

Topic 1: Data models and query languages

Unit 4: Datalog (continued)

Lecture 09

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp22)

<https://northeastern-datalab.github.io/cs7240/sp22/>

2/15/2022

Pre-class conversations

- Last class Datalog recapitulation
- Calendar: two more classes on data models & query languages
- Reminder: It is time to start to hand in your first scribes
- today:
 - Recursion (Datalog)

Syntax of rules

- evaluates to true when relation R_i contains the tuple described by args_i
 - e.g. $\text{Actor}(344759, \text{'Douglas'}, \text{'Fowley'})$ is true
- $R_i(\text{args}_i)$: relational predicate with arguments ("atom")
- arithmetic predicate

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

head

(or consequent)
single IDB subgoal

body

(or antecedent)
conjunction of subgoals

{f,l}: head variables
{u,x,y,z}: existential variables

Alternative notation: $Q(\text{args}) \leftarrow R_1(\text{args}) \text{ AND } R_2(\text{args}) \dots$ / or variables begin with a capital, predicates with lower-case (problem: can't have "Boston")

Source: Dan Suciu, CSE 554, 2018.

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>

Logical interpretation of a single rule



Actor(id, fname, lname)
Plays(aid, mid)
Movie(id, name, year)

$Q(y) :- \text{Movie}(x,y,z), z < '1940'.$

Meaning of a datalog rule is a logical statement:

?

Logical interpretation of a single rule



Actor(id, fname, lname)
Plays(aid, mid)
Movie(id, name, year)

$Q(y) :- \text{Movie}(x,y,z), z < '1940'.$

Meaning of a datalog rule is a logical statement:

For all x, y, z : if $(x, y, z) \in \text{Movies}$ and $z < '1940'$ then y is in Q (i.e. is part of the answer)

$\forall x, y, z [(\text{Movie}(x, y, z) \wedge z < 1940) \Rightarrow Q(y)]$

logically equivalent to

?

Logical interpretation of a single rule



Actor(id, fname, lname)
Plays(aid, mid)
Movie(id, name, year)

$Q(y) :- \text{Movie}(x,y,z), z < '1940'.$

Meaning of a datalog rule is a logical statement:

For all x, y, z : if $(x, y, z) \in \text{Movies}$ and $z < '1940'$ then y is in Q (i.e. is part of the answer)

$\forall x, y, z [(\text{Movie}(x, y, z) \wedge z < 1940) \Rightarrow Q(y)]$

logically equivalent to

$\forall y [\exists x, z [\text{Movie}(x, y, z) \wedge z < 1940] \Rightarrow Q(y)]$

Thus, non-head variables are called "existential variables"

compare with RC

?

Logical interpretation of a single rule



Actor(id, fname, lname)
Plays(aid, mid)
Movie(id, name, year)

$Q(y) :- \text{Movie}(x,y,z), z < '1940'.$

Meaning of a datalog rule is a logical statement:

For all x, y, z : if $(x, y, z) \in \text{Movies}$ and $z < '1940'$ then y is in Q (i.e. is part of the answer)

$$\forall x, y, z [(\text{Movie}(x, y, z) \wedge z < 1940) \Rightarrow Q(y)]$$

logically equivalent to

$$\forall y [\exists x, z [\text{Movie}(x, y, z) \wedge z < 1940] \Rightarrow Q(y)]$$

Thus, non-head variables are called "existential variables"

compare with RC

$$\{(y) \mid \exists x, z [\text{Movie}(x, y, z) \wedge z < 1940]\}$$

We want the smallest set Q with this property (why?)



Syntactic Constraints

$$Q(\mathbf{x}) :- R_1(\mathbf{x}_1, \mathbf{y}_1), \dots, R_m(\mathbf{x}_m, \mathbf{y}_m).$$

The rule stands for the following logical formula:

$$\forall \mathbf{x} [Q(\mathbf{x}) \leftarrow \exists \mathbf{y} [R_1(\mathbf{x}_1, \mathbf{y}_1) \wedge \dots \wedge R_m(\mathbf{x}_m, \mathbf{y}_m)]]$$

Two restrictions:

1. **Safety**: every head variable should occur in the body at least once

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

?

$$\mathbf{x}_i \subseteq \mathbf{x}, \mathbf{y}_i \subseteq \mathbf{y}$$

(bold = vector notation)

Recall we want the smallest set Q with this property



Syntactic Constraints

$$Q(\mathbf{x}) :- R_1(\mathbf{x}_1, \mathbf{y}_1), \dots, R_m(\mathbf{x}_m, \mathbf{y}_m).$$

$$\mathbf{x}_i \subseteq \mathbf{x}, \mathbf{y}_i \subseteq \mathbf{y}$$

(bold = vector notation)

The rule stands for the following logical formula:

$$\forall \mathbf{x} [Q(\mathbf{x}) \leftarrow \exists \mathbf{y} [R_1(\mathbf{x}_1, \mathbf{y}_1) \wedge \dots \wedge R_m(\mathbf{x}_m, \mathbf{y}_m)]]$$

Two restrictions:

1. **Safety**: every head variable should occur in the body at least once

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

forbidden rule: z not in body

2. The head predicate must be an **IDB (Intensional)** predicate
(Body can include both EDBs and IDBs)

$$\text{Arc}(x, y) :- \text{Arc}(x, z), \text{Arc}(z, y).$$

?



Syntactic Constraints

$$Q(\mathbf{x}) :- R_1(\mathbf{x}_1, \mathbf{y}_1), \dots, R_m(\mathbf{x}_m, \mathbf{y}_m).$$
$$\mathbf{x}_i \subseteq \mathbf{x}, \mathbf{y}_i \subseteq \mathbf{y}$$

(bold = vector notation)

The rule stands for the following logical formula:

$$\forall \mathbf{x} [Q(\mathbf{x}) \leftarrow \exists \mathbf{y} [R_1(\mathbf{x}_1, \mathbf{y}_1) \wedge \dots \wedge R_m(\mathbf{x}_m, \mathbf{y}_m)]]$$

Two restrictions:

1. **Safety**: every head variable should occur in the body at least once

~~$R(x, z) :- S(x, y), R(y, x).$~~

forbidden rule: z not in body

2. The head predicate must be an **IDB (Intensional)** predicate
(Body can include both EDBs and IDBs)

~~$\text{Arc}(x, y) :- \text{Arc}(x, z), \text{Arc}(z, y).$~~

assuming Arc is EDB

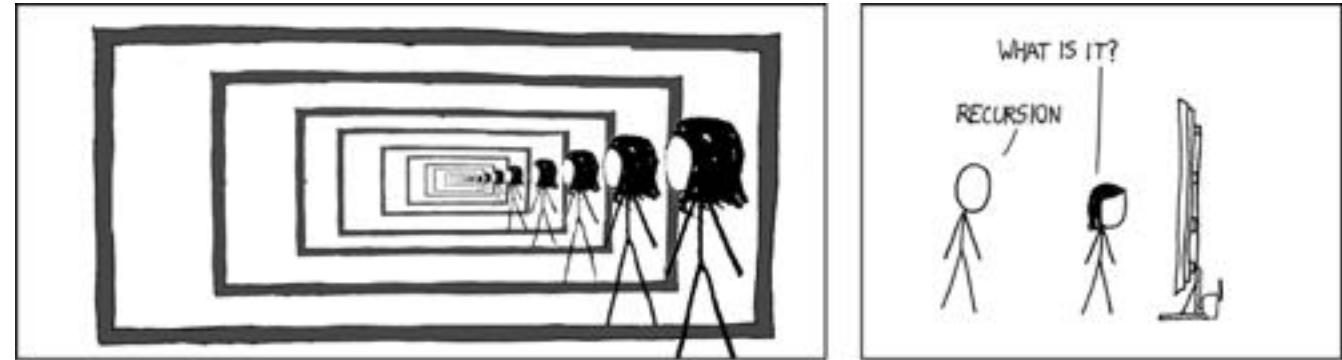
Outline: T1-4: Datalog

- Datalog
 - Datalog rules
 - Recursion
 - Semantics
 - Datalog \neg : Negation, stratification
 - Datalog \pm
 - Stable model semantics (Answer set programming)
 - Datalog vs. RA
 - Naive and Semi-naive evaluation (incl. Incremental View Maintenance)

Datalog program

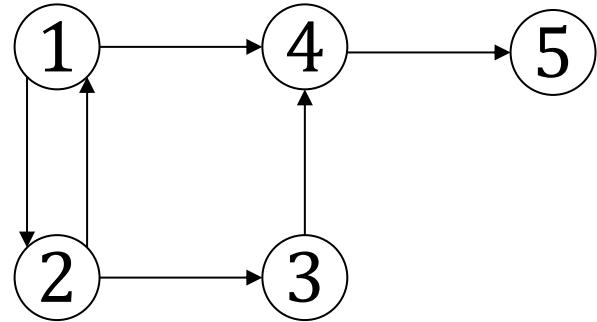
- A datalog program consists of several rules
- Usually there is one **distinguished predicate** that's the output
- Importantly, rules may be **recursive!**

Recursion occurs when
a thing is defined
in terms of itself
(self-repetition).



