

CSE544

Data Management

Lectures 7: Datalog

Announcements

- HW2 is graded (thanks Kyle!)
- HW3 is due this coming Sunday
- P2 Project Milestone due next Sunday
- HW4 is posted and due on Friday, 3/7

Motivation

- Data processing today require iteration.
 - Common solution: external driver
 - The adventurous: SQL WITH RECURSIVE
- Datalog is a language designed for recursive queries

Datalog

- Designed in the 80's: simple, concise
- Used for research on recursive queries
- Only research prototypes exists...
- ...in HW4 we use Souffle

Outline

- Syntax
- Single rule
- Multiple rules
- Recursion
- Naïve algorithm
- Semi-naïve algorithm
- Non-monotone operations

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

← Schema

Datalog: Facts and Rules

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, year)

Datalog: Facts and Rules

Facts = tuples in the database

Rules = queries

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, 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).

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, year)

Datalog: Facts and Rules

Facts = tuples in the database

```
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).
```

Rules = queries

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

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, year)

Datalog: Facts and Rules

Facts = tuples in the database

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).

Rules = queries

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

Find titles of movies made in 1940

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, year)

Datalog: Facts and Rules

Facts = tuples in the database

```
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).
```

Rules = queries

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

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

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, year)

Datalog: Facts and Rules

Facts = tuples in the database

```
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).
```

Rules = queries

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

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

Find first and last names of Actors
who acted in Movies made in 1940

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, year)

Datalog: Facts and Rules

Facts = tuples in the database

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).

Rules = queries

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)

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, year)

Datalog: Facts and Rules

Facts = tuples in the database

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).

Rules = queries

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

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, year)

Datalog: Facts and Rules

Facts = tuples in the database

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).

Rules = queries

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

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, year)

Datalog: Facts and Rules

Facts = tuples in the database

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).

Rules = queries

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

Outline

- Syntax
- Single rule
- Multiple rules
- Recursion
- Naïve algorithm
- Semi-naïve algorithm
- Non-monotone operations

Semantics of One Rule

SQL and single Datalog rule have similar semantics:

- SQL: nested loops of tuples
- Datalog rule: nested loops of variables:
 - Map variables to values in the database
 - Check that all atoms occur in the input

Semantics of One Rule

Actor

id	fname	lname
344759	Douglas	Fowley

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

Casts

pid	mid
344759	29851
355713	29000

Movie

id	title	year
7909	A Night in Armour	1910
29000	Arizona	1940
29445	Ave Maria	1940

Semantics of One Rule

Actor

id	fname	lname
344759	Douglas	Fowley

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

Casts

pid	mid
344759	29851
355713	29000

x	y	z
29000	Arizona	1940
29445	Ave Maria	1940

Movie

id	title	year
7909	A Night in Armour	1910
29000	Arizona	1940
29445	Ave Maria	1940

Semantics of One Rule

Actor

id	fname	lname
344759	Douglas	Fowley

Casts

pid	mid
344759	29851
355713	29000

Movie

id	title	year
7909	A Night in Armour	1910
29000	Arizona	1940
29445	Ave Maria	1940

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

x	y	z
29000	Arizona	1940
29445	Ave Maria	1940

Answer

y
Arizona
Ave Maria

Semantics of One Rule

Actor

id	fname	lname
344759	Douglas	Fowley

Casts

pid	mid
344759	29851
355713	29000

Movie

id	title	year
7909	A Night in Armour	1910
29000	Arizona	1940
29445	Ave Maria	1940

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

x	y	z
29000	Arizona	1940
29445	Ave Maria	1940

Answer

y
Arizona
Ave Maria

Usually displayed like this:

Q1('Arizona')
Q1('Ave Maria')

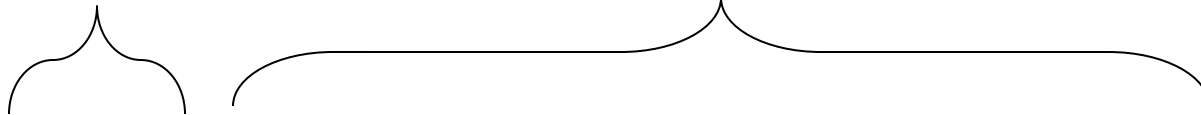
Anatomy of a Rule

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

Anatomy of a Rule

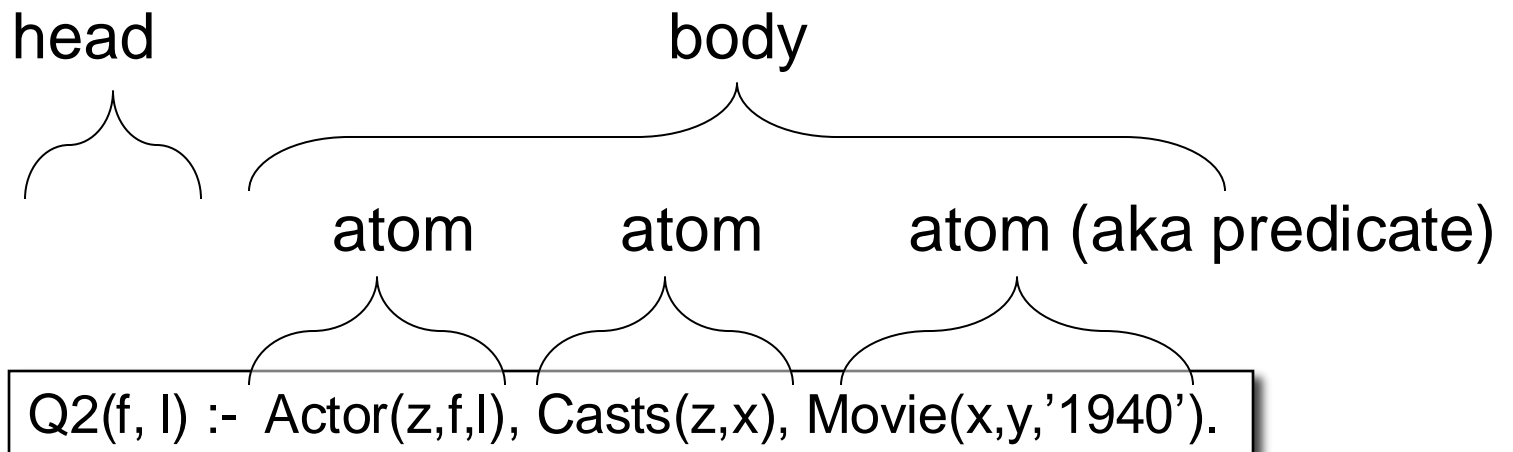
head

body

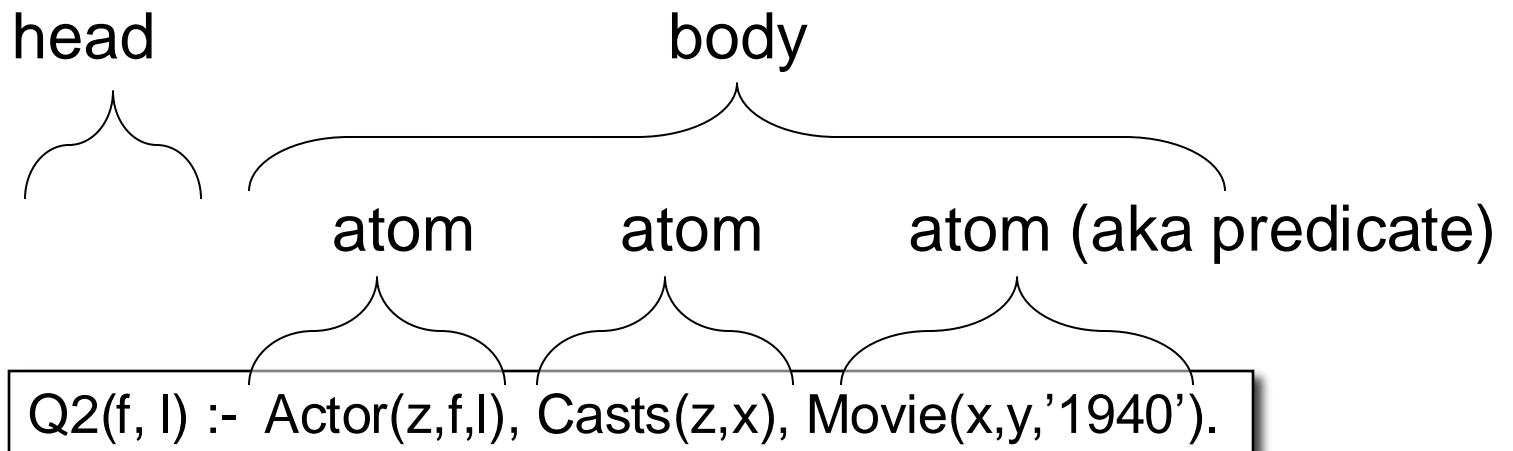


```
Q2(f, l) :- Actor(z,f,l), Casts(z,x), Movie(x,y,'1940').
```


Anatomy of a Rule



Anatomy of a Rule



f, l = head variables

x,y,z = existential variables

Discussion

- Rule body: a select-join expression
- Rule head specifies what to return
- Datalog program = multiple rules

Outline

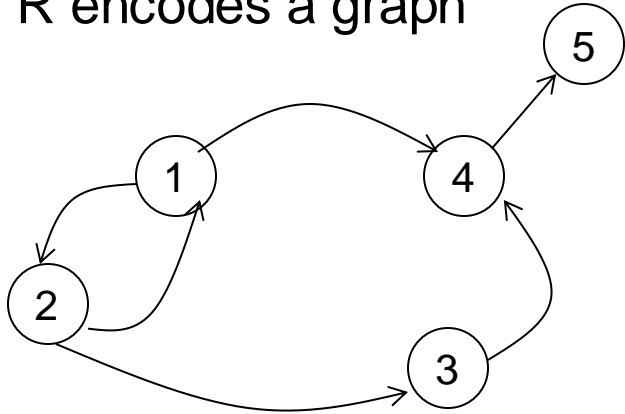
- Syntax
- Single rule
- Multiple rules
- Recursion
- Naïve algorithm
- Semi-naïve algorithm
- Non-monotone operations

Multiple Rules

- Rule body may have EDBs and IDBs
- Same IDB may be in the head of multiple rules

