

CSE544: Principles of Database Systems

Datalog

Announcements

- Review 7 (Datalog) due next Monday
- Project Milestone due next Friday
- HW3 is posted, due the following Monday

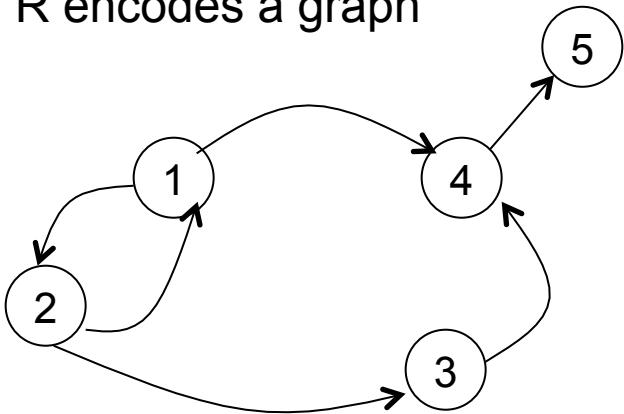
Datalog

Review the following basic concepts from Lecture 5:

- Fact
- Rule
- Head and body of a rule
- Existential variable
- Head variable

Simple datalog programs

R encodes a graph



R=

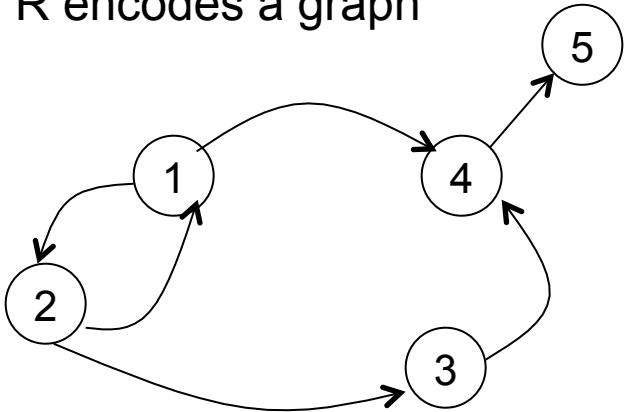
1	2
2	1
2	3
1	4
3	4
4	5

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

What does it compute?

Simple datalog programs

R encodes a graph



R=

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

Initially:
T is empty.

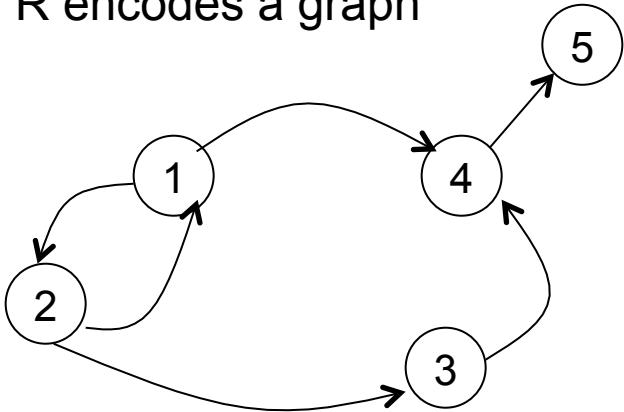


```
T(x,y) :- R(x,y)  
T(x,y) :- R(x,z), T(z,y)
```

What does
it compute?

Simple datalog programs

R encodes a graph



$R =$

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

Initially:
 T is empty.



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

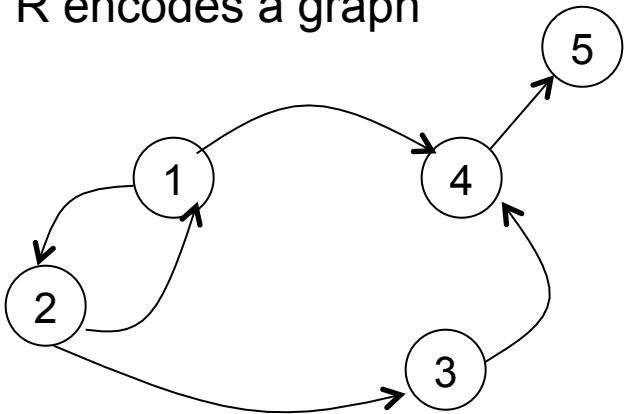
What does
it compute?