- (We later also see how to use recursion in SQL)

Example



EDB
IDB

$P(x,y) :- A(x,y).$
 $P(x,y) :- A(x,z), P(z,y).$

recursion due to
head in rule body

A(S,T)



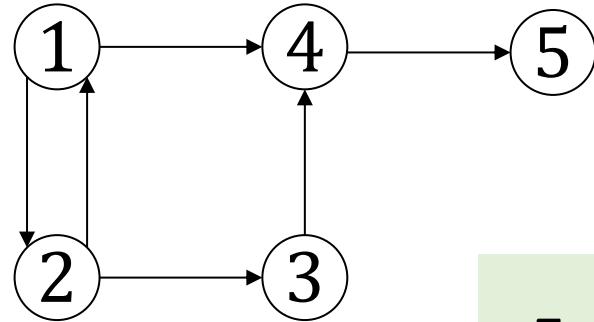
What does this query compute?

?

A

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

Example



EDB
IDB

$P(x,y) :- A(x,y).$
 $P(x,y) :- A(x,z), P(z,y).$

recursion due to
head in rule body

$x \rightarrow z \rightsquigarrow y$

Calculates all paths (transitive closure)

For all nodes x and y :

If there is an **arc** from x to y ,
then there is a **path** from x to y .

A

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

For all nodes x , z , and y :

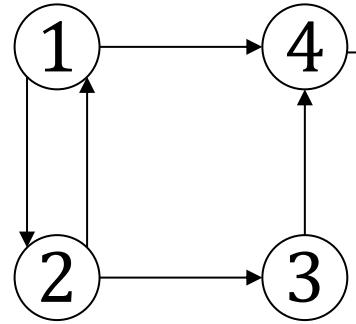
If there is an **arc** from x to z , and there is a **path** from z to y
then there is a **path** from x to y .

A(S,T)



Example

A(S,T)



EDB
IDB

$P(x,y) :- A(x,y).$
 $P(x,y) :- A(x,z), P(z,y).$

recursion due to
head in rule body

Initially: P is empty

1st iteration

P

?

A

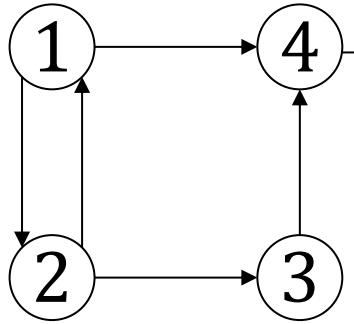
1	2
2	1
2	3
1	4
3	4
4	5

Example



A(S,T)

501



EDB
IDB

$P(x,y) :- A(x,y).$
 $P(x,y) :- A(x,z), P(z,y).$

recursion due to
head in rule body

Initially: P is empty

1st iteration

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

2nd iteration

P

?

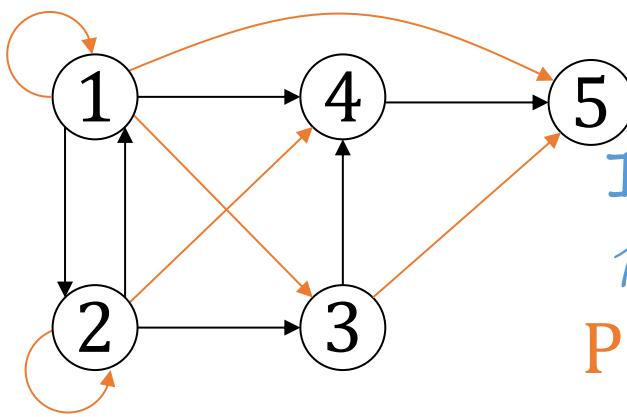
$P = A$ from
1st rule

A

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

2nd rule generates
nothing (because
 P is empty)

Example



EDB
IDB

$P(x,y) :- A(x,y).$
$P(x,y) :- A(x,z), P(z,y).$

recursion due to
head in rule body

Initially: P is empty

1st iteration

P

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

$P = A$ from
1st rule

2nd iteration

P

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

1st rule

A

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

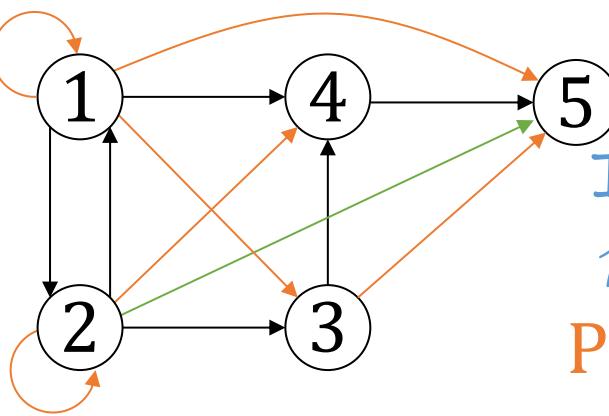
2nd rule generates
nothing (because
 P is empty)

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

2nd rule

New facts from 2nd rule

Example



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

EDB
IDB

P(x,y) :- A(x,y).
P(x,y) :- A(x,z), P(z,y).

recursion due to
head in rule body

Initially: P is empty

1st iteration

P

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

P=A from
1st rule

2nd iteration

P

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

recall set semantics!

(NO new facts)

3rd iteration = 4th iteration

P

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

1st + 2nd rule

1st rule

2nd rule generates
nothing (because
P is empty)

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

2nd rule

1	1
2	2
1	3
2	4
1	5

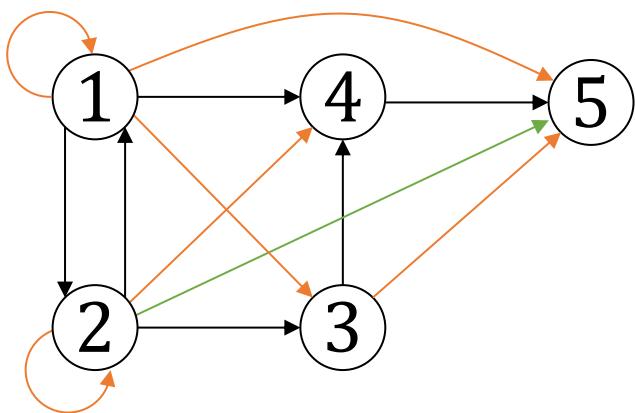
2nd rule

1	1
2	2
1	3
2	4
3	5

2	5
---	---

New facts from 2nd rule

Example



EDB
IDB

$P(x,y) :- A(x,y).$
 $P(x,y) :- A(x,z), P(z,y).$

A(S,T)



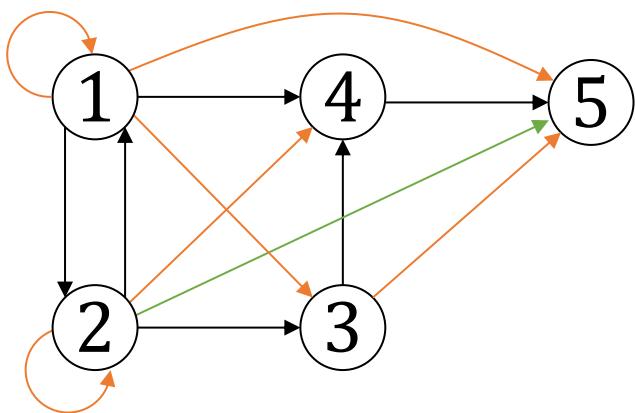
	S	T
	2	1
	2	3
	1	4
	3	4
	4	5

In SQL ?

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

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

Example



EDB
IDB

$P(x,y) :- A(x,y).$
 $P(x,y) :- A(x,z), P(z,y).$

A(S,T)



	S	T
	2	1
	2	3
	1	4
	3	4
	4	5

WITH RECURSIVE P AS (

?

UNION

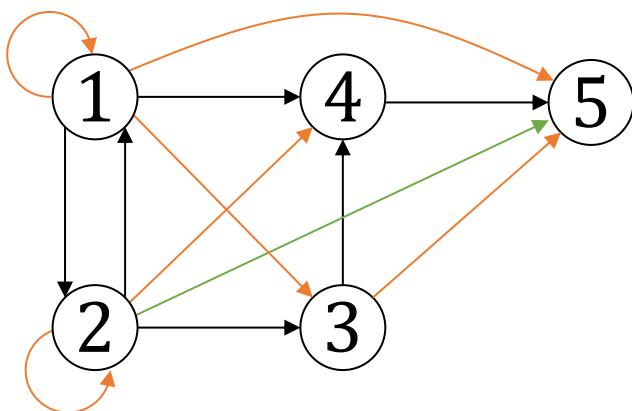
?