Multiple Rules

R encodes a graph

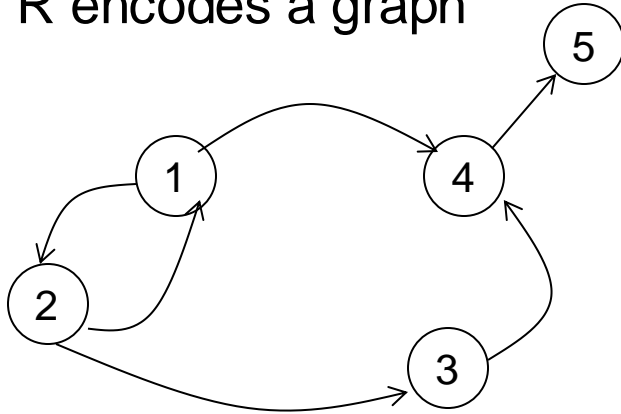


R=

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

Multiple Rules

R encodes a graph



R=

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

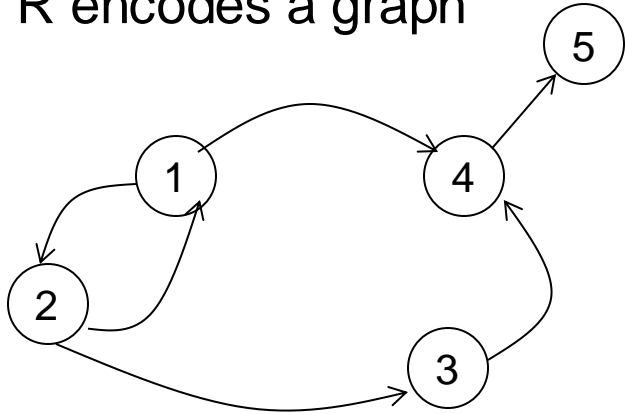
Entering data

```
R(1,2).  
R(2,1).  
...  
R(4,5).
```

As facts

Multiple Rules

R encodes a graph



R=

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

Entering data

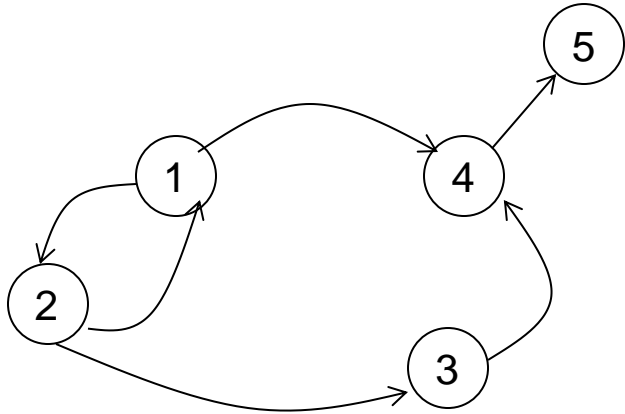
```
R(1,2).  
R(2,1).  
...  
R(4,5).
```

As facts

```
.input R("file-name")
```

From a file

Multiple Rules



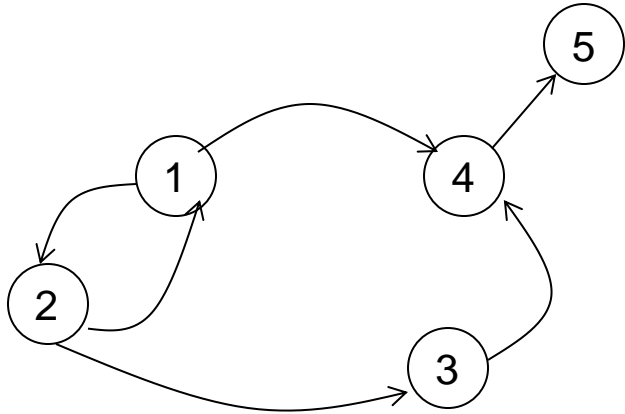
$A(x) :- R(1,z),R(z,x)$

R=

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

Answer??

Multiple Rules



R=

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

$A(x) \text{ :- } R(1,z),R(z,x)$

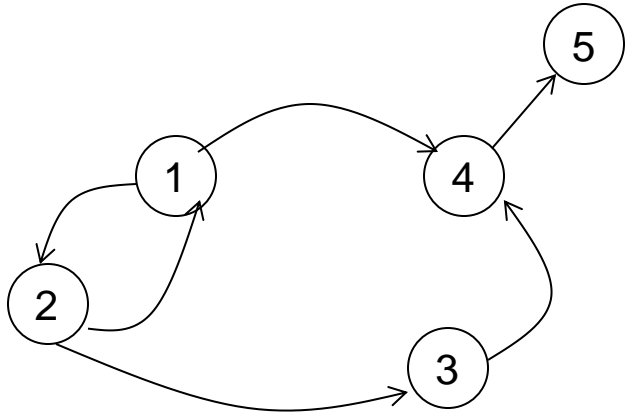
Answer:

A(1)

A(3)

A(5)

Multiple Rules



R=

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

$A(x) :- R(1,z),R(z,x)$

$B(x) :- A(z),R(z,x)$

Answer:

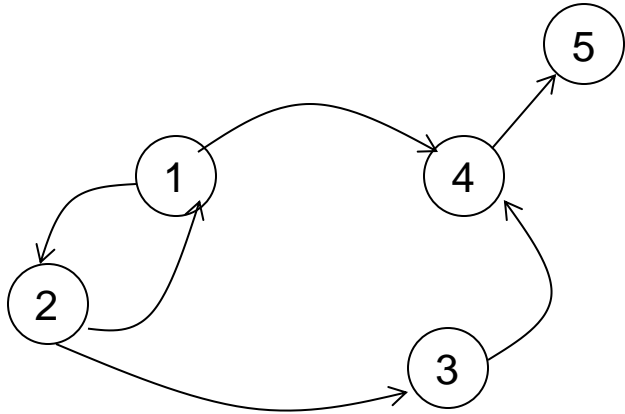
A(1)

A(3)

A(5)

Answer??

Multiple Rules



R=

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

$A(x) :- R(1,z),R(z,x)$

$B(x) :- A(z),R(z,x)$

Answer:

A(1)

A(3)

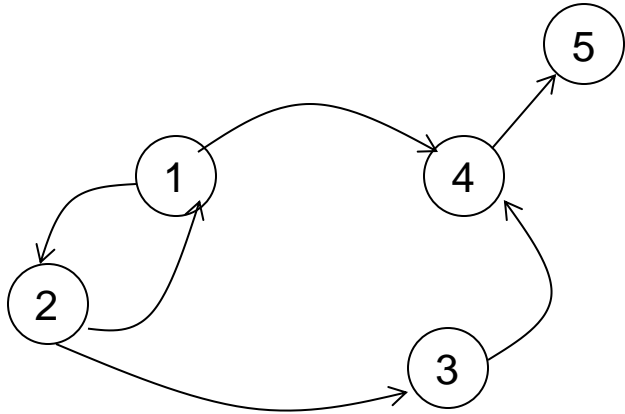
A(5)

B(2)

B(4)

Answer??

Multiple Rules



R=

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

$A(x) :- R(1,z),R(z,x)$

$B(x) :- A(z),R(z,x)$

Answer:

A(1)

A(3)

A(5)

B(2)

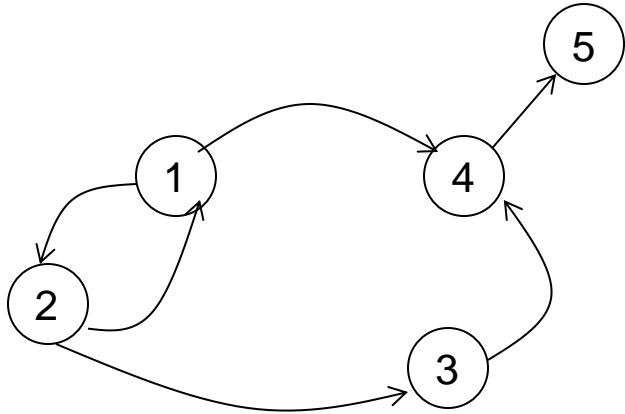
B(4)

Answer??

Set semantics:
B(4) occurs
only once

Multiple Rules

Order of rules does not matter



R=

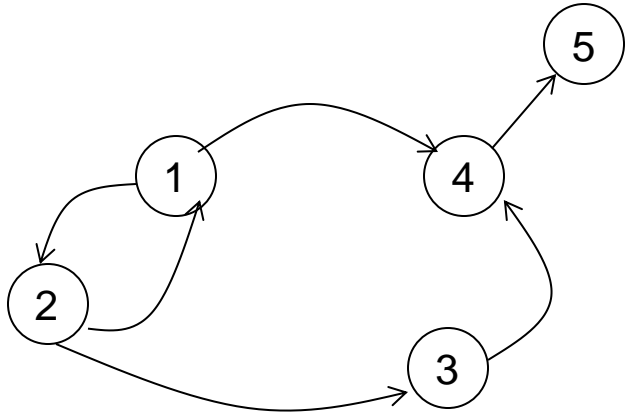
1	2
2	1
2	3
1	4
3	4
4	5

$B(x) :- A(z), R(z, x)$
 $A(x) :- R(1, z), R(z, x)$

Answer:

A(1) B(2)
A(3) B(4)
A(5)

Multiple Rules



$C(x) :- R(1,x),R(x,3)$

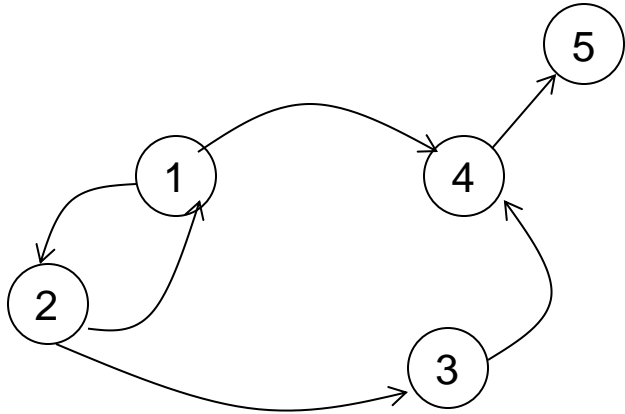
$C(x) :- R(1,x),R(3,x)$

R=

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

Multiple Rules

Multiple rules
with same head
means "union"



R=

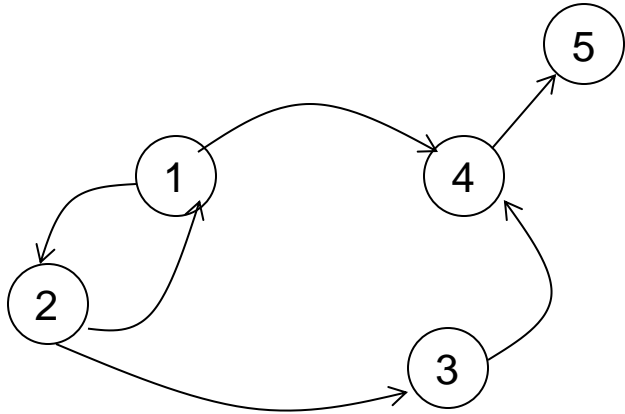
1	2
2	1
2	3
1	4
3	4
4	5

$C(x) :- R(1,x),R(x,3)$

$C(x) :- R(1,x),R(3,x)$

Multiple Rules

Multiple rules
with same head
means "union"



$C(x) :- R(1,x),R(x,3)$

$C(x) :- R(1,x),R(3,x)$

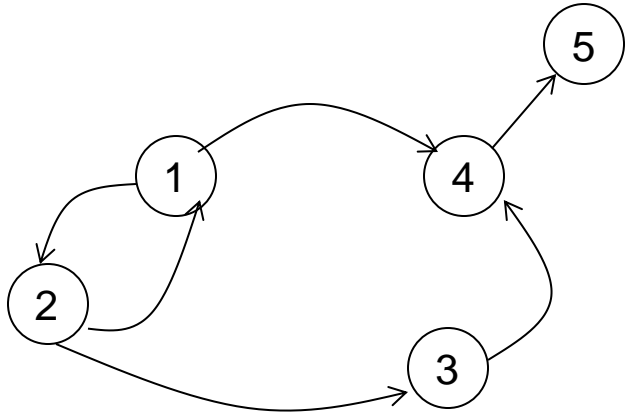
R=

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

Answer??