First iteration:
 $T =$

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

Simple datalog programs

R encodes a graph



$R =$

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

Initially:
 T is empty.



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

What does
it compute?

First iteration:
 $T =$

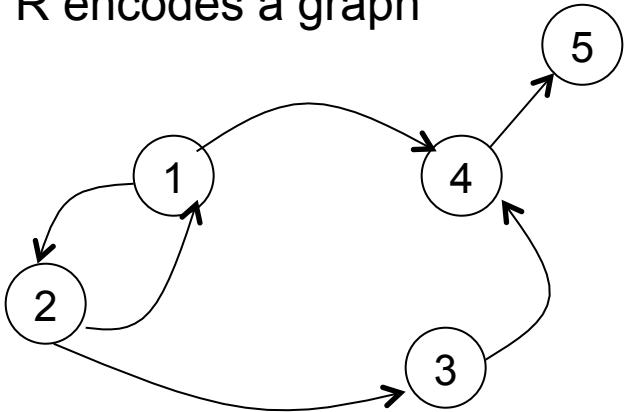
1	2
2	1
2	3
1	4
3	4

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

Simple datalog programs

R encodes a graph



R =

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

Initially:
T is empty.



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

What does
it compute?

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

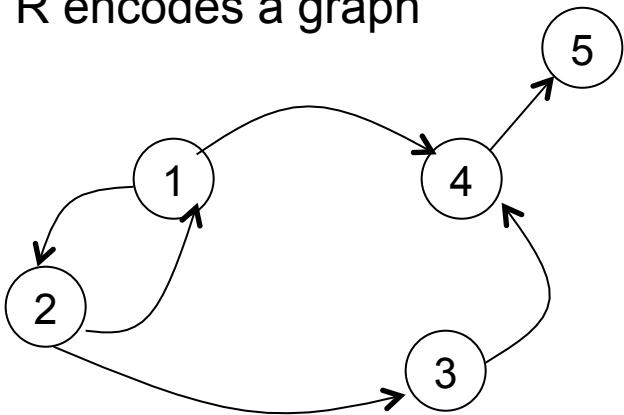
Third iteration:
T =

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

Done

Simple datalog programs

R encodes a graph



R =

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

Initially:
T is empty.



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

What does
it compute?

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

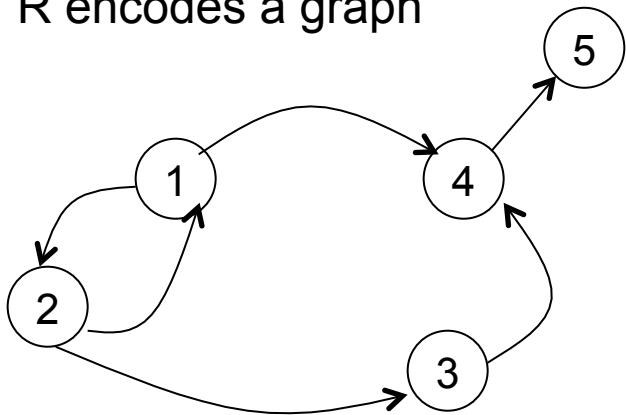
Discovered
3 times!

Discovered
twice

Done

Simple datalog programs

R encodes a graph



$R =$

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

Alternative ways to compute TC:

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

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

Right linear

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

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

Left linear

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

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

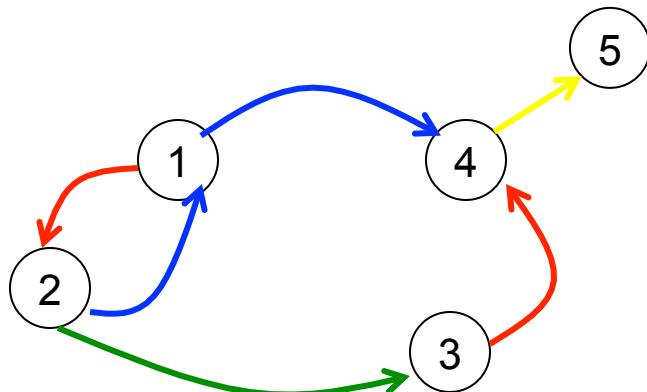
Non-linear

Discuss pros/cons in class

Simple datalog programs

Compute TC (ignoring color):