SELECT *

FROM P

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

Example



EDB
IDB

$P(x,y) :- A(x,y).$
 $P(x,y) :- A(x,z), P(z,y).$

A(S,T)

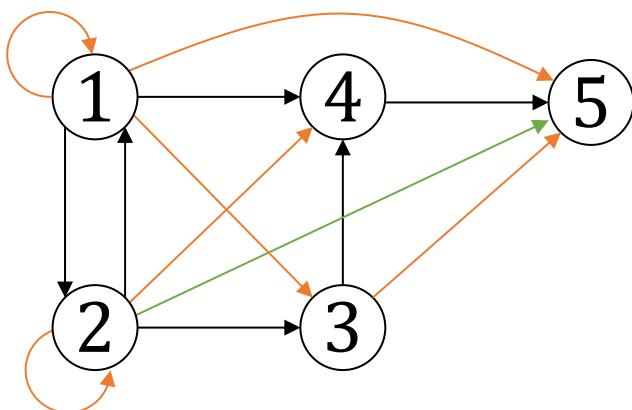


	S	T
A	1	2
1	2	1
2	3	
3	4	
4	5	

WITH RECURSIVE P AS (
 SELECT S, T
 FROM A
 UNION
 SELECT A.S, P.T
 FROM A, P
 WHERE A.T = P.S)
 SELECT *
 FROM P

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

Example



EDB
IDB

$P(x,y) :- A(x,y).$
 $P(x,y) :- A(x,z), A(z,y).$

A(S,T)



Strictly speaking, this process is iteration, not recursion:

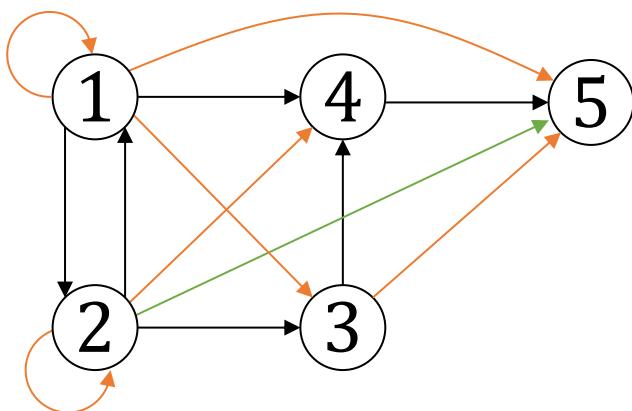
```
WITH RECURSIVE P AS (
    SELECT S, T
    FROM A
    UNION
    SELECT A.S, P.T
    FROM A, P
    WHERE A.T = P.S)
SELECT *
FROM P
```

A	S	T
1	2	
2	1	
2	3	
1	4	
3	4	
4	5	

Recursion and Iteration both repeatedly execute a set of instructions.

- Recursion (self-similarity) is when a statement in a function calls itself repeatedly.
- Iteration (repetition) is when a loop repeatedly executes until the controlling condition becomes false.

Example



A	S	T
1	2	
2	1	
2	3	
1	4	
3	4	
4	5	

EDB
IDB

$P(x,y) :- A(x,y).$
 $P(x,y) :- A(x,z), A(z,y).$

A(S,T)

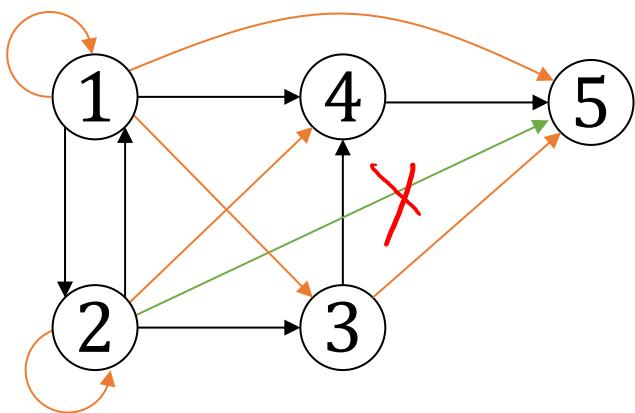


```
WITH RECURSIVE P AS (
    SELECT S, T
    FROM A
    UNION
    SELECT A1.S, A2.T
    FROM A A1, A A2
    WHERE A1.T = A2.S)
SELECT *
FROM P
```

P	S	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	

?

Example



EDB
IDB

$P(x,y) :- A(x,y).$
 $P(x,y) :- A(x,z), A(z,y).$

A(S,T)



	S	T
1	2	1
2	3	
3	1	4
4	3	4
5	4	5

```
WITH RECURSIVE P AS (
    SELECT S, T
    FROM A
    UNION
    SELECT A1.S, A2.T
    FROM A A1, A A2
    WHERE A1.T = A2.S)
SELECT *
FROM P
```

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

Challenge



- Write a query that finds the shortest path to each node from a starting node
- Create an interesting minimum database instance
- Show interesting variations
- <https://www.postgresql.org/docs/14/queries-with.html>

?



Which of these programs is recursive?

1

Local(x) :- Person(x,y,'MA').

Relative(x,x) :- Person(x,y,z).

Relative(x,y) :- Relative(x,z), Parent(z,y).

Relative(x,y) :- Relative(x,z), Parent(y,z).

Relative(x,y) :- Relative(x,z), Spouse(z,y).

Invited(y) :- Relative('myself',y), Local(y).

?

2

Local(x) :- Person(x,y,'MA').

Relative(x,x) :- Person(x,y,z).

Invited(y) :- Relative('myself',y), Local(y).

?

3

MayLike(x,y) :- Close(x,z), Likes(z,y).

Visit(x,y) :- MayLike(x,y).

Close(x,z) :- Visit(x,y), Visit(z,y).

?

Dependency Graph

- The **dependency graph** of a Datalog program is the directed graph (V,E) where
 - V is the set of **IDB** predicates (relation names)
 - E contains an arc $S \rightarrow T$ whenever there is a rule with T in the head and S in the body
- A Datalog program is **recursive** if its dependency graph contains a **cycle**



Which of these programs is recursive?

1

```
Local(x) :- Person(x,y,'MA').  
Relative(x,x) :- Person(x,y,z).  
Relative(x,y) :- Relative(x,z), Parent(z,y).  
Relative(x,y) :- Relative(x,z), Parent(y,z).  
Relative(x,y) :- Relative(x,z), Spouse(z,y).  
Invited(y) :- Relative('myself',y), Local(y).
```

?

2

```
Local(x) :- Person(x,y,'MA').  
Relative(x,x) :- Person(x,y,z).  
Invited(y) :- Relative('myself',y), Local(y).
```

?

3

```
MayLike(x,y) :- Close(x,z), Likes(z,y).  
Visit(x,y) :- MayLike(x,y).  
Close(x,z) :- Visit(x,y), Visit(z,y).
```

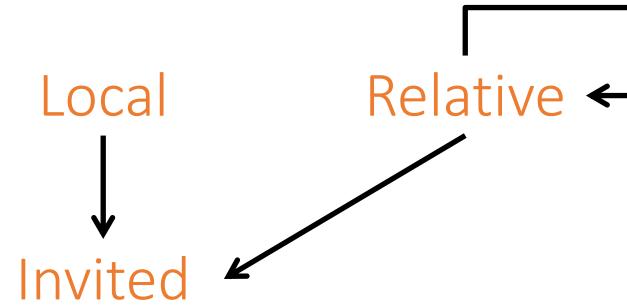
?



Which of these programs is recursive?

1

```
Local(x) :- Person(x,y,'MA').  
Relative(x,x) :- Person(x,y,z).  
Relative(x,y) :- Relative(x,z), Parent(z,y).  
Relative(x,y) :- Relative(x,z), Parent(y,z).  
Relative(x,y) :- Relative(x,z), Spouse(z,y).  
Invited(y) :- Relative('myself',y), Local(y).
```



2

```
Local(x) :- Person(x,y,'MA').  
Relative(x,x) :- Person(x,y,z).  
Invited(y) :- Relative('myself',y), Local(y).
```

?

3

```
MayLike(x,y) :- Close(x,z), Likes(z,y).  
Visit(x,y) :- MayLike(x,y).  
Close(x,z) :- Visit(x,y), Visit(z,y).
```

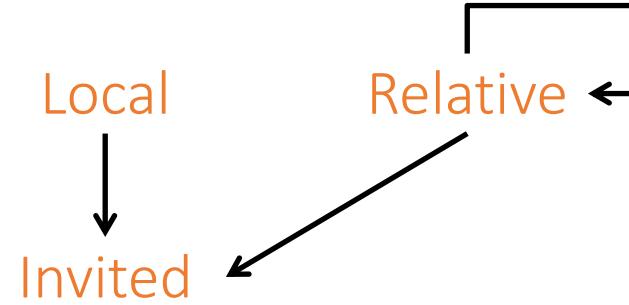
?