Multiple Rules

Multiple rules with same head means "union"



$C(x) :- R(1,x),R(x,3)$
 $C(x) :- R(1,x),R(3,x)$

R=

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

Answer??

C(2)

C(4)

Discussion

- Rule body: EDBs and/or IDBs
- Rule heads w/ the same IDB: union
- Rule order does not matter
- Rules may be recursive: next

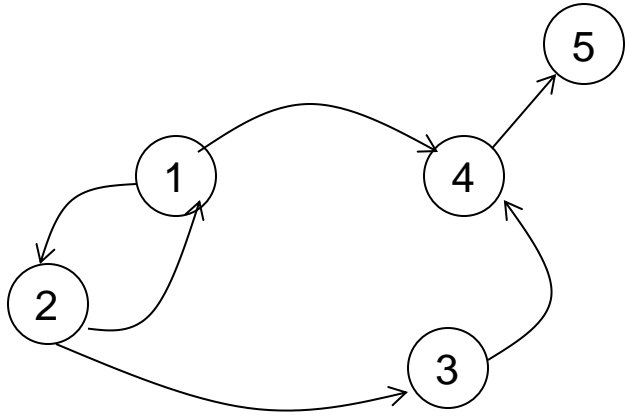
Outline

- Syntax
- Single rule
- Multiple rules
- Recursion
- Naïve algorithm
- Semi-naïve algorithm
- Non-monotone operations

Recursion in Datalog

- Rules may be recursive
- Bottom up evaluation:
 - Start with empty IDBs
 - Compute rules repeatedly, increasing IDBs
 - Stop when no more change

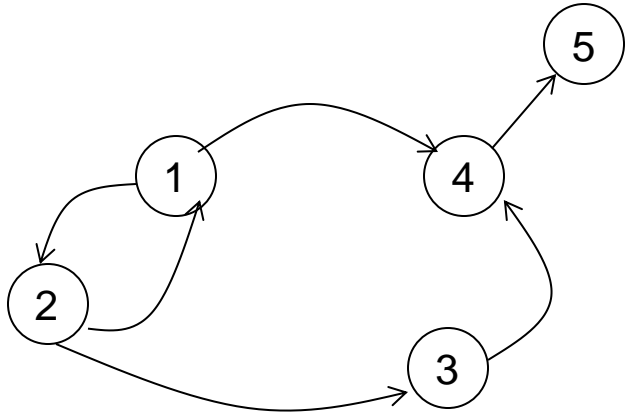
Example



R=

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

Example



R=

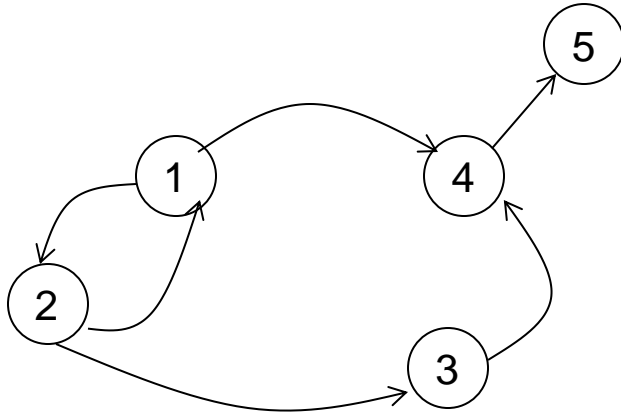
1	2
2	1
2	3
1	4
3	4
4	5

$T(x,y) \text{ :- } R(x,y)$

$T(x,y) \text{ :- } R(x,z), T(z,y)$

What does it compute?

Example



R=

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

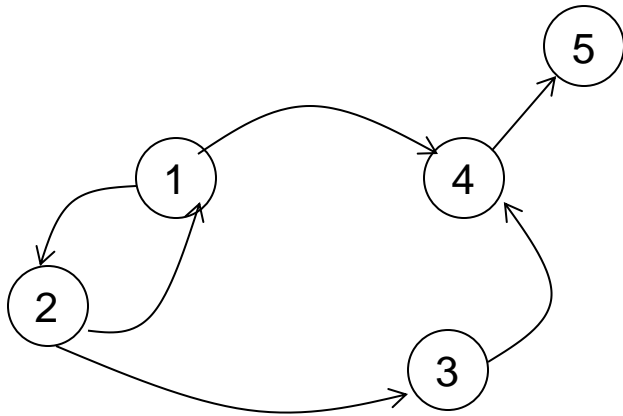
Initially:
T is empty.



$T(x,y) \text{ :- } R(x,y)$
 $T(x,y) \text{ :- } R(x,z), T(z,y)$

What does it compute?

Example



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

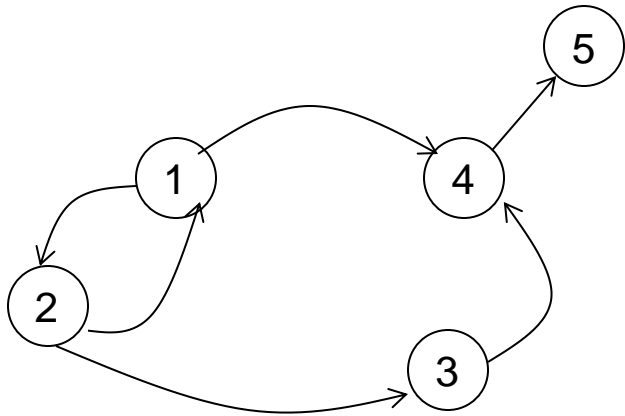
First rule generates this

Second rule
generates nothing
(because T is empty)

What does
it compute?

Example

What does it compute?



$T(x,y) \text{ :- } R(x,y)$
 $T(x,y) \text{ :- } R(x,z), T(z,y)$

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 =

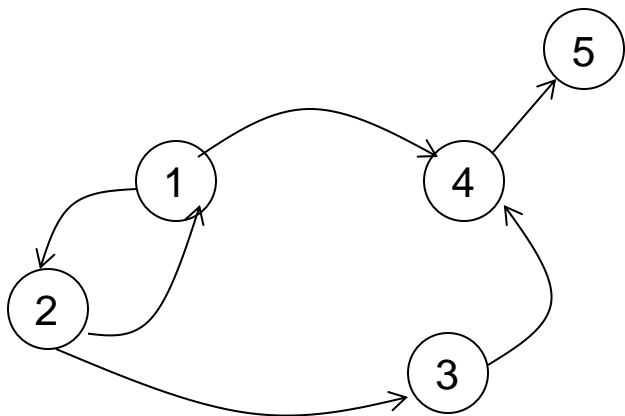
1	2
2	1
2	3
1	4
3	4
4	5
1	1
2	2
1	3
2	4
1	5
3	5

First rule generates this

Second rule generates this

New facts

Example



$T(x,y) \text{ :- } R(x,y)$
 $T(x,y) \text{ :- } R(x,z), T(z,y)$

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

New fact

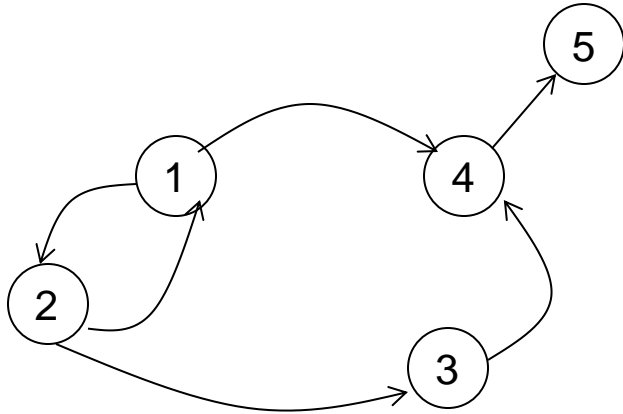
Third iteration:

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

Both rules
First rule
Second rule

Example



$T(x,y) \text{ :- } R(x,y)$
 $T(x,y) \text{ :- } R(x,z), T(z,y)$

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

Third iteration:

T =

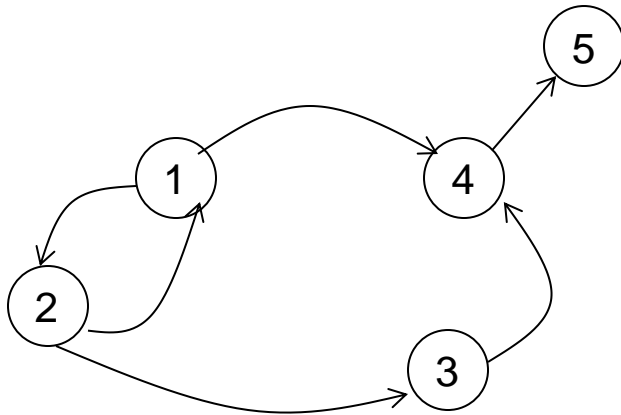
1	2
2	1
2	3
1	4
3	4
4	5
1	1
2	2
1	3
2	4
1	5
3	5
2	5

Fourth iteration

T = (same)

No new facts.
DONE

Example



$T(x,y) \text{ :- } R(x,y)$
 $T(x,y) \text{ :- } R(x,z), T(z,y)$

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

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

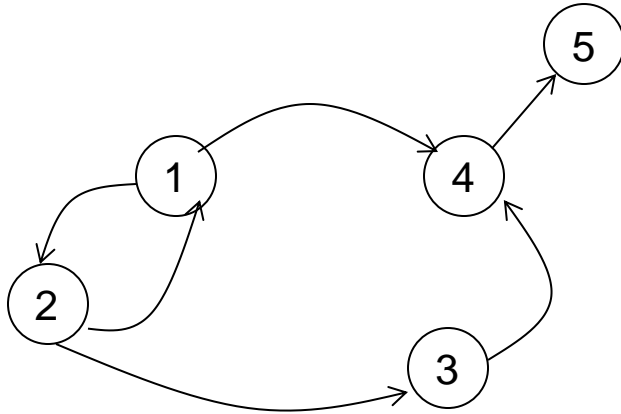
T = (same)