R encodes a **colored** graph



Compute pairs of nodes connected by the same color (e.g. (2,4))

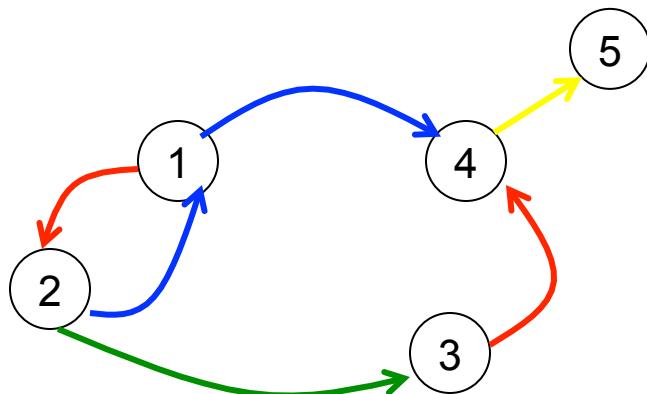
R=

1	Red	2
2	Blue	1
2	Green	3
1	Blue	4
3	Red	4
4	Yellow	5

Simple datalog programs

Compute TC (ignoring color):

R encodes a **colored** graph



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

Compute pairs of nodes connected by the same color (e.g. (2,4))

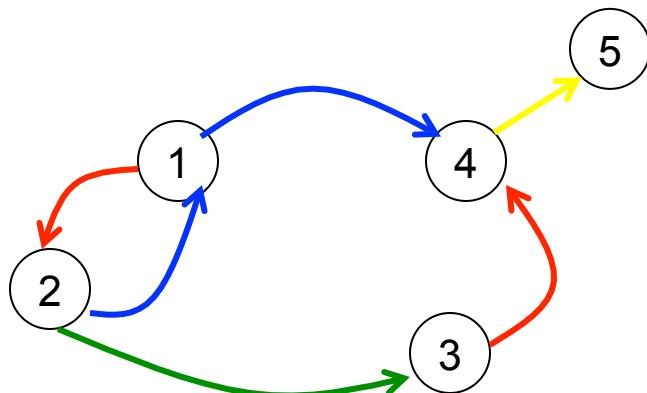
R=

1	Red	2
2	Blue	1
2	Green	3
1	Blue	4
3	Red	4
4	Yellow	5

Simple datalog programs

Compute TC (ignoring color):

R encodes a **colored** graph



R=

1	Red	2
2	Blue	1
2	Green	3
1	Blue	4
3	Red	4
4	Yellow	5

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

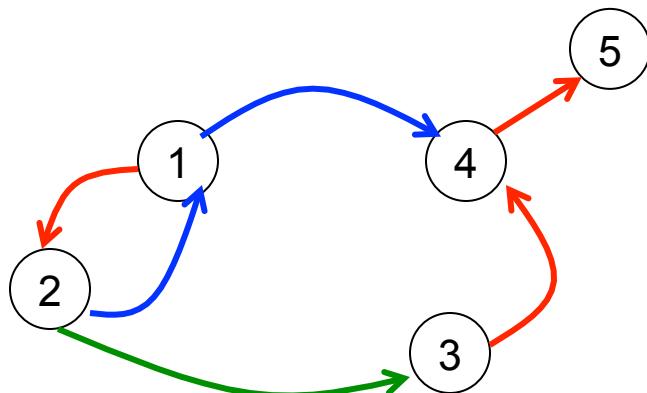
Compute pairs of nodes connected by the same color (e.g. (2,4))

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

Simple datalog programs

R, G, B encodes a 3-colored graph

What does this program compute in general?



R=

1	2
3	4
4	5

G=

2	3
---	---

B=

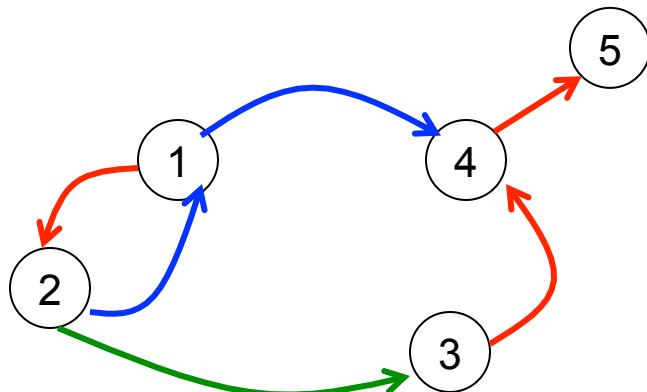
2	1
1	4

```
S(x,y) :- B(x,y)
S(x,y) :- T(x,z), B(z,y)
T(x,y) :- S(x,z), R(z,y)
T(x,y) :- S(x,z), G(z,y)
Answer(x,y) :- T(x,y)
```

Simple datalog programs

R, G, B encodes a 3-colored graph

What does this program compute in general?