Which of these programs is recursive?

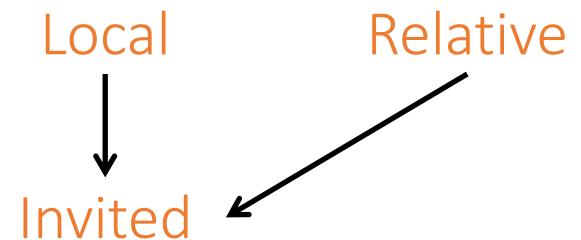
1

```
Local(x) :- Person(x,y,'MA').  
Relative(x,x) :- Person(x,y,z).  
Relative(x,y) :- Relative(x,z), Parent(z,y).  
Relative(x,y) :- Relative(x,z), Parent(y,z).  
Relative(x,y) :- Relative(x,z), Spouse(z,y).  
Invited(y) :- Relative('myself',y), Local(y).
```



2

```
Local(x) :- Person(x,y,'MA').  
Relative(x,x) :- Person(x,y,z).  
Invited(y) :- Relative('myself',y), Local(y).
```



3

```
MayLike(x,y) :- Close(x,z), Likes(z,y).  
Visit(x,y) :- MayLike(x,y).  
Close(x,z) :- Visit(x,y), Visit(z,y).
```

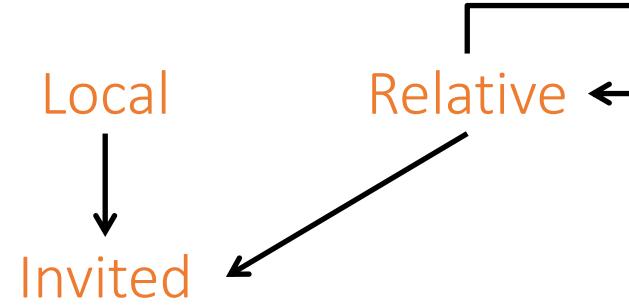
?



Which of these programs is recursive?

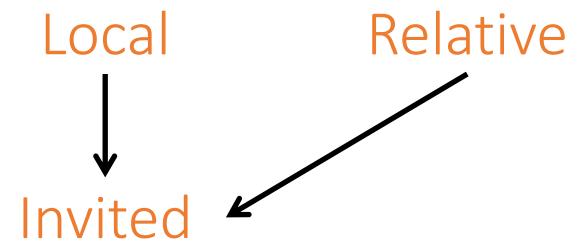
1

```
Local(x) :- Person(x,y,'MA').  
Relative(x,x) :- Person(x,y,z).  
Relative(x,y) :- Relative(x,z), Parent(z,y).  
Relative(x,y) :- Relative(x,z), Parent(y,z).  
Relative(x,y) :- Relative(x,z), Spouse(z,y).  
Invited(y) :- Relative('myself',y), Local(y).
```



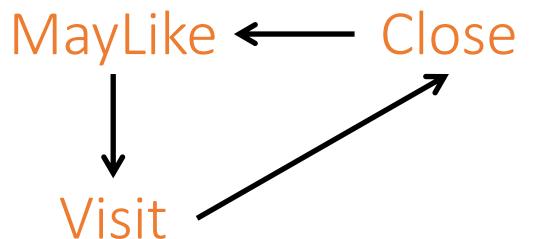
2

```
Local(x) :- Person(x,y,'MA').  
Relative(x,x) :- Person(x,y,z).  
Invited(y) :- Relative('myself',y), Local(y).
```



3

```
MayLike(x,y) :- Close(x,z), Likes(z,y).  
Visit(x,y) :- MayLike(x,y).  
Close(x,z) :- Visit(x,y), Visit(z,y).
```



Outline: T1-4: Datalog

- Datalog
 - Datalog rules
 - Recursion
 - Semantics
 - Datalog \neg : Negation, stratification
 - Datalog \pm
 - Stable model semantics (Answer set programming)
 - Datalog vs. RA
 - Naive and Semi-naive evaluation (incl. Incremental View Maintenance)

Semantics of Datalog Programs

- Let \mathbf{S} be a schema, \mathbf{D} a database over \mathbf{S} , and \mathbf{P} be a Datalog program over \mathbf{S} (i.e., all EDBs predicates belong to \mathbf{S})
- The result of evaluating \mathbf{P} over \mathbf{D} is a database \mathbf{I} over the IDB schema of \mathbf{P}
- We give 2 definitions:
 1. Fixpoint semantics *operative (think procedural)*
 2. model-theoretic *declarative*

1. Fixpoint semantics via the chase (operative definition)

Pseudo-code of a **chase** procedure:

Chase(P, D)

```
I := empty
repeat {
    if( $D \cup I$  satisfies all the rules of  $P$ ), then return  $I$ 
    Find a rule  $\text{head}(x) :- \text{body}(x,y)$  and constants  $a,b$ 
        s.t. that  $D \cup I$  contains  $\text{body}(a,b)$  but not  $\text{head}(a)$ 
     $I := I \cup \{\text{head}(a)\}$ 
}
```

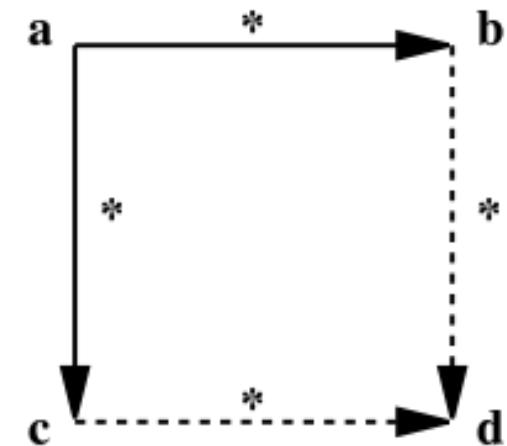
Notice since rules are monotone, I is also monotonically increasing

Nondeterminism

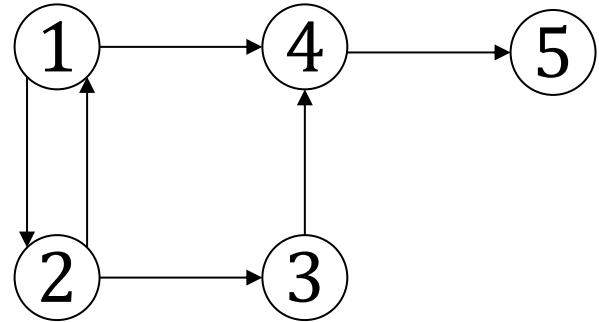
- Note: the chase is underspecified (i.e., not fully defined)
 - There can be many ways of choosing the next violation to handle
 - And each choice can lead to new violations, and so on
- We can view the choice of a new violation as **nondeterministic**

Church-Rosser property (defined for term reduction):

*If term **a** can be reduced to both **b** and **c**, then there must be a further term **d** (possibly equal to either **b** or **c**) to which both **b** and **c** can be reduced.*



Example



```
Path(x,y) :- Arc(x,y).  
Path(x,y) :- Arc(x,z), Path(z,y).  
Reachable(y) :- Path('1',y).
```

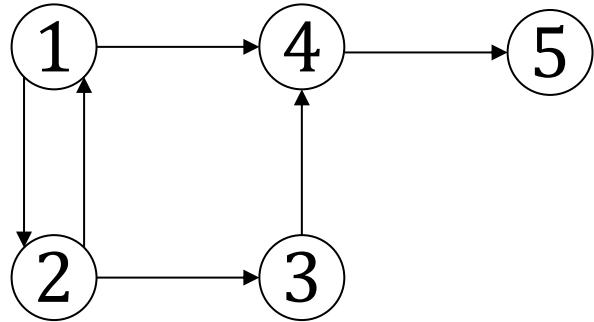
Arc

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

Path

Reachable

Example



```
Path(x,y) :- Arc(x,y).  
Path(x,y) :- Arc(x,z), Path(z,y).  
Reachable(y) :- Path('1',y).
```

Arc



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

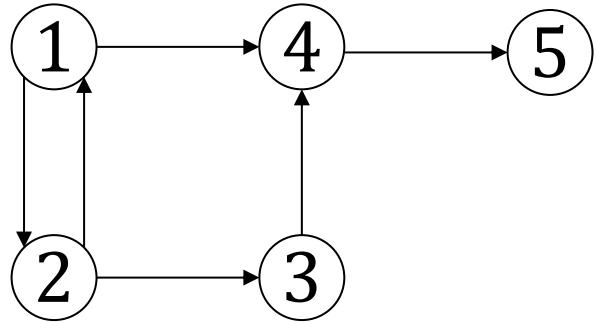
Path



1	2
---	---

Reachable

Example

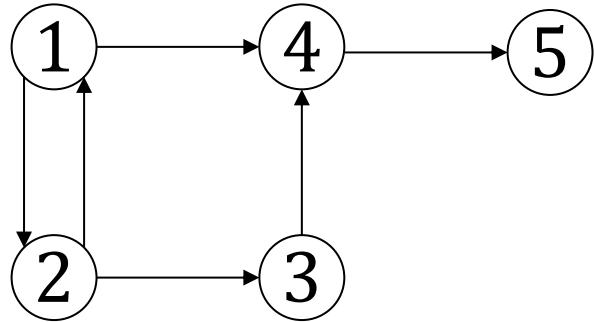


```
Path(x,y) :- Arc(x,y).  
Path(x,y) :- Arc(x,z), Path(z,y).  
Reachable(y) :- Path('1',y).
```

Arc	Path	Reachable
1 2	1 2	
2 1	2 1	
2 3		
1 4		
3 4		
4 5		

Red arrows point from the 'Arc' table to the 'Path' table, and from the 'Path' table to the 'Reachable' table.