No new facts.
DONE

Iteration k computes pairs (x,y) connected by path of length ≤ k

Example

Computes the transitive closure of R



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

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

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

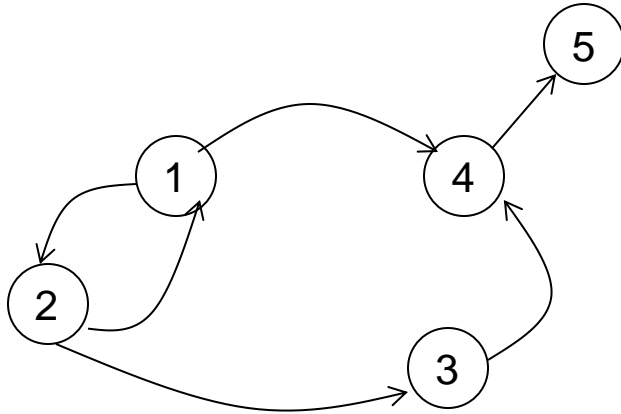
Fourth iteration

T = (same)

No new facts.
DONE

Iteration k computes pairs (x,y) connected by path of length $\leq k$

Three Equivalent Programs



R=

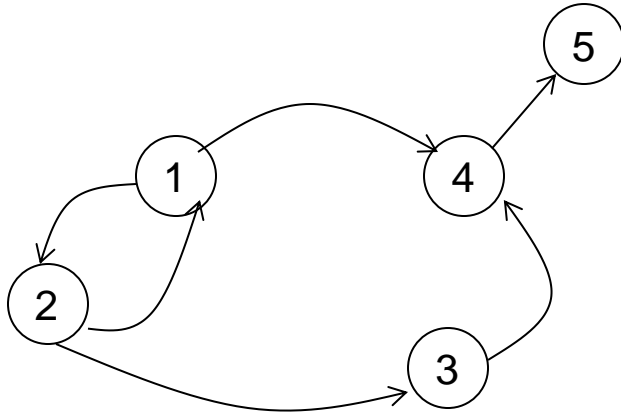
1	2
2	1
2	3
1	4
3	4
4	5

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

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

Right linear

Three Equivalent Programs



R=

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

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

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

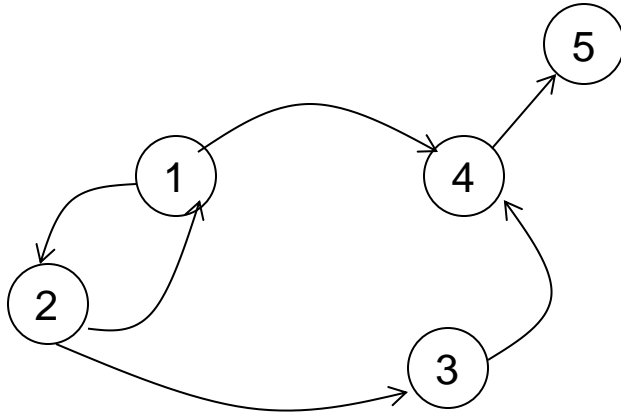
Right linear

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

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

Left linear

Three Equivalent Programs



R=

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

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

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

Right linear

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

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

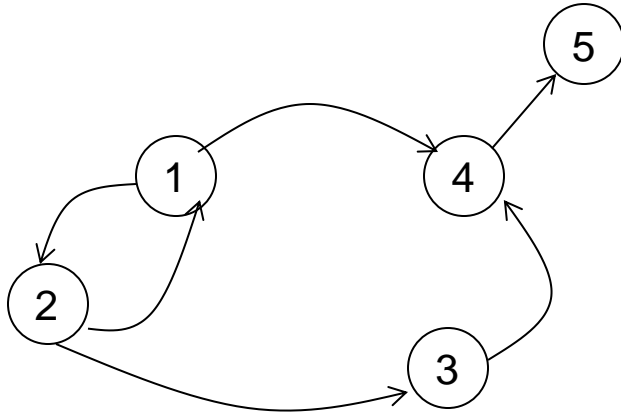
Left linear

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

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

Non-linear

Three Equivalent Programs



R=

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

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

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

Right linear

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

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

Left linear

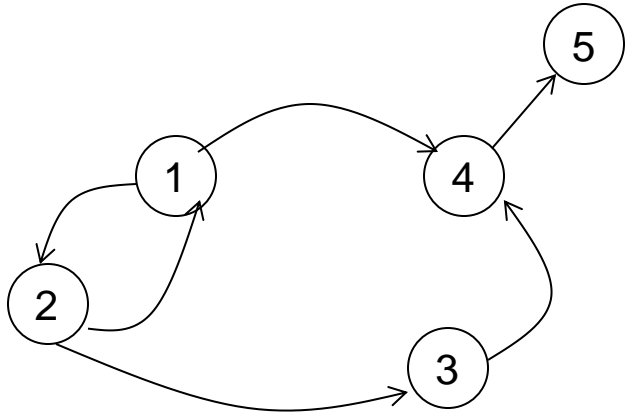
$T(x,y) :- R(x,y)$

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

Non-linear

Question: how many iterations does each require?

Three Equivalent Programs



R=

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

#iterations =
diameter

#iterations =
log(diameter)

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

Right linear

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

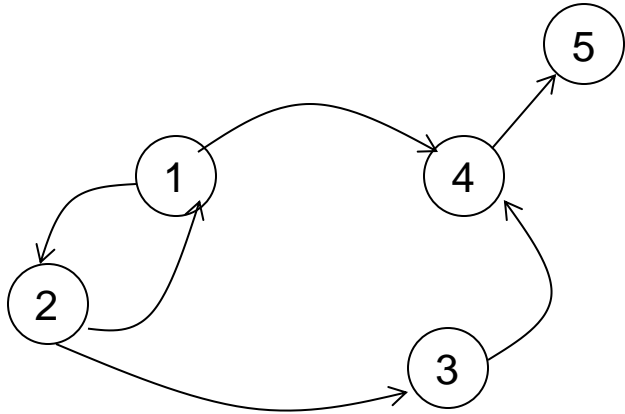
Left linear

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

Non-linear

Question: how many iterations does each require?

Multiple IDBs

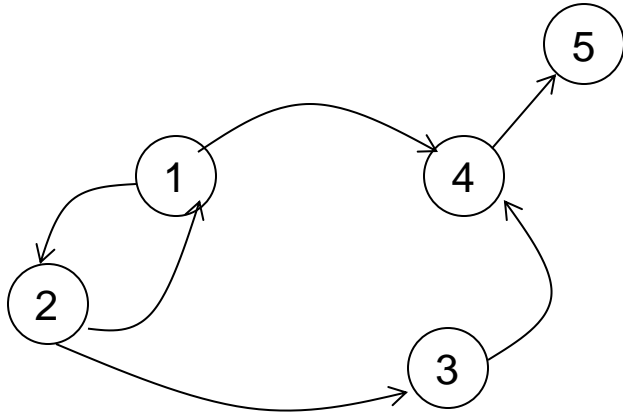


Find pairs of nodes (x,y) connected by a path of even length

R=

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

Multiple IDBs



Find pairs of nodes (x,y) connected by a path of even length

R=

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

Odd $(x,y) :- R(x,y)$

Even $(x,y) :- Odd(x,z), R(z,y)$

Odd $(x,y) :- Even(x,z), R(z,y)$

Two IDBs: Odd (x,y) and Even (x,y)

Recap

- Rules may be recursive
- If every rule body has at most one recursive IDB, then the program is called **linear**. Otherwise: **non-linear**.
- Evaluation is bottom-up
This is called the Naïve Algorithm (next)

Outline

- Syntax
- Single rule
- Multiple rules
- Recursion
- Naïve algorithm
- Semi-naïve algorithm
- Non-monotone operations

Naïve Evaluation Algorithm

- Every rule* \rightarrow SPJ query

*SPJ = select-project-join

+USPJ = union-select-project-join

Naïve Evaluation Algorithm

- Every rule* \rightarrow SPJ query

$T(x,z) \text{ :- } R(x,y), T(y,z), C(y, \text{'green'})$

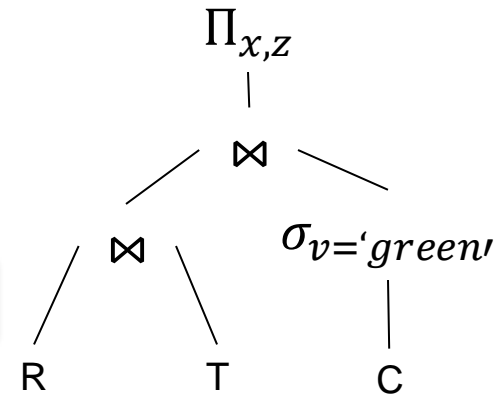
*SPJ = select-project-join

+USPJ = union-select-project-join

Naïve Evaluation Algorithm

- Every rule* \rightarrow SPJ query

$T(x,z) :- R(x,y), T(y,z), C(y,'green')$



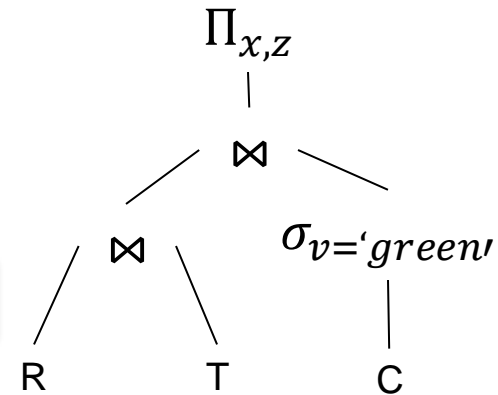
*SPJ = select-project-join

+USPJ = union-select-project-join

Naïve Evaluation Algorithm

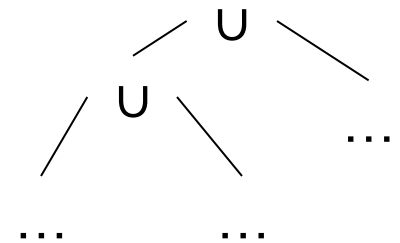
- Every rule* \rightarrow SPJ query

$T(x,z) :- R(x,y), T(y,z), C(y,'green')$



- Multiple rules same head+ \rightarrow USPJ

$T(x,y) :- \dots$
 $T(x,y) :- \dots$
 \dots



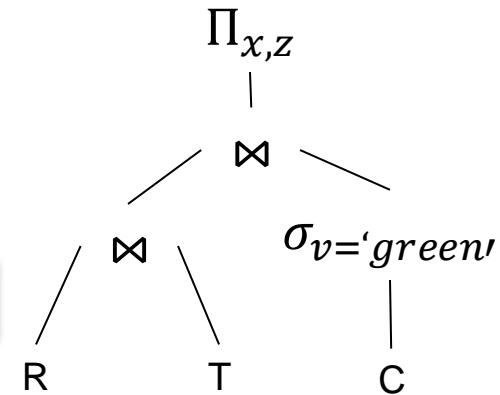
*SPJ = select-project-join

+USPJ = union-select-project-join

Naïve Evaluation Algorithm

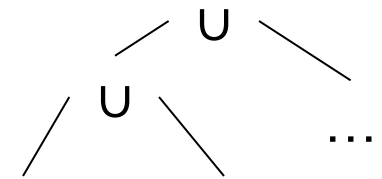
- Every rule* \rightarrow SPJ query

$T(x,z) :- R(x,y), T(y,z), C(y,'green')$



- Multiple rules same head+ \rightarrow USPJ

$T(x,y) :- \dots$
 $T(x,y) :- \dots$
 \dots



- Naïve Algorithm:

$IDB_0 := \emptyset$
repeat $IDB_{t+1} := USPJ(IDB_t)$
until no more change

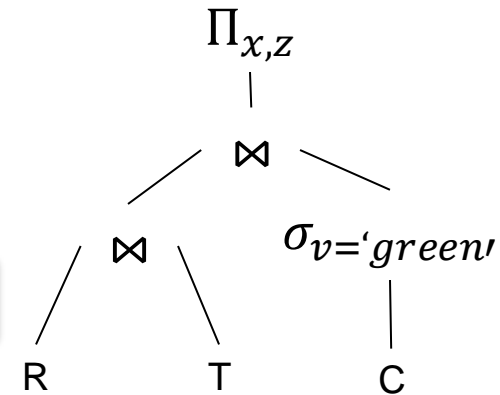
*SPJ = select-project-join

+USPJ = union-select-project-join

Naïve Evaluation Algorithm

- Every rule* \rightarrow SPJ query

$T(x,z) :- R(x,y), T(y,z), C(y,'green')$



SQL optimizes the plan.