R=

1	2
3	4
4	5

G=

2	3
---	---

B=

2	1
1	4

```
S(x,y) :- B(x,y)
S(x,y) :- T(x,z), B(z,y)
T(x,y) :- S(x,z), R(z,y)
T(x,y) :- S(x,z), G(z,y)
Answer(x,y) :- T(x,y)
```

Answer: it computes pairs of nodes connected by a path spelling out these regular expressions:

- $S = (B.(R \text{ or } G))^* \cdot B$
- $T = (B.(R \text{ or } G))^+$

Syntax of Datalog Programs

The schema consists of two sets of relations:

- Extensional Database (EDB): R_1, R_2, \dots
- Intentional Database (IDB): P_1, P_2, \dots

A datalog program P has the form:

$P:$
$$\begin{array}{l} P_{i1}(x_{11}, x_{12}, \dots) :- \text{body}_1 \\ P_{i2}(x_{21}, x_{22}, \dots) :- \text{body}_2 \\ \dots \end{array}$$

- Each head predicate P_i is an IDB
- Each body is a conjunction of IDB and/or EDB predicates
- See lecture 2

Note: no negation (yet)! Recursion OK.

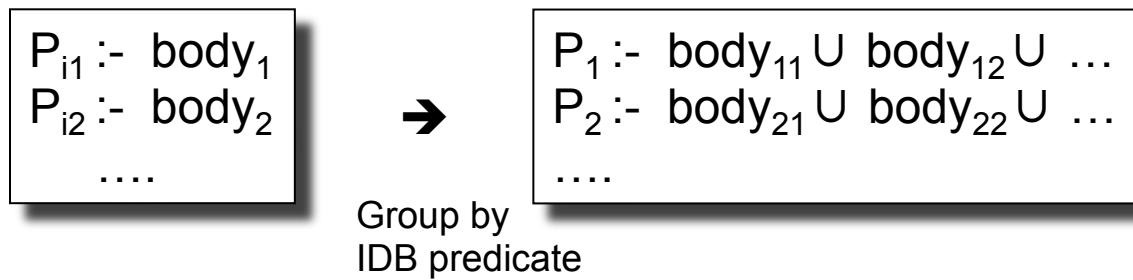
Naïve Datalog Evaluation Algorithm

Datalog program:

```
Pi1 :- body1
Pi2 :- body2
...
...
```

Naïve Datalog Evaluation Algorithm

Datalog program:



Naïve Datalog Evaluation Algorithm

Datalog program:

```
Pi1 :- body1  
Pi2 :- body2  
....
```



```
P1 :- body11 ∪ body12 ∪ ...  
P2 :- body21 ∪ body22 ∪ ...  
....
```

Group by
IDB predicate



```
P1 :- SPJU1  
P2 :- SPJU2  
....
```

Each rule is a
Select-Project-Join-Union query

Naïve Datalog Evaluation Algorithm

Datalog program:

```
Pi1 :- body1  
Pi2 :- body2  
....
```



```
P1 :- body11 ∪ body12 ∪ ...  
P2 :- body21 ∪ body22 ∪ ...  
....
```

Group by
IDB predicate



```
P1 :- SPJU1  
P2 :- SPJU2  
....
```

Each rule is a
Select-Project-Join-Union query

Example:

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



?