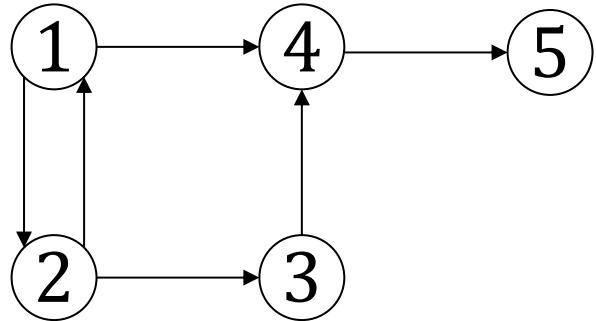
Example



```
Path(x,y) :- Arc(x,y).  
Path(x,y) :- Arc(x,z), Path(z,y).  
Reachable(y) :- Path('1',y).
```

Arc	Path	Reachable																		
<table border="1"><tr><td>1</td><td>2</td></tr><tr><td>2</td><td>1</td></tr><tr><td>2</td><td>3</td></tr><tr><td>1</td><td>4</td></tr><tr><td>3</td><td>4</td></tr><tr><td>4</td><td>5</td></tr></table>	1	2	2	1	2	3	1	4	3	4	4	5	<table border="1"><tr><td>1</td><td>2</td></tr><tr><td>2</td><td>1</td></tr><tr><td>2</td><td>3</td></tr></table>	1	2	2	1	2	3	
1	2																			
2	1																			
2	3																			
1	4																			
3	4																			
4	5																			
1	2																			
2	1																			
2	3																			

Example



```
Path(x,y) :- Arc(x,y).  
Path(x,y) :- Arc(x,z), Path(z,y).  
Reachable(y) :- Path('1',y).
```

Arc

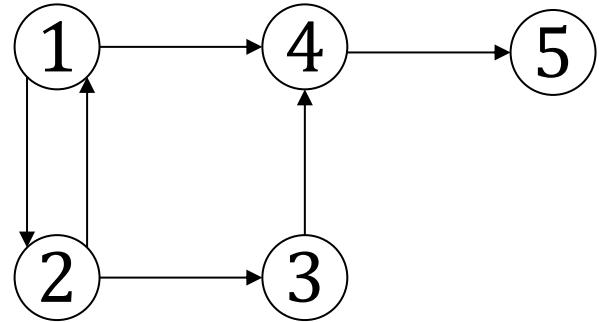
1	2
2	1
2	3
1	4
3	4
4	5

Path

1	2
2	1
2	3
1	4

Reachable

Example



Path(x,y) :- Arc(x,y).
Path(x,y) :- Arc(x,z), Path(z,y).
Reachable(y) :- Path('1',y).

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

→

1	2
2	1
2	3

→

2

2. Minimal model semantics (model-theoretic definition)

- We say that $\text{IDB } I$ is a **model** of Datalog program P (w.r.t. **EDB** D) if I satisfies all the rules of P

$$\forall \text{var} [\text{Head}(\text{IDB}) \Leftarrow \text{Body}(\text{EDB}, \text{IDB})]$$

- We say that I is a **minimal model** if I does not properly contain any other model
- Theorem: there exists one minimal model

Illustration with our example

Path(x,y) :- Arc(x,y).

Path(x,y) :- Arc(x,z), Path(z,y).

1. Fixpoint semantics



2. Minimal model semantics: smallest Path s.t.

Illustration with our example

Path(x,y) :- Arc(x,y).
Path(x,y) :- Arc(x,z), Path(z,y).

1. Fixpoint semantics

Path⁽⁰⁾ := \emptyset , t:=0

Repeat {

inc(t)

Path^(t)(x,y) := Arc(x,y) $\cup \prod_{xy} (\text{Arc}(x,z) \bowtie \text{Path}^{(t-1)}(z,y))$

until Path^(t) = Path^(t-1)}



"immediate consequence operator"

2. Minimal model semantics: smallest relation Path s.t.

Illustration with our example

Path(x,y) :- Arc(x,y).
Path(x,y) :- Arc(x,z), Path(z,y).

1. Fixpoint semantics

- $\text{Path}^{(0)} := \emptyset, t:=0$
 - Repeat {
 - inc(t)
 - $\text{Path}^{(t)}(x,y) := \text{Arc}(x,y) \cup \prod_{xy}(\text{Arc}(x,z) \bowtie \text{Path}^{(t-1)}(z,y))$
 - until $\text{Path}^{(t)} = \text{Path}^{(t-1)}$ }
- "immediate consequence operator"*

2. Minimal model semantics: smallest relation Path s.t.

$$\forall x,y [\text{Arc}(x,y) \Rightarrow \text{Path}(x,y)] \wedge$$

$$\forall x,y,z [\text{Arc}(x,z) \wedge \text{Path}(z,y) \Rightarrow \text{Path}(x,y)]$$

Datalog Semantics & equivalence b/w the definitions

(nondeterministic)

1. The **fixpoint** semantics tells us **how** to compute a Datalog query
2. The **minimal model** semantics is more declarative: only says **what we get**

THEOREM: For all Datalog programs P and DBs D
there is a **unique minimal model**,
and **every chase** returns this model

Proof sketch:

1. If I_1 and I_2 are models, so are $I_1 \cap I_2$
2. Every chase returns a model (finite)
3. Pick a chase and prove by induction: If I' is a model,
then every intermediate I is contained in I' (monotonicity)

The minimal model is the *result*, denoted $P(D)$

Semantics Summary

1. Fixpoint-theoretic

- Most "operational": Based on the immediate consequence operator for a Datalog program.

2. Model-theoretic

- Most "declarative": Based on model-theoretic semantics of first order logic.
View rules as logical constraints.

Semantics Summary

1. Fixpoint-theoretic

- Most "operational": Based on the immediate consequence operator for a Datalog program.
- Least fixpoint is reached after finitely many iterations of the immediate consequence operator.
- Basis for practical, bottom-up evaluation strategy.

2. Model-theoretic

- Most "declarative": Based on model-theoretic semantics of first order logic.
View rules as logical constraints.
- Given input DB D and Datalog program P, find the smallest possible DB instance D' that extends D and satisfies all constraints in P.

Expressiveness of Non-recursive Datalog

THEOREM: Non-recursive Datalog with built-in predicates ($<$, $>$, \leq , \geq , \neq) has the same expressive power as the positive algebra $\{\sigma, \pi, \times, U\}$

If we restrict selection to $\sigma_=($ i.e. selection with a single equality), this fragment is also called at times UCQs (Union of Conjunctive Queries) or USPJ (Union-Select-Project-Join) queries.

Monotonicity

- Can Datalog express **difference**?
 - Answer: **No!**
- Proof: Datalog is **monotone**, difference is not
 - That is, if D and D' are such that every relation of D is contained in the corresponding relation of D' ($D \subseteq D'$), then $P(D) \subseteq P(D')$

$$D \subseteq D' \Rightarrow P(D) \subseteq P(D')$$

Outline: T1-4: Datalog

- Datalog
 - Datalog rules
 - Recursion
 - Semantics
 - Datalog \neg : Negation, stratification
 - Datalog $^\pm$
 - Stable model semantics (Answer set programming)
 - Datalog vs. RA
 - Naive and Semi-naive evaluation (incl. Incremental View Maintenance)



What should be the Semantics?

Friend(x,y) :- Likes(x,y), \neg Parent(y,x).

Likes (1,2). Parent (2,1).
Likes (1,3).

?