Souffle does not optimize the plan ☹️

Computes joins left-to-right.

Equivalently, left-deep plan.

*SPJ = select-project-join

+USPJ = union-select-project-join

Naïve Evaluation Algorithm

- Naïve Algorithm:

```
IDB0 := ∅  
repeat IDBt+1 := USPJ(IDBt)  
until no more change
```

Called the
Immediate Consequence Operator
ICO

*SPJ = select-project-join
+USPJ = union-select-project-join

Example

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

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

Example

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

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

$T := \emptyset;$

repeat

$T := R \cup \Pi_{x,y}(R \bowtie T);$

until [no more change]

Example

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

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

$T := \emptyset;$

repeat

$T := R \cup \Pi_{x,y}(R \bowtie T);$

until [no more change]

A light blue speech bubble with a tail pointing towards the bottom right of the slide, containing the text 'ICO'.

ICO

Example

- When multiple IDBs: need to compute their new values *in parallel*:

Odd(x,y) :- R(x,y)

Even(x,y) :- Odd(x,z),R(z,y)

Odd(x,y) :- Even(x,z),R(z,y)

Example

- When multiple IDBs: need to compute their new values *in parallel*:

```
Odd(x,y) :- R(x,y)
```

```
Even(x,y) :- Odd(x,z),R(z,y)
```

```
Odd(x,y) :- Even(x,z),R(z,y)
```

```
Odd :=  $\emptyset$ ; Even :=  $\emptyset$ ;
```

```
repeat
```

```
Evennew :=  $\Pi_{x,y}(\text{Odd} \bowtie R)$ ;
```

```
Oddnew :=  $R \cup \Pi_{x,y}(\text{Even} \bowtie R)$ ;
```

Example

- When multiple IDBs: need to compute their new values *in parallel*:

```
Odd(x,y) :- R(x,y)
```

```
Even(x,y) :- Odd(x,z),R(z,y)
```

```
Odd(x,y) :- Even(x,z),R(z,y)
```

```
Odd :=  $\emptyset$ ; Even :=  $\emptyset$ ;
```

```
repeat
```

```
  Evennew :=  $\Pi_{x,y}(\text{Odd} \bowtie R)$ ;
```

```
  Oddnew :=  $R \cup \Pi_{x,y}(\text{Even} \bowtie R)$ ;
```

```
Odd:=Oddnew
```

```
Even:=Evennew
```

Example

- When multiple IDBs: need to compute their new values in parallel:

```
Odd(x,y) :- R(x,y)
```

```
Even(x,y) :- Odd(x,z),R(z,y)
```

```
Odd(x,y) :- Even(x,z),R(z,y)
```

```
Odd :=  $\emptyset$ ; Even :=  $\emptyset$ ;
```

```
repeat
```

```
  Evennew :=  $\Pi_{x,y}(\text{Odd} \bowtie R)$ ;
```

```
  Oddnew :=  $R \cup \Pi_{x,y}(\text{Even} \bowtie R)$ ;
```

```
  if Odd=Oddnew  $\wedge$  Even=Evennew  
    then break
```

```
  Odd:=Oddnew
```

```
  Even:=Evennew
```

Example

- When multiple IDBs: need to compute their new values *in parallel*:

```
Odd(x,y) :- R(x,y)
Even(x,y) :- Odd(x,z),R(z,y)
Odd(x,y) :- Even(x,z),R(z,y)
```

```
Odd :=  $\emptyset$ ; Even :=  $\emptyset$ ;
```

```
repeat
```

```
Evennew :=  $\Pi_{x,y}(\text{Odd} \bowtie R)$ ;
```

```
Oddnew :=  $R \cup \Pi_{x,y}(\text{Even} \bowtie R)$ ;
```

```
if Odd=Oddnew  $\wedge$  Even=Evennew  
then break
```

```
Odd:=Oddnew
```

```
Even:=Evennew
```



ICO

Example

Theorem The Naïve Algorithm terminates after a number of iteration that is at most a polynomial in the size of the database

Before we prove this, a review of **monotone queries**

Monotone Queries

- A query with input relations R, S, T, \dots is called monotone if, whenever we increase a relation, the query answer also increases (or stays the same)
- Increase here means larger set

Monotone Queries

- A query with input relations R, S, T, \dots is called monotone if, whenever we increase a relation, the query answer also increases (or stays the same)
- Increase here means larger set
- Mathematically

If $R \subseteq R', S \subseteq S', \dots$ then $Q(R, S, \dots) \subseteq Q(R', S', \dots)$

Supplier(sno,sname,scity,sstate)

Supply(sno,pno,price)

Which Queries are Monotone?

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno = 2
```

Supplier(sno,sname,scity,sstate)

Supply(sno,pno,price)

Which Queries are Monotone?

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno = 2
```

MONOTONE

Supplier(sno,sname,scity,sstate)

Supply(sno,pno,price)

Which Queries are Monotone?

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno = 2
```

MONOTONE

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno != 2
```

Supplier(sno,sname,scity,sstate)

Supply(sno,pno,price)

Which Queries are Monotone?

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno = 2
```

MONOTONE

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno != 2
```

MONOTONE

Supplier(sno,sname,scity,sstate)

Supply(sno,pno,price)

Which Queries are Monotone?

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno = 2
```

MONOTONE

```
SELECT x.city, count(*)  
FROM Supplier x  
GROUP BY x.city
```

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno != 2
```

MONOTONE

Supplier(sno,sname,scity,sstate)

Supply(sno,pno,price)

Which Queries are Monotone?

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno = 2
```

MONOTONE

```
SELECT x.city, count(*)  
FROM Supplier x  
GROUP BY x.city
```

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno != 2
```

MONOTONE

NON-MONOTONE

Supplier(sno,sname,scity,sstate)

Supply(sno,pno,price)

Which Queries are Monotone?

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno = 2
```

MONOTONE

```
SELECT x.city, count(*)  
FROM Supplier x  
GROUP BY x.city
```

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno != 2
```

MONOTONE

NON-MONOTONE

```
SELECT x.sno, x.sname FROM Supplier x  
WHERE x.sno IN (SELECT y.sno  
                FROM Supply y  
                WHERE y.pno = 2 )
```


Supplier(sno,sname,scity,sstate)

Supply(sno,pno,price)

Which Queries are Monotone?

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno = 2
```

MONOTONE

```
SELECT x.city, count(*)  
FROM Supplier x  
GROUP BY x.city
```

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno != 2
```

MONOTONE

NON-MONOTONE

```
SELECT x.sno, x.sname FROM Supplier x  
WHERE x.sno IN (SELECT y.sno  
                FROM Supply y  
                WHERE y.pno = 2 )
```

MONOTONE

Supplier(sno,sname,scity,sstate)

Supply(sno,pno,price)

Which Queries are Monotone?

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno = 2
```

MONOTONE

```
SELECT x.city, count(*)  
FROM Supplier x  
GROUP BY x.city
```

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno != 2
```

MONOTONE

NON-MONOTONE

```
SELECT x.sno, x.sname FROM Supplier x  
WHERE x.sno IN (SELECT y.sno  
                FROM Supply y  
                WHERE y.pno = 2 )
```

MONOTONE

```
SELECT x.sno, x.sname FROM Supplier x  
WHERE x.sno NOT IN (SELECT y.sno  
                   FROM Supply y  
                   WHERE y.pno != 2 )
```

Supplier(sno,sname,scity,sstate)

Supply(sno,pno,price)

Which Queries are Monotone?

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno = 2
```

MONOTONE

```
SELECT x.city, count(*)  
FROM Supplier x  
GROUP BY x.city
```

```
SELECT DISTINCT x.sno, x.name  
FROM Supplier x, Supply y  
WHERE x.sno = y.sno and y.pno != 2
```

MONOTONE

NON-MONOTONE

```
SELECT x.sno, x.sname FROM Supplier x  
WHERE x.sno IN (SELECT y.sno  
                FROM Supply y  
                WHERE y.pno = 2 )
```

MONOTONE

