y).
```

Example

For each person, compute the total number of descendants

```
// for each person, compute his/her descendants
D(x,y) :- ParentChild(x,y).
D(x,z) :- D(x,y), ParentChild(y,z).
```

Example

For each person, compute the total number of descendants

```
// for each person, compute his/her descendants
D(x,y) :- ParentChild(x,y).
D(x,z) :- D(x,y), ParentChild(y,z).

// For each person, count the number of descendants
```

Example

For each person, compute the total number of descendants

```
// for each person, compute his/her descendants
D(x,y) :- ParentChild(x,y).
D(x,z) :- D(x,y), ParentChild(y,z).

// For each person, count the number of descendants
T(p,c) :- D(p,_), c = count : { D(p,y) }.
```

Example

For each person, compute the total number of descendants

```
// for each person, compute his/her descendants
D(x,y) :- ParentChild(x,y).
D(x,z) :- D(x,y), ParentChild(y,z).

// For each person, count the number of descendants
T(p,c) :- D(p,_), c = count : { D(p,y) }.

// Find the number of descendants of Alice
```

Example

For each person, compute the total number of descendants

```
// for each person, compute his/her descendants
D(x,y) :- ParentChild(x,y).
D(x,z) :- D(x,y), ParentChild(y,z).

// For each person, count the number of descendants
T(p,c) :- D(p,_), c = count : { D(p,y) }.

// Find the number of descendants of Alice
Q(d) :- T(p,d), p = "Alice".
```

ParentChild(p,c)

Negation: use “!”

Find all descendants of Alice,
who are not descendants of Bob

```
// for each person, compute his/her descendants
D(x,y) :- ParentChild(x,y).
D(x,z) :- D(x,y), ParentChild(y,z).

// Compute the answer: notice the negation
Q(x) :- D("Alice",x), !D("Bob",x).
```

CSE 414 - Autumn 2018 43

ParentChild(p,c)

Safe Datalog Rules

Here are unsafe datalog rules. What's “unsafe” about them ?

```
U1(x,y) :- ParentChild("Alice",x), y != "Bob"
U2(x) :- ParentChild("Alice",x), !ParentChild(x,y)
```

CSE 414 - Autumn 2018 44

ParentChild(p,c)

Safe Datalog Rules

Here are unsafe datalog rules. What's “unsafe” about them ?

```
U1(x,y) :- ParentChild("Alice",x), y != "Bob"
U2(x) :- ParentChild("Alice",x), !ParentChild(x,y)
```

Hold for every y other than “Bob”
U1 = infinite!

CSE 414 - Autumn 2018 45

ParentChild(p,c)

Safe Datalog Rules

Here are unsafe datalog rules. What's “unsafe” about them ?

```
U1(x,y) :- ParentChild("Alice",x), y != "Bob"
U2(x) :- ParentChild("Alice",x), !ParentChild(x,y)
```

Hold for every y other than “Bob”
U1 = infinite!

Want Alice's childless children,
but we get all children x (because
there exists some y that x is not parent of y)

CSE 414 - Autumn 2018 46

ParentChild(p,c)

Safe Datalog Rules

Here are unsafe datalog rules. What's “unsafe” about them ?

```
U1(x,y) :- ParentChild("Alice",x), y != "Bob"
U2(x) :- ParentChild("Alice",x), !ParentChild(x,y)
```

A datalog rule is safe if every variable appears
in some positive relational atom

CSE 414 - Autumn 2018 47

Stratified Datalog

- Recursion does not cope well with aggregates or negation
- Example: what does this mean?


```
A() :- !B().
B() :- !A().
```
- A datalog program is stratified if it can be partitioned into *strata*
 - Only IDB predicates defined in strata 1, 2, ..., n may appear under ! or agg in stratum n+1.
- Many Datalog DBMSs (including souffle) accepts only stratified Datalog.

CSE 414 - Autumn 2018 48

Stratified Datalog

<pre>D(x,y) :- ParentChild(x,y). D(x,z) :- D(x,y), ParentChild(y,z). T(p,c) :- D(p,_), c = count : { D(p,y) }. Q(d) :- T(p,d), p = "Alice".</pre>	<p>Stratum 1</p> <hr/> <p>Stratum 2</p>
---	---

May use D
in an agg since it was
defined in previous
stratum

CSE 414 - Autumn 2018

49

Stratified Datalog

<pre>D(x,y) :- ParentChild(x,y). D(x,z) :- D(x,y), ParentChild(y,z). T(p,c) :- D(p,_), c = count : { D(p,y) }. Q(d) :- T(p,d), p = "Alice".</pre>	<p>Stratum 1</p> <hr/> <p>Stratum 2</p>
---	---

May use D
in an agg since it was
defined in previous
stratum

<pre>D(x,y) :- ParentChild(x,y). D(x,z) :- D(x,y), ParentChild(y,z). Q(x) :- D("Alice",x), !D("Bob",x).</pre>	<p>Stratum 1</p> <hr/> <p>Stratum 2</p>
---	---

May use !D

<pre>A() :- !B(). B() :- !A().</pre>	<p>Non-stratified</p>
--------------------------------------	-----------------------

Cannot use !A

50

Stratified Datalog

- If we don't use aggregates or negation, then the Datalog program is already stratified
- If we do use aggregates or negation, it is usually quite natural to write the program in a stratified way

CSE 414 - Autumn 2018

51