Naïve Datalog Evaluation Algorithm

Datalog program:

```
Pi1 :- body1  
Pi2 :- body2  
....
```



```
P1 :- body11 ∪ body12 ∪ ...  
P2 :- body21 ∪ body22 ∪ ...  
....
```

Group by
IDB predicate



```
P1 :- SPJU1  
P2 :- SPJU2  
....
```

Each rule is a
Select-Project-Join-Union query

Example:

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



```
T(x,y) :- R(x,y) ∪ Πxy(R(x,z) ⋙ T(z,y))
```

Naïve Datalog Evaluation Algorithm

Datalog program:

```
Pi1 :- body1  
Pi2 :- body2  
....
```



```
P1 :- body11 ∪ body12 ∪ ...  
P2 :- body21 ∪ body22 ∪ ...  
....
```

Group by
IDB predicate



```
P1 :- SPJU1  
P2 :- SPJU2  
....
```

Each rule is a
Select-Project-Join-Union query

Naïve datalog evaluation algorithm:

```
P1 = P2 = ... = ∅
```

Loop

```
NewP1 = SPJU1; NewP2 = SPJU2; ...  
if (NewP1 = P1 and NewP2 = P2 and ...) then exit
```

```
P1 = NewP1; P2 = NewP2; ...
```

Endloop

Example:

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



```
T(x,y) :- R(x,y) ∪ Πxy(R(x,z) ▷ T(z,y))
```

Naïve Datalog Evaluation Algorithm

Datalog program:

```
Pi1 :- body1  
Pi2 :- body2  
....
```



```
P1 :- body11 ∪ body12 ∪ ...  
P2 :- body21 ∪ body22 ∪ ...  
....
```

Group by
IDB predicate



```
P1 :- SPJU1  
P2 :- SPJU2  
....
```

Each rule is a
Select-Project-Join-Union query

Naïve datalog evaluation algorithm:

```
P1 = P2 = ... = ∅
```

Loop

```
NewP1 = SPJU1; NewP2 = SPJU2; ...  
if (NewP1 = P1 and NewP2 = P2 and ...) then exit
```

```
P1 = NewP1; P2 = NewP2; ...
```

Endloop

Example:

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



```
T(x,y) :- R(x,y) ∪ Πxy(R(x,z) ▷ T(z,y))
```

T = ∅

Loop

```
NewT(x,y) = R(x,y) ∪ Πxy(R(x,z) ▷ T(z,y))
```

```
if (NewT = T)  
then exit
```

```
T = NewT
```

Endloop

Discussion

- A datalog program always terminates (why?)
- What is the running time of a datalog program as a function of the input database?

Problem with the Naïve Algorithm

- The same facts are discovered over and over again
- The semi-naïve algorithm tries to reduce the number of facts discovered multiple times

Background: Incremental View Maintenance

Let V be a view computed by one datalog rule (no recursion)

$V :- \text{body}$

If (some of) the relations are updated: $R_1 \leftarrow R_1 \cup \Delta R_1, R_2 \leftarrow R_2 \cup \Delta R_2, \dots$

Then the view is also modified as follows: $V \leftarrow V \cup \Delta V$

Incremental view maintenance:

Compute ΔV without having to recompute V

Background: Incremental View Maintenance

Example 1:

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

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

Background: Incremental View Maintenance

Example 1:

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

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

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

Background: Incremental View Maintenance

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 $\Delta V(x,y)$?

Background: Incremental View Maintenance

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 $\Delta V(x,y)$?

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

Background: Incremental View Maintenance

Example 3:

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

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

Background: Incremental View Maintenance

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

Separate the Datalog program into the non-recursive, and the recursive part.
Each P_i defined by non-recursive-SPJU $_i$ and (recursive-)SPJU $_i$.

$P_1 = \Delta P_1 = \text{non-recursive-SPJU}_1, P_2 = \Delta P_2 = \text{non-recursive-SPJU}_2, \dots$

Loop

$\Delta P_1 = \Delta \text{SPJU}_1; \Delta P_2 = \Delta \text{SPJU}_2; \dots$

if ($\Delta P_1 = \emptyset$ and $\Delta P_2 = \emptyset$ and ...)

 then break

$P_1 = P_1 \cup \Delta P_1; P_2 = P_2 \cup \Delta P_2; \dots$

Endloop

Semi-naïve Evaluation Algorithm

Separate the Datalog program into the non-recursive, and the recursive part.
Each P_i defined by non-recursive-SPJU $_i$ and (recursive-)SPJU $_i$.