```
SELECT x.sno, x.sname FROM Supplier x  
WHERE x.sno NOT IN (SELECT y.sno  
                   FROM Supply y  
                   WHERE y.pno != 2 )
```

NON-MONOTONE

Back to Datalog

Theorem The Naïve Algorithm terminates after a number of iteration that is at most a polynomial in the size of the database

Proof next

Naïve Evaluation Algorithm

The ICO is monotone (uses only σ, Π, \bowtie, U)

Naïve Evaluation Algorithm

The ICO is monotone (uses only $\sigma, \Pi, \bowtie, \cup$)

Fact: the IDBs increase: $IDB_t \subseteq IDB_{t+1}$

Proof: by induction

Naïve Evaluation Algorithm

The ICO is monotone (uses only $\sigma, \Pi, \bowtie, \cup$)

Fact: the IDBs increase: $IDB_t \subseteq IDB_{t+1}$

Proof: by induction $IDB_0 (= \emptyset) \subseteq IDB_1$

Naïve Evaluation Algorithm

The ICO is monotone (uses only $\sigma, \Pi, \bowtie, \cup$)

Fact: the IDBs increase: $IDB_t \subseteq IDB_{t+1}$

Proof: by induction $IDB_0 (= \emptyset) \subseteq IDB_1$

Assuming $IDB_t \subseteq IDB_{t+1}$ we have:
 $USPJ(IDB_t) \subseteq USPJ(IDB_{t+1})$

Naïve Evaluation Algorithm

The ICO is monotone (uses only $\sigma, \Pi, \bowtie, \cup$)

Fact: the IDBs increase: $IDB_t \subseteq IDB_{t+1}$

Proof: by induction $IDB_0 (= \emptyset) \subseteq IDB_1$

Assuming $IDB_t \subseteq IDB_{t+1}$ we have:

$$IDB_{t+1} = \text{USPJ}(IDB_t) \subseteq \text{USPJ}(IDB_{t+1}) = IDB_{t+2}$$

Naïve Evaluation Algorithm

Number of iterations of naïve algorithm:

$$\leq \sum_{R \in IDB} n^{\text{arity}(R)} = O(n^{\max(\text{arity})})$$

- n = number of distinct values in EDBs

Proof: there are $\leq \sum_{R \in IDB} n^{\text{arity}(R)}$
possible facts that can be inserted to IDBs

Take-Away

- Datalog always terminates in PTIME
- Only holds for monotone programs
- And cannot compute new values:
$$T(x, v) :- T(x, w), v=w+1$$

Outline

- Syntax
- Single rule
- Multiple rules
- Recursion
- Naïve algorithm
- Semi-naïve algorithm
- Non-monotone operations

Problem with the Naïve Algorithm

- Same facts are discovered repeatedly
- The semi-naïve algorithm reduces the number of rediscovered facts
- Uses incremental view maintenance

Incremental View Maintenance

- Materialized view: $V = Q(R, S, T, \dots)$
- A table gets updated: $R \leftarrow R \cup \Delta R$
- Want to update the view: $V \leftarrow V \cup \Delta V$

IVM: compute $\Delta V = \Delta Q(\Delta R, R, S, T, \dots)$

IVM

Example 1:

$V(x,y) :- R(x,z), S(z,y)$

If $R \leftarrow R \cup \Delta R$
then what is $\Delta V(x,y)$?

IVM

Example 1:

$V(x,y) :- R(x,z), S(z,y)$

If $R \leftarrow R \cup \Delta R$
then what is $\Delta V(x,y)$?

$\Delta V(x,y) :- \Delta R(x,z), S(z,y)$

IVM

Example 2:

$V(x,y) :- R(x,z), S(z,y)$

If $R \leftarrow R \cup \Delta R$ and $S \leftarrow S \cup \Delta S$
then what is $\Delta V(x,y)$?

IVM

Example 2:

$V(x,y) :- R(x,z), S(z,y)$

If $R \leftarrow R \cup \Delta R$ and $S \leftarrow S \cup \Delta S$
then what is $\Delta V(x,y)$?

$\Delta V(x,y) :- \Delta R(x,z), S(z,y)$
 $\Delta V(x,y) :- R(x,z), \Delta S(z,y)$
 $\Delta V(x,y) :- \Delta R(x,z), \Delta S(z,y)$

We use datalog convention
to represent the union
of three queries

IVM

Example 3:

$V(x,y) :- T(x,z), T(z,y)$

If $T \leftarrow T \cup \Delta T$
then what is $\Delta V(x,y)$?

IVM

Example 3:

$V(x,y) :- T(x,z), T(z,y)$

If $T \leftarrow T \cup \Delta T$
then what is $\Delta V(x,y)$?

$\Delta V(x,y) :- \Delta T(x,z), T(z,y)$

$\Delta V(x,y) :- T(x,z), \Delta T(z,y)$

$\Delta V(x,y) :- \Delta T(x,z), \Delta T(z,y)$

Semi-Naïve Evaluation Algorithm

Key idea:

- Use IVM in the inner loop of the naïve algorithm

Applies only to a monotone program

Semi-Naïve Evaluation Algorithm

```
IDB :=  $\emptyset$   
repeat  
  IDB := USPJ(IDB)  
until no more change
```

Naive

Semi-Naïve Evaluation Algorithm

```
IDB :=  $\emptyset$   
repeat  
  IDB := USPJ(IDB)  
until no more change
```

Naive

Semi-naive

```
IDB :=  $\Delta$  := NonRecursiveUSPJ  
repeat  
   $\Delta$  :=  $\Delta$ USPJ(IDB,  $\Delta$ ) - IDB  
  if  $\Delta = \emptyset$  then break  
  IDB := IDB  $\cup$   $\Delta$ 
```

Example

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

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

$T := \emptyset;$

repeat

$T := R \cup \Pi_{x,y}(R \bowtie T);$

until [no more change]

Naïve

Example

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

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

$T := \emptyset;$

repeat

$T := R \cup \Pi_{x,y}(R \bowtie T);$

until [no more change]

Naïve

Semi-naïve

$T := \Delta T := R;$

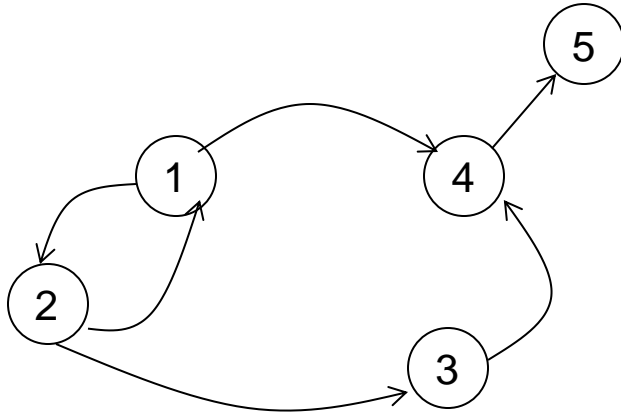
repeat

$\Delta T := \Pi_{x,y}(R \bowtie \Delta T) - T;$

if $\Delta T = \emptyset$ **then** break

$T := T \cup \Delta T;$

Example

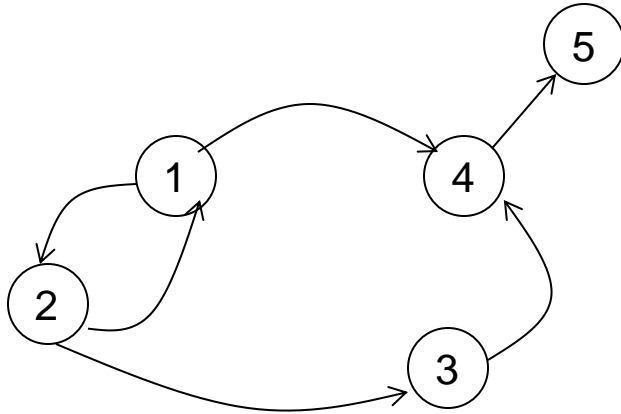


$T(x,y) \text{ :- } R(x,y)$
 $T(x,y) \text{ :- } R(x,z), T(z,y)$

R=

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

Example



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

$T := \Delta T := R;$

repeat

$\Delta T := \Pi_{x,y}(R \bowtie \Delta T) - T;$

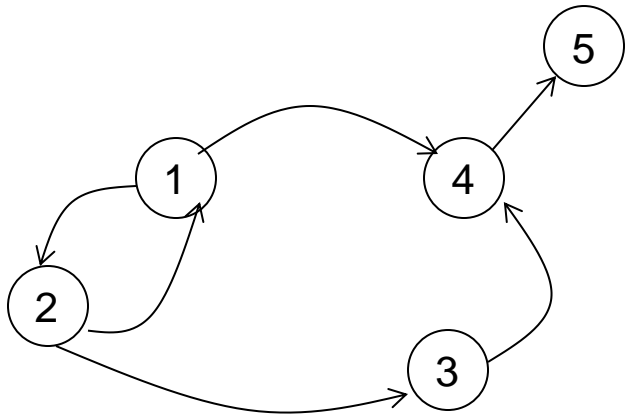
if $\Delta T = \emptyset$ **then break**

$T := T \cup \Delta T;$

R=

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

Example



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

```

T := ΔT := R;
repeat
  ΔT := Πx,y(R ⋈ ΔT) - T;
  if ΔT = ∅ then break
  T := T ∪ ΔT;
  
```

R=

Initially:

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

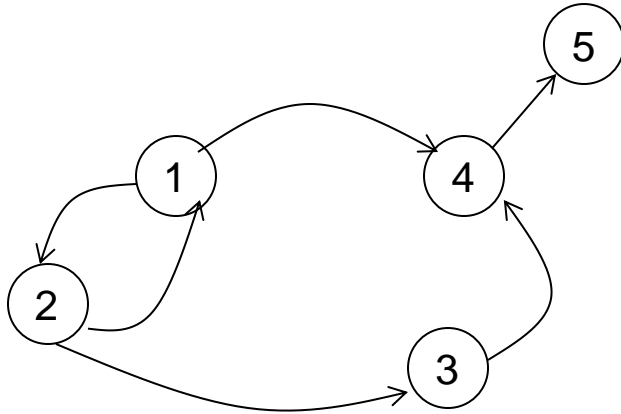
ΔT=

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

T=

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

Example



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

```

T := ΔT := R;
repeat
  ΔT := Πx,y(R ⋈ ΔT) - T;
  if ΔT = ∅ then break
  T := T ∪ ΔT;
  
```

First iteration:

R=

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

Initially:

ΔT=

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

T=

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

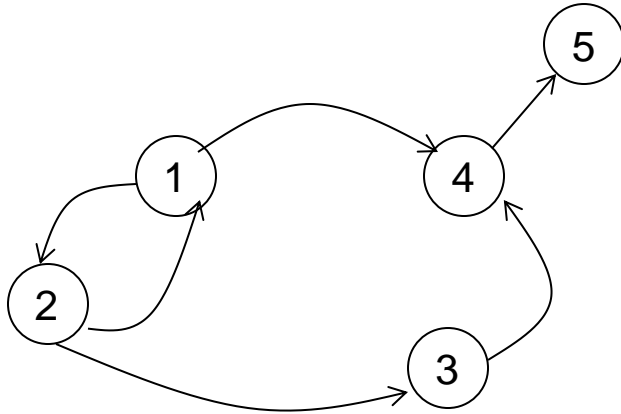
T=

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

ΔT= paths of length 2

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

Example



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

```

T := ΔT := R;
repeat
  ΔT := Πx,y(R ⋈ ΔT) - T;
  if ΔT = ∅ then break
  T := T ∪ ΔT;
  
```

First iteration:

Second iteration:

R=

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

Initially:

ΔT=

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

T=

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

T=

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

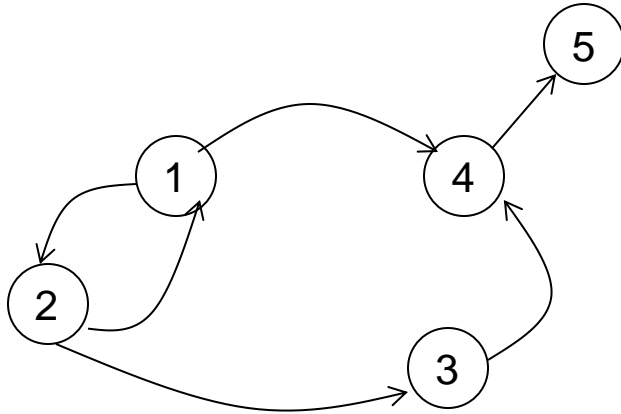
ΔT= paths of length 2

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

ΔT= paths of length 3

1	2
1	4
2	1
2	3
2	5

Example



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