Alternative notations are “! Parent(y,x)” or “NOT Parent(y,x)” instead of “ \neg Parent(y,x)”

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>



What should be the Semantics?

Friend(x,y) :- Likes(x,y), \neg Parent(y,x).

Likes (1,2). Parent (2,1).
Likes (1,3). \longrightarrow Friend(1,3)

Box(x) :- Item(x), \neg Box(x).

Item('ball') \longrightarrow ?



What should be the Semantics?

Friend(x,y) :- Likes(x,y), \neg Parent(y,x).

Likes (1,2). Parent (2,1).
Likes (1,3). Friend(1,3)

Box(x) :- Item(x), \neg Box(x).

Item('ball') Box('ball') ???

LeftBox(x) :- Item(x), \neg RightBox(x).

Item('ball') ?

RightBox(x) :- \neg LeftBox(x).



What should be the Semantics?

```
Friend(x,y) :- Likes(x,y), not Parent(y,x).
```

Likes (1,2). Parent (2,1).
Likes (1,3). \longrightarrow Friend(1,3)

```
Box(x) :- Item(x), not Box(x).
```

Item('ball') \longrightarrow Box('ball') ???

```
LeftBox(x) :- Item(x), not RightBox(x).
```

```
RightBox(x) :- not LeftBox(x).
```

Item('ball') \longrightarrow LeftBox('ball') ???
unsafe!

```
LeftBox(x) :- Item(x), not RightBox(x).
```

```
RightBox(x) :- Item(x), not LeftBox(x).
```

Item('ball') ?



What should be the Semantics?

Friend(x,y) :- Likes(x,y), \neg Parent(y,x).

Likes (1,2). Parent (2,1).
Likes (1,3). \longrightarrow Friend(1,3)

Box(x) :- Item(x), \neg Box(x).

Item('ball') \longrightarrow Box('ball') ???

LeftBox(x) :- Item(x), \neg RightBox(x).

Item('ball') \longrightarrow LeftBox('ball') ???
unsafe!

RightBox(x) :- \neg LeftBox(x).

LeftBox(x) :- Item(x), \neg RightBox(x).

Item('ball') \longrightarrow LeftBox('ball')
 \longrightarrow RightBox('ball')

RightBox(x) :- Item(x), \neg LeftBox(x).

\Rightarrow Adding negation to Datalog is not straightforward!

Alternative notations to " \neg Parent(y,x)" are "!" Parent(y,x)" or "NOT Parent(y,x)"

Wolfgang Gatterbauer. Principles of scalable data management: <https://northeastern-datalab.github.io/cs7240/>



What should be the Semantics?

```
Friend(x,y) :- Likes(x,y),  $\neg$  Parent(y,x).
```

Likes (1,2). Parent (2,1).
Likes (1,3). \longrightarrow Friend(1,3)

```
Box(x) :- Item(x),  $\neg$  Box(x).
```

Item('ball') \longrightarrow ~~Box('ball')~~
no model

```
LeftBox(x) :- Item(x),  $\neg$  RightBox(x).
```

```
RightBox(x) :-  $\neg$  LeftBox(x).
```

Item('ball') \longrightarrow ~~LeftBox('ball')~~
unsafe!

```
LeftBox(x) :- Item(x),  $\neg$  RightBox(x).
```

```
RightBox(x) :- Item(x),  $\neg$  LeftBox(x).
```

Item('ball') \longrightarrow LeftBox('ball')
 \longrightarrow RightBox('ball')

Later discussed "stable model" semantics

Negation in Datalog

- Various semantics have been proposed for supporting negation in Datalog that still allow tractability
- We will first look at two:
 - 1. Semipositive Datalog[¬] (restricted): PTIME
 - 2. Stratified Datalog[¬] (standard): PTIME
- We will later look at a more powerful (but intractable) semantics
 - Stable Models semantics (or answer set programming, ASP): NP-complete

1. Semipositive Programs and Safety



Friend(x,y) :- Likes(x,y), \neg Parent(y,x).

Likes $- \pi_{y,x}$ Parent

A **semipositive** program is a program where **only EDBs** may be negated

- Semantics: same as ordinary Datalog programs
- Safety: rule is safe if every variable occurs in a positive (= unnegated) relational atom (ensures domain independence: the results of programs are finite and depend only on the actual contents of the database)

Exercise: Are following rules safe?

S(x) :- T(y), Arc(z,y), \neg Arc(x,y).

?

S(x) :- T(y), \neg T(x).

?

1. Semipositive Programs and Safety



Friend(x,y) :- Likes(x,y), \neg Parent(y,x).

Likes $- \pi_{y,x}$ Parent

A **semipositive** program is a program where **only EDBs** may be negated

- Semantics: same as ordinary Datalog programs
- Safety: rule is safe if every variable occurs in a positive (= unnegated) relational atom (ensures domain independence: the results of programs are finite and depend only on the actual contents of the database)

Exercise: Are following rules safe?

S(x) :- T(y), Arc(z,y), \neg Arc(x,y).

unsafe

P(\rightarrow)

S(x) :- T(y), \neg T(x).

unsafe

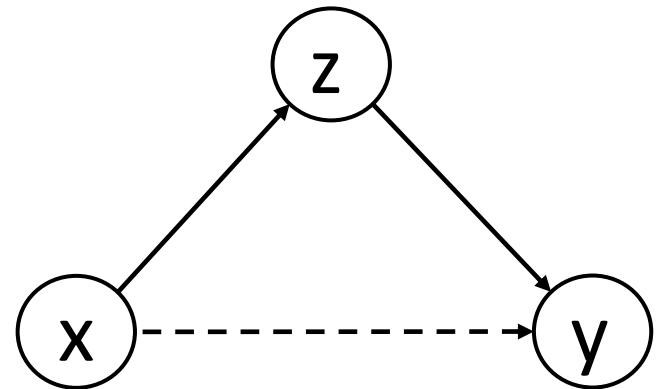
1. Semipositive: Negated Atoms

Arc(Source,Target)



- We may put \neg or $!$ or **not** in front of an atom, to negate its meaning.
- EXAMPLE: For any node x , return all nodes y that are two hops away, where y is not an immediate neighbor of x .

?



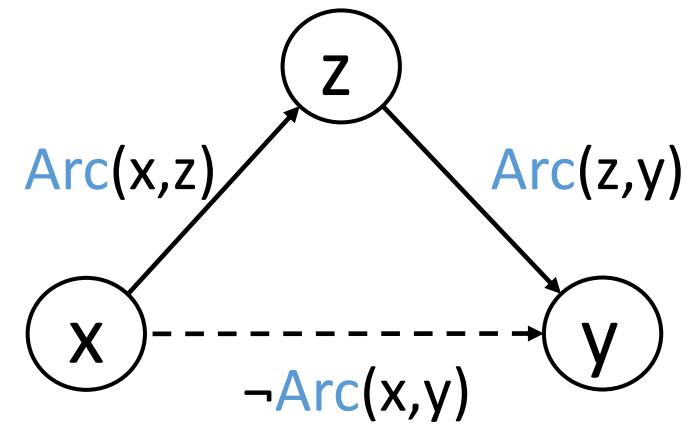
1. Semipositive: Negated Atoms

Arc(Source,Target)



- We may put \neg or ! or not in front of an atom, to negate its meaning.
- EXAMPLE: For any node x, return all nodes y that are two hops away, where y is not an immediate neighbor of x.

TwoHopsAway(x,y) :- Arc(x,z), Arc(z,y), \neg Arc(x,y).



A(S,T)



501

SQL ?

1. Semipositive: Negated Atoms

Arc(Source,Target)



- We may put \neg or ! or not in front of an atom, to negate its meaning.
- EXAMPLE: For any node x, return all nodes y that are two hops away, where y is not an immediate neighbor of x.

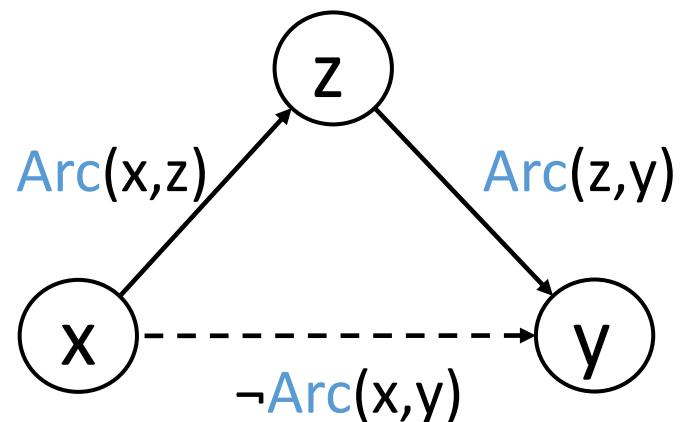
TwoHopsAway(x,y) :- Arc(x,z), Arc(z,y), \neg Arc(x,y).