$P_1 = \Delta P_1 = \text{non-recursive-SPJU}_1, P_2 = \Delta P_2 = \text{non-recursive-SPJU}_2, \dots$

Loop

$\Delta P_1 = \Delta \text{SPJU}_1; \Delta P_2 = \Delta \text{SPJU}_2; \dots$

if ($\Delta P_1 = \emptyset$ and $\Delta P_2 = \emptyset$ and ...)

then break

$P_1 = P_1 \cup \Delta P_1; P_2 = P_2 \cup \Delta P_2; \dots$

Endloop

Example:

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

$T = \Delta T = ?$ (non-recursive rule)

Loop

$\Delta T(x,y) = ?$ (recursive Δ -rule)

if ($\Delta T = \emptyset$)

then break

$T = T \cup \Delta T$

Endloop

Semi-naïve Evaluation Algorithm

Separate the Datalog program into the non-recursive, and the recursive part.
Each P_i defined by non-recursive-SPJU $_i$ and (recursive-)SPJU $_i$.

$P_1 = \Delta P_1 = \text{non-recursive-SPJU}_1, P_2 = \Delta P_2 = \text{non-recursive-SPJU}_2, \dots$

Loop

$\Delta P_1 = \Delta \text{SPJU}_1; \Delta P_2 = \Delta \text{SPJU}_2; \dots$

if ($\Delta P_1 = \emptyset$ and $\Delta P_2 = \emptyset$ and ...)

then break

$P_1 = P_1 \cup \Delta P_1; P_2 = P_2 \cup \Delta P_2; \dots$

Endloop

Example:

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

$T(x,y) = R(x,y), \Delta T(x,y) = R(x,y)$

Loop

$\Delta T(x,y) = R(x,z), \Delta T(z,y)$

if ($\Delta T = \emptyset$)

then break

$T = T \cup \Delta T$

Endloop

Semi-naïve Evaluation Algorithm

Separate the Datalog program into the non-recursive, and the recursive part.
Each P_i defined by non-recursive-SPJU $_i$ and (recursive-)SPJU $_i$.

$P_1 = \Delta P_1 = \text{non-recursive-SPJU}_1, P_2 = \Delta P_2 = \text{non-recursive-SPJU}_2, \dots$

Loop

$\Delta P_1 = \Delta \text{SPJU}_1; \Delta P_2 = \Delta \text{SPJU}_2; \dots$

if ($\Delta P_1 = \emptyset$ and $\Delta P_2 = \emptyset$ and ...)

then break

$P_1 = P_1 \cup \Delta P_1; P_2 = P_2 \cup \Delta P_2; \dots$

Endloop

Example:

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

$T(x,y) = R(x,y), \Delta T(x,y) = R(x,y)$

Loop

$\Delta T(x,y) = R(x,z), \Delta T(z,y)$

if ($\Delta T = \emptyset$)

then break

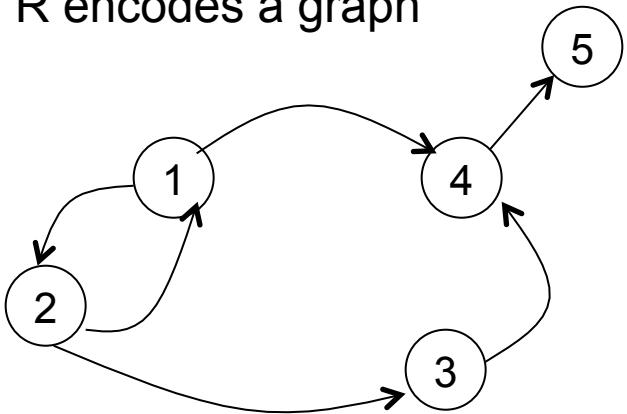
$T = T \cup \Delta T$

Endloop

Note: for any linear datalog programs,
the semi-naïve algorithm has only
one Δ -rule for each rule!

Simple datalog programs

R encodes a graph



R=

Initially:

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

$\Delta T =$

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

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

T= R, $\Delta T = R$

Loop

$\Delta T(x,y) = R(x,z), \Delta T(z,y)$

if ($\Delta T = \emptyset$)