```

T := ΔT := R;
repeat
  ΔT := Πx,y(R ⋈ ΔT) - T;
  if ΔT = ∅ then break
  T := T ∪ ΔT;
  
```

First iteration:

Second iteration:

R=

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

Initially:

ΔT=

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

T=

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

ΔT= paths of length 2

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

T=

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

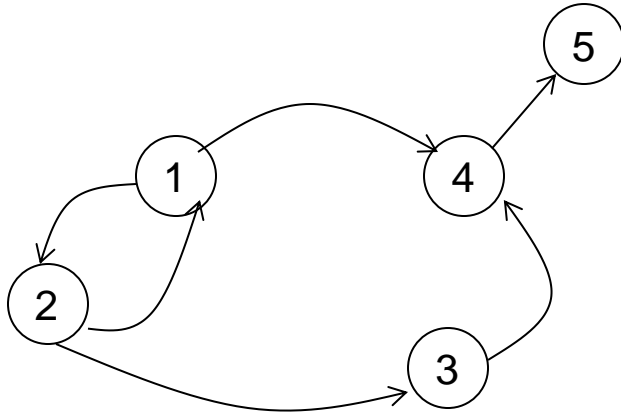
ΔT= paths of length 3

1	2
1	4
2	1
2	3
2	5

...this difference

Removed by the difference

Example



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

```

T := ΔT := R;
repeat
  ΔT := Πx,y(R ⋈ ΔT) - T;
  if ΔT = ∅ then break
  T := T ∪ ΔT;

```

First iteration:

Second iteration:

R=

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

Initially:

ΔT=

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

T=

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

T=

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

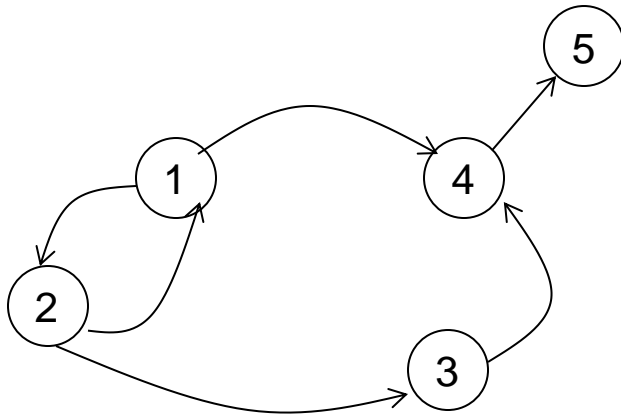
ΔT= paths of length 2

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

ΔT= paths of length 3

2	5
---	---

Example



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

$T := \Delta T := R;$

repeat

$\Delta T := \Pi_{x,y}(R \bowtie \Delta T) - T;$

if $\Delta T = \emptyset$ **then break**

$T := T \cup \Delta T;$

R=

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

Initially:

$\Delta T =$

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

T=

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

First iteration:

$\Delta T =$

paths of length 2

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

T=

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

Second iteration:

$\Delta T =$

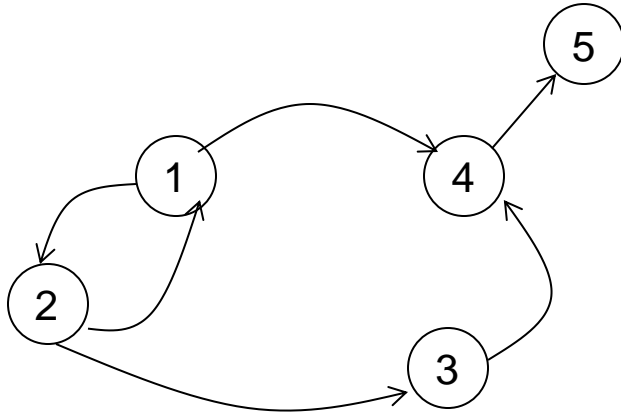
paths of length 3

2	5
---	---

T=

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

Example



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

```

T := ΔT := R;
repeat
  ΔT := Πx,y(R ⋈ ΔT) - T;
  if ΔT = ∅ then break
  T := T ∪ ΔT;
  
```

R=

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

Initially:

ΔT=

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

T=

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

First iteration:

ΔT= paths of length 2

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

T=

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

Second iteration:

ΔT= paths of length 3

2	5
---	---

T=

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

Third iteration:

ΔT= paths of length 4

--	--

Discussion

- Datalog systems, SQL run semi-naïve
- When rule is non-linear, IVM more complicated:
 $T(x,y) :- T(x,z), T(z,y)$
 becomes:
 $\Delta T(x,y) :- T(x,z), \Delta T(z,y)$
 $\Delta T(x,y) :- \Delta T(x,z), T(z,y)$
 $\Delta T(x,y) :- \Delta T(x,z), \Delta T(z,y)$
- SQL allows only linear recursion
- Postgres does not implement yet mutual recursion

Outline

- Syntax
- Single rule
- Multiple rules
- Recursion
- Naïve algorithm
- Semi-naïve algorithm
- Non-monotone operations

Non-monotone Extensions

- Aggregates

No standard syntax

We will follow Souffle

- Grouping

- Negation

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, year)

Aggregates

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

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

Meaning (in SQL)

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

Aggregates in Souffle:

- count
- min
- max
- sum

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, 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
```

Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, year)

Grouping

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

Here we bind
the year y

Meaning (in SQL)

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


Actor(id, fname, lname)

Casts(pid, mid)

Movie(id, title, year)

Grouping

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

Here we bind
the year y

We count
movies with
year = y

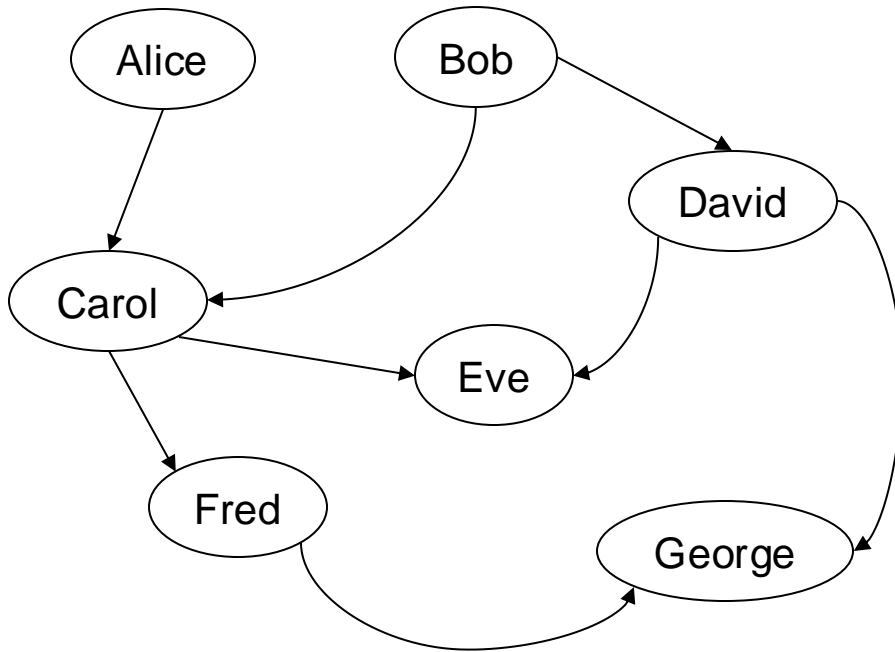
Meaning (in SQL)

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

Stratified datalog

- The (semi-)naïve algorithm does not work with non-monotone rules
- Datalog programs must be stratified:
 - Rules can be grouped into strata...
 - ...such that IDBs that occur in a body have been defined in earlier strata

Genealogy Database

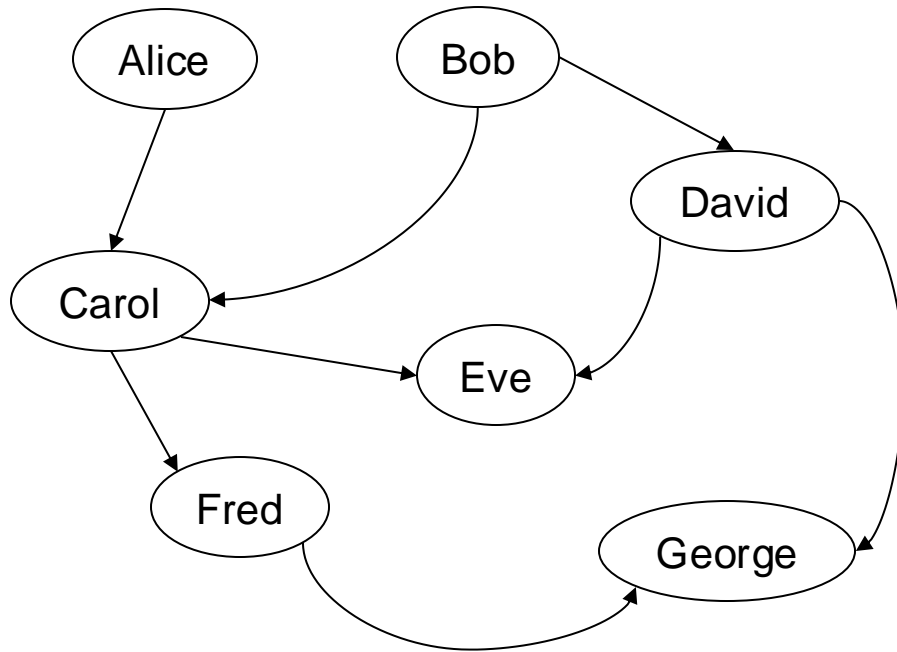


ParentChild

p	c
Alice	Carol
Bob	Carol
Bob	David
Carol	Eve
...	