A(S,T)



501

```
SELECT A1.S, A2.T  
FROM A A1, A A2  
WHERE A1.T = A2.S  
AND NOT EXISTS  
(SELECT *  
FROM A A3  
WHERE A3.S = A1.S  
AND A3.T = A2.T)
```



Example: beyond Semipositive

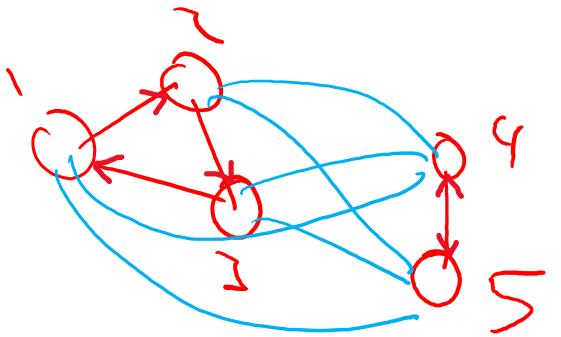
Arc(Source,Target)
Node(id)



Compute all pairs of disconnected nodes in a graph.

?

Node is basically ADom:
 $\text{Node}(x) := \text{Arc}(x,y)$
 $\text{Node}(y) := \text{Arc}(x,y)$



Example: beyond Semipositive

Arc(Source,Target)
Node(id)



Compute all pairs of disconnected nodes in a graph.

```
Reachable(x,y) :- Arc(x,y).
```

```
Reachable(x,y) :- Arc(x,z), Reachable(z,y).
```

Example: beyond Semipositive

Arc(Source,Target)
Node(id)

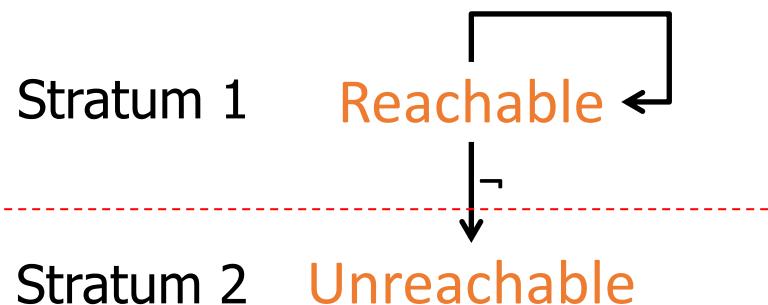


Compute all pairs of disconnected nodes in a graph.

Reachable(x,y) :- Arc(x,y).

Reachable(x,y) :- Arc(x,z), Reachable(z,y).

Unreachable(x,y) :- Node(x), Node(y), \neg Reachable(x,y).



- Straightforward syntactic restriction.
- When the Datalog program is stratified, we can evaluate IDB predicates stratum-by-stratum
- Once evaluated, treat it as EDB for higher strata.

Precedence graph

- Nodes = IDB predicates
 - Arc p → q if predicate q depends on p
 - Label this arc " \neg " if predicate p is negated
- think: "topological sort"

Non-stratified example: LeftBox(x) :- \neg LeftBox(x), Item(x).

?

Example: beyond Semipositive

Arc(Source,Target)
Node(id)

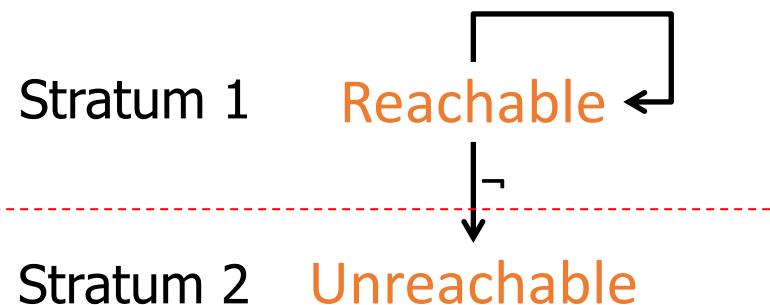


Compute all pairs of disconnected nodes in a graph.

Reachable(x,y) :- Arc(x,y).

Reachable(x,y) :- Arc(x,z), Reachable(z,y).

Unreachable(x,y) :- Node(x), Node(y), \neg Reachable(x,y).

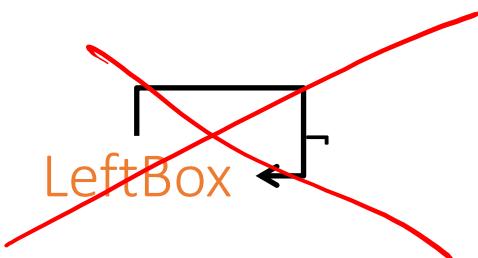


- Straightforward syntactic restriction.
- When the Datalog program is stratified, we can evaluate IDB predicates stratum-by-stratum
- Once evaluated, treat it as EDB for higher strata.

Precedence graph

- Nodes = IDB predicates
 - Arc p \rightarrow q if predicate q depends on p
 - Label this arc " \neg " if predicate p is negated
- think: "topological sort"

Non-stratified example: LeftBox(x) :- \neg LeftBox(x), Item(x).



2. Stratified Programs: Definition and Semantics

- DEFINITION: Let P be a Datalog program, E be the set of EDB predicates, and I be the set of IDB predicates. A **stratification** of P is a **partitioning** of the IDB predicates into disjoint sets I_1, \dots, I_k such that:
 - For $i=1, \dots, k$, every rule with head in I_i has possible body predicates only from E, I_1, \dots, I_i
 - For $i=1, \dots, k$, every rule with head in I_i has **negated** body predicates only from E, I_1, \dots, I_{i-1}
- SEMANTICS:
 - For $i=1, \dots, k$:
 - Compute the IDBs of the stratum I_i , possibly via recursion
 - Add computed IDBs to the EDBs
 - Due to the definition of stratification, each E_i can be viewed as **semipositive**

2. Theorems on Stratification

- THEOREM 1: A program has a **stratification** if and only if its **dependency graph** does not contain a **cycle with a "negated edge"**
 - Dependency graph is defined as previously, except that edges can be labeled with negation
 - Hence, we can test for stratifiability efficiently, via graph reachability

```
A(x) :- B(x).  
B(x) :- C(x).  
C(x) :-  $\neg$ A(x).
```

Stratified

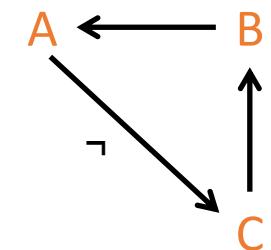
?

2. Theorems on Stratification

- THEOREM 1: A program has a **stratification** if and only if its **dependency graph** does not contain a **cycle with a "negated edge"**

- Dependency graph is defined as previously, except that edges can be labeled with negation
 - Hence, we can test for stratifiability efficiently, via graph reachability

$A(x) :- B(x).$
 $B(x) :- C(x).$
 $C(x) :- \neg A(x).$



- THEOREM 2: **Non-recursive Datalog with negation** can always be stratified in the topological order
- THEOREM 3: **Non-recursive Datalog with negation** has the same expressive power as the algebra $\{\sigma_-, \pi, \times, U, -\}$
 - Extendable to RA if we add the comparison predicates $<$, $>$, \neq , \leq , \geq

2. Stratification practice

Parent(P,C)



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

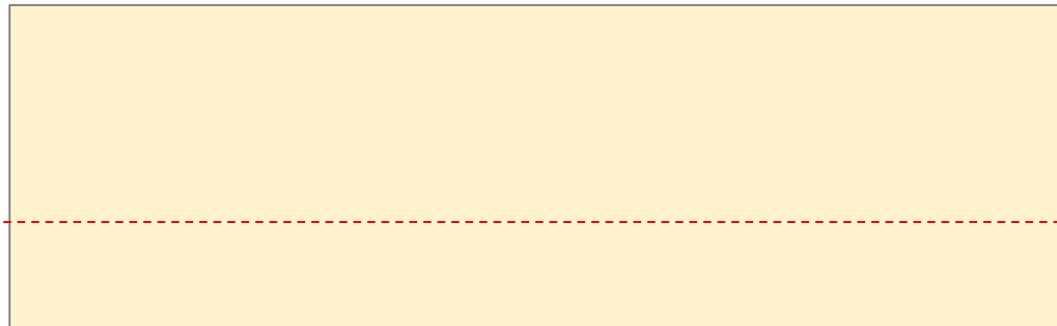
?

2. Stratification practice

Parent(P,C)



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



first compute for each
person their descendants

then use negation

2. Stratification practice

Parent(P,C)

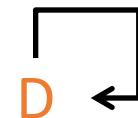


Q: Find all descendants of Alice,
who are not descendants of Bob

$D(x,y) :- \text{Parent}(x,y).$

$D(x,z) :- \text{Parent}(y,z), D(x,y).$

first compute for each
person their descendants



then use negation

?

2. Stratification practice

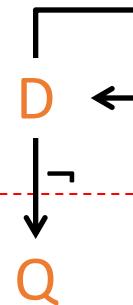
Parent(P,C)



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

```
D(x,y) :- Parent(x,y).  
D(x,z) :- Parent(y,z), D(x,y).  
Q(x) :- D('Alice',x),  $\neg$ D('Bob',x).
```

first compute for each
person their descendants
then use negation



2. Stratification practice

Parent(P,C)



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

$D(x,y) :- \text{Parent}(x,y).$

$D(x,z) :- \text{Parent}(y,z), D(x,y).$

$Q(x) :- D('Alice',x), \neg D('Bob',x).$

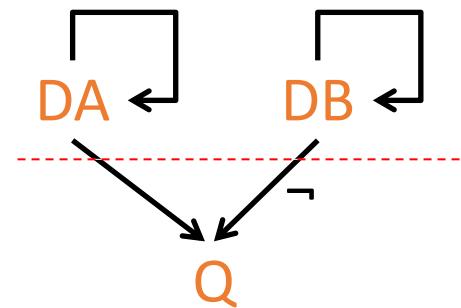
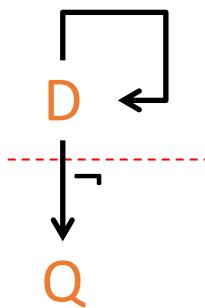
$DA(y) :- \text{Parent}('Alice',y).$

$DA(y) :- \text{Parent}(x,y), DA(x).$

$DB(y) :- \text{Parent}('Bob',y).$

$DB(y) :- \text{Parent}(x,y), DB(x).$

$Q(x) :- DA(x), \neg DB(x).$



Outline: T1-4: Datalog

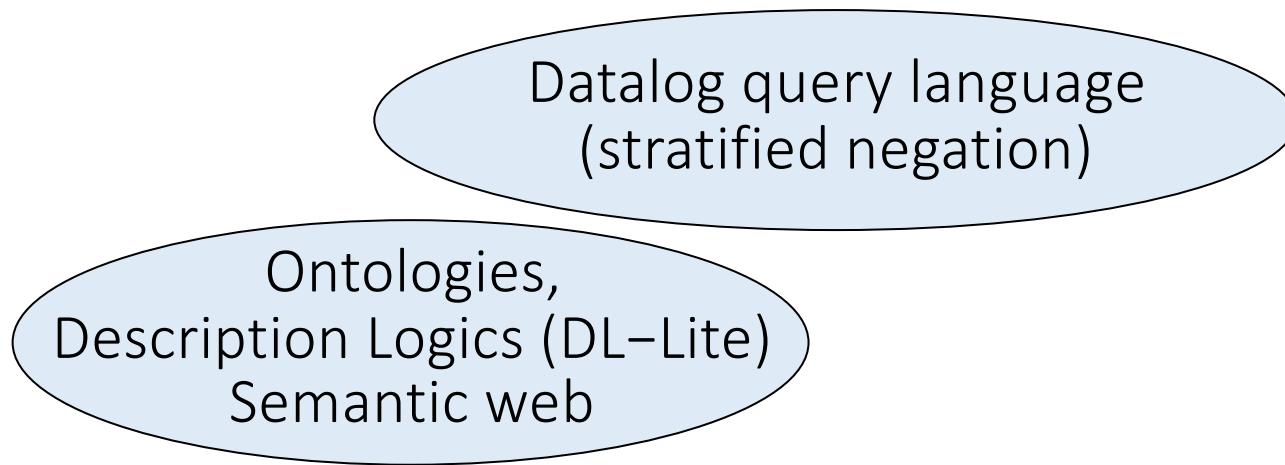
- Datalog
 - Datalog rules
 - Recursion
 - Semantics
 - Datalog \neg : Negation, stratification
 - Datalog \pm
 - Stable model semantics (Answer set programming)
 - Datalog vs. RA
 - Naive and Semi-naive evaluation (incl. Incremental View Maintenance)

Datalog[±]: background

Datalog query language
(stratified negation)

- Much is possible with Datalog

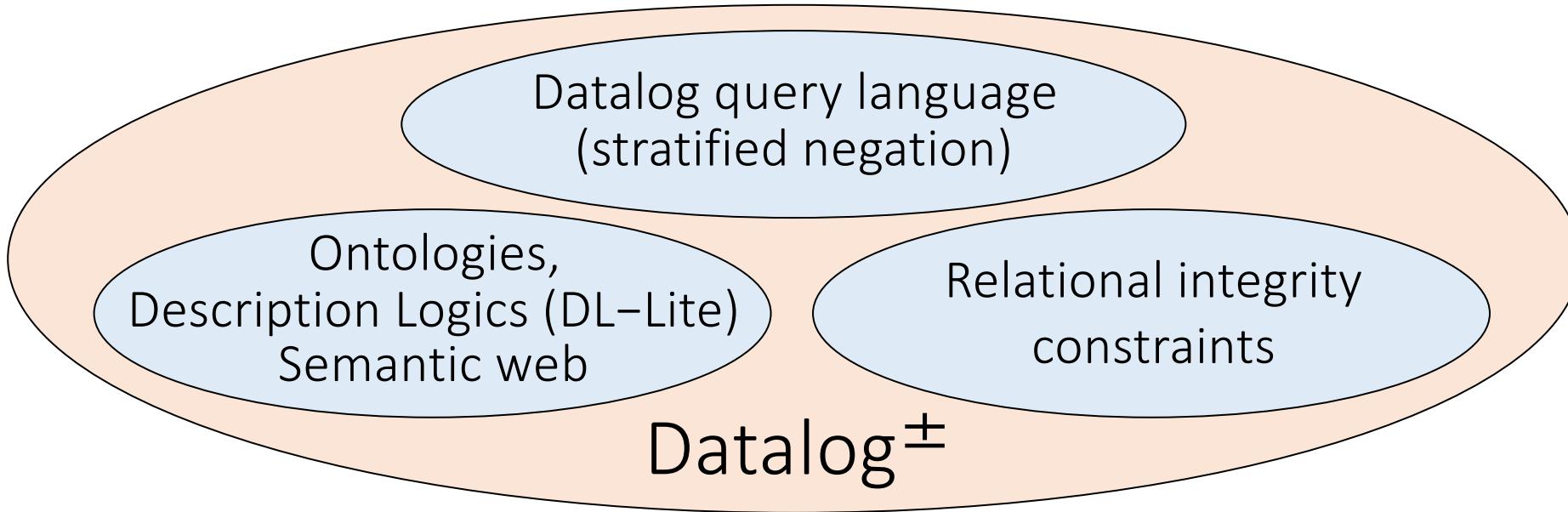
Datalog[±]: background



- Much is possible with Datalog
- Much is not (observed e.g. by [Patel-Schneider, Horrocks 2006])

Patel-Schneider, Horrocks. Position paper: A comparison of two modelling paradigms in the Semantic Web. WWW (Semantic Web track). 2006. <https://dl.acm.org/doi/10.1145/1135777.1135784>

Datalog[±]: goal



- Much is possible with Datalog
- Much is not (observed e.g. by Patel-Schneider, Horrocks 2006)
- **Datalog[±]** is a framework that extends Datalog with:
 - **value invention** (\exists -variables in the head): **TGDs** (Tuple-Generating Dependencies)
 - **equality predicate** in the head: **EGDs** (Equality Generating Dependencies)
 - **constant \perp** in the head: **negative constraints** (disjointness)

Datalog and expressiveness for ontological reasoning

Assertion type	Datalog rule
Inclusion	$\text{emp}(X) \rightarrow \text{person}(X)$
(Inverse) role inclusion	$\text{reportsTo}(X, Y) \rightarrow \text{manages}(Y, X)$
Reflexive expansion	$\text{boss}(X) \rightarrow \text{manages}(X, X)$
Transitivity	$\text{manages}(X, Y), \text{manages}(Y, Z) \rightarrow \text{manages}(X, Z)$
Concept product	$\text{seniorEmp}(X), \text{emp}(Y) \rightarrow \text{higher}(X, Y)$
Participation	?
Disjointness	?
Functionality	?

Ontology assertion	Datalog [±] rule
Participation	$\text{boss}(X) \rightarrow \exists Y \text{ reports}(Y, X)$
Disjointness	$\text{customer}(X), \text{boss}(X) \rightarrow \perp$
Functionality	$\text{reports}(X, Y_1), \text{reports}(X, Y_2) \rightarrow Y_1 = Y_2$

Datalog[±] vs. DL

*FiveStar(X) → Hotel(X),
FiveStar(X), notPool(X, Y) → ∃Z Beach(X, Z),