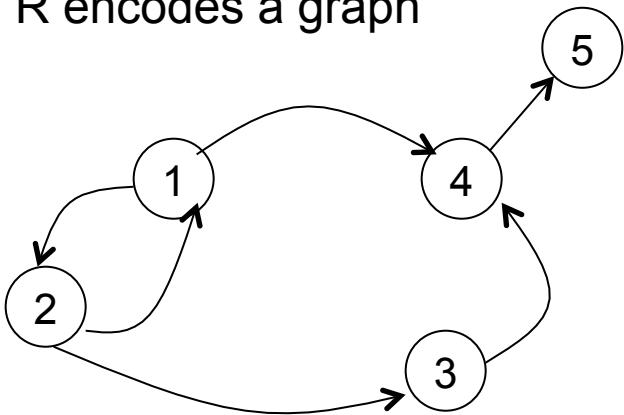
then break

T = T \cup ΔT

Endloop

Simple datalog programs

R encodes a graph



$R =$

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

Initially:

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

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

```

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

First iteration:

$T =$

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

$\Delta T =$
paths of
length 2

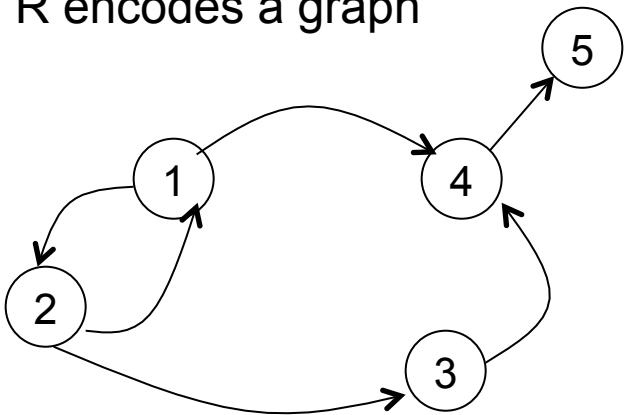
1	1
1	3
1	5
2	2
2	4
3	5

```

T = R, ΔT = R
Loop
ΔT(x,y) = R(x,z), ΔT(z,y)
if (ΔT = ∅)
  then break
T = T ∪ ΔT
Endloop
  
```

Simple datalog programs

R encodes a graph



$R =$

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

Initially:

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

First iteration:

$T =$

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

$\Delta T =$
paths of
length 2

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

Second iteration:

$T =$

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

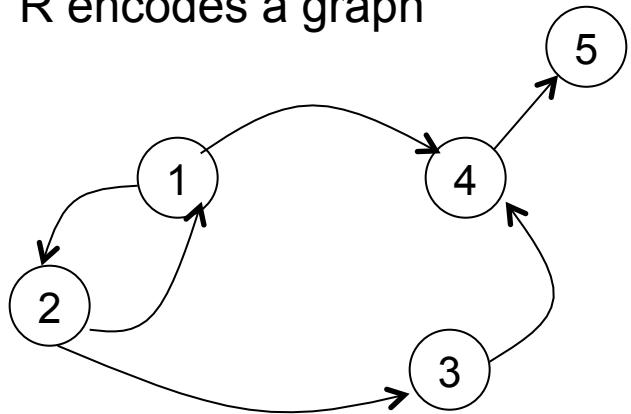
$\Delta T =$
paths of
length 3

```

T= R, ΔT = R
Loop
ΔT(x,y) = R(x,z), ΔT(z,y)
if (ΔT = ∅)
    then break
T = T ∪ ΔT
Endloop
    
```

Simple datalog programs

R encodes a graph



$R =$

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

Initially:

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

First iteration:

$T =$

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

$\Delta T =$
paths of
length 2

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

Second iteration:

$T =$

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

$\Delta T =$
paths of
length 3

1	2
1	4
2	1
2	3
3	5

Third iteration:

$\Delta T =$
paths of
length 4

$T = R, \Delta T = R$

Loop

$\Delta T(x,y) = R(x,z), \Delta T(z,y)$

if ($\Delta T = \emptyset$)

then break

$T = T \cup \Delta T$

Endloop

Discussion of Semi-Naïve Algorithm

- Avoids re-computing some tuples, but not all tuples
- Easy to implement, no disadvantage over naïve
- A rule is called *linear* if its body contains only one recursive IDB predicate:
 - A linear rule always results in a single incremental rule
 - A non-linear rule may result in multiple incremental rules

Summary of Datalog

- Simple syntax for expressing queries with recursion
- Bottom-up evaluation – always terminates
 - Naïve evaluation
 - Semi-naïve evaluation
- Next lecture:
 - Discuss the paper on datalog
 - Datalog semantics
 - Datalog with negation