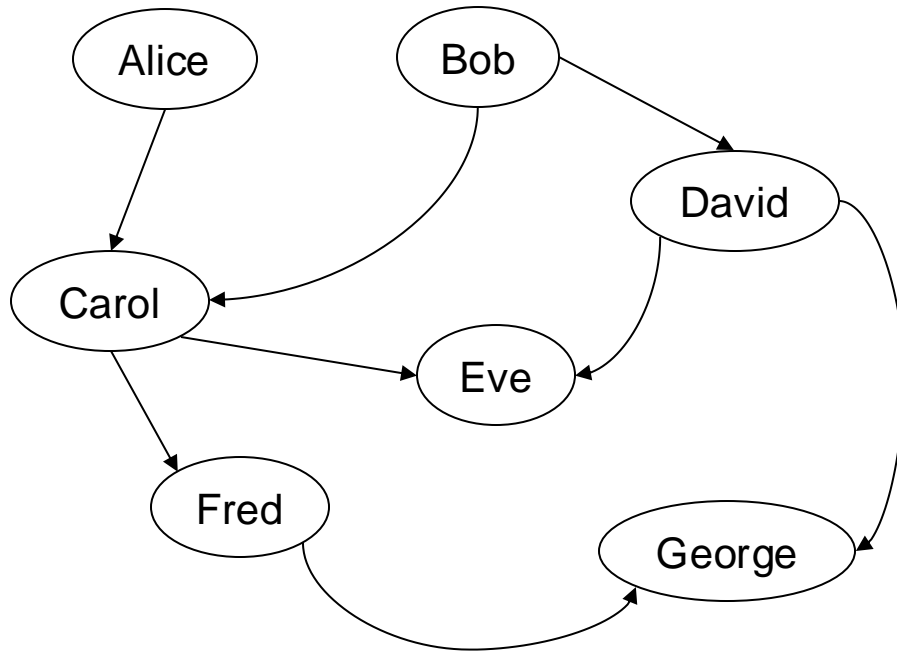
Count Descendants

For each person, count his/her descendants



Count Descendants

For each person, count his/her descendants

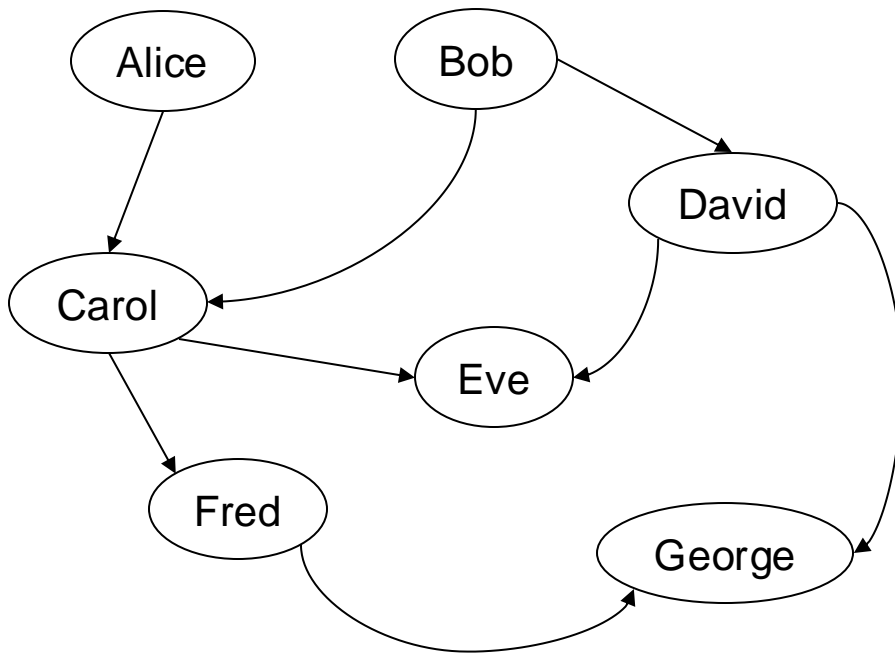


Answer

p	cnt
Alice	4
Bob	5
Carol	3
David	2
Fred	1

Count Descendants

For each person, count his/her descendants



Answer

p	cnt
Alice	4
Bob	5
Carol	3
David	2
Fred	1

Note: Eve and George do not appear in the answer (why?)

Count Descendants

Compute transitive closure of ParentChild

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

Count Descendants

Compute transitive closure of ParentChild

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


Count Descendants

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).
```

Count Descendants

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,_) }.
```

Count Descendants

How many descendants does Alice have?

```
// 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,_) }.
```

Count Descendants

How many descendants does Alice have?

```
// 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,_) }.

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

Count Descendants

How many descendants does Alice have?

Stratum 1

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

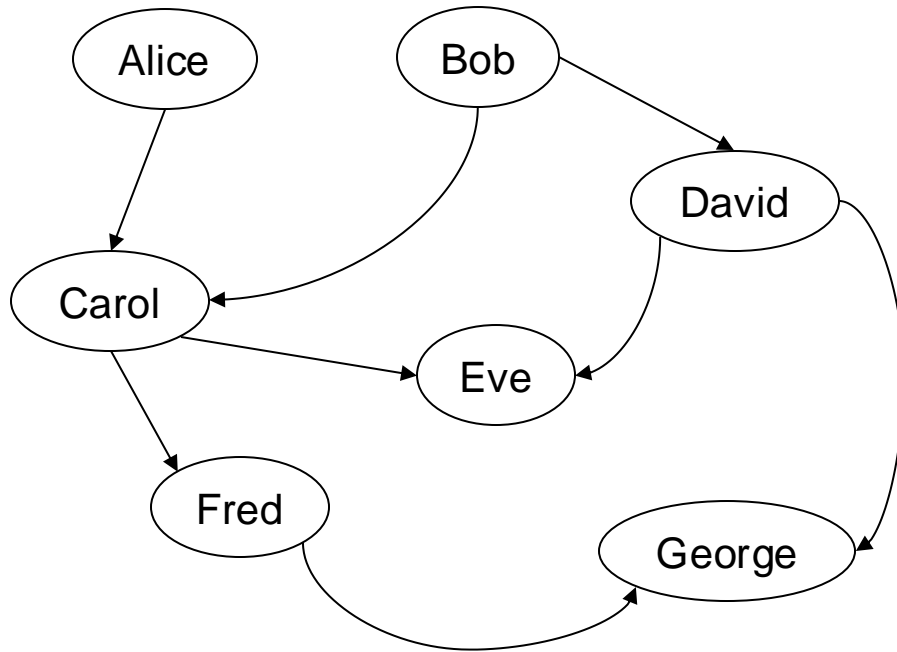
```
// For each person, count the number of descendants
T(p,c) :- D(p,_), c = count : { D(p,_) }.

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

Stratum 2

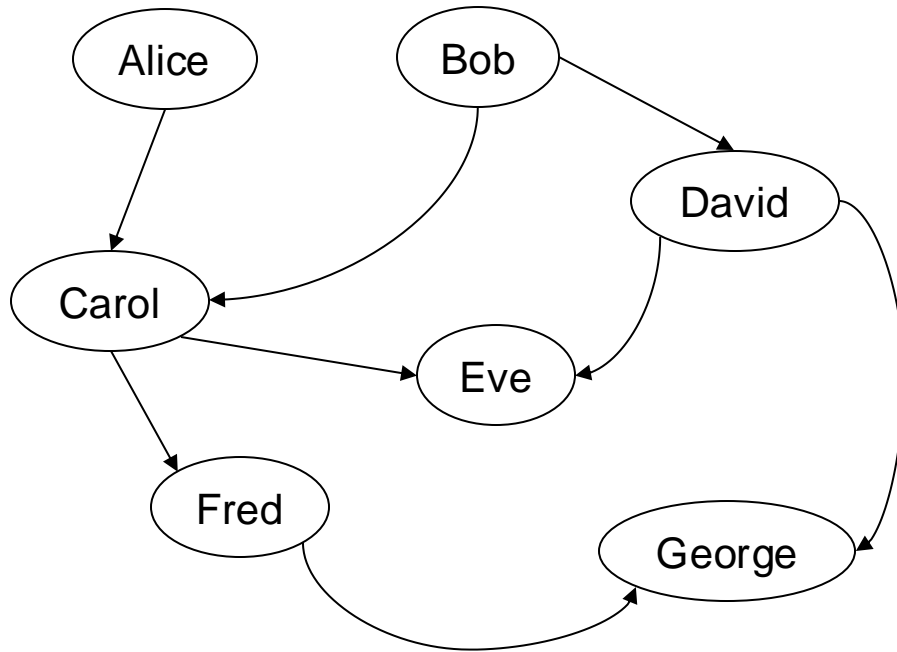
Negation: use “!”

Find all descendants of Bob that are not descendants of Alice



Negation: use “!”

Find all descendants of Bob that are not descendants of Alice



Answer

x
David

Negation: use “!”

Find all descendants of Bob that are not descendants of Alice

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


Negation: use “!”

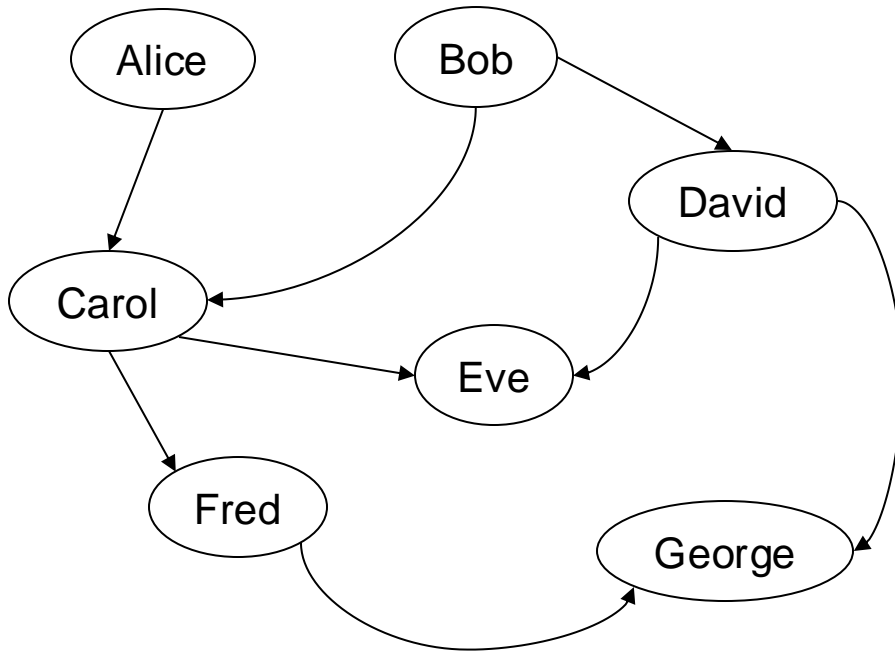
Find all descendants of Bob that are not descendants of Alice

```
// 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(“Bob”,x), !D(“Alice”,x).
```

Same Generation

Two people are in the same generation if they are descendants at the same generation of some common ancestor



SG

p1	p2
Carol	David
Eve	George
Fred	George
Fred	Eve

Same Generation

Compute pairs of people at the same generation

```
// common parent
```

Same Generation

Compute pairs of people at the same generation

```
// common parent  
SG(x,y) :- ParentChild(p,x), ParentChild(p,y)
```

Same Generation

Compute pairs of people at the same generation

```
// common parent
SG(x,y) :- ParentChild(p,x), ParentChild(p,y)

// parents at the same generation
```

Same Generation

Compute pairs of people at the same generation

```
// common parent
SG(x,y) :- ParentChild(p,x), ParentChild(p,y)

// parents at the same generation
SG(x,y) :- ParentChild(p,x), ParentChild(q,y), SG(p,q)
```

Same Generation

Compute pairs of people at the same generation

```
// common parent
SG(x,y) :- ParentChild(p,x), ParentChild(p,y)

// parents at the same generation
SG(x,y) :- ParentChild(p,x), ParentChild(q,y), SG(p,q)
```

Problem: this includes answers like SG(Carol, Carol)

And also SG(Eve, George), SG(George, Eve)

How to fix?

Same Generation

Compute pairs of people at the same generation

```
// common parent
SG(x,y) :- ParentChild(p,x), ParentChild(p,y), x < y

// parents at the same generation
SG(x,y) :- ParentChild(p,x), ParentChild(q,y),
           SG(p,q), x < y
```


Summary of Datalog

- Simple, concise syntax
- (Semi)-naïve algorithm for monotone programs; otherwise, need stratification
- Ad-hoc restrictions in SQL:
 - Linear recursion only
 - Postgres: no mutual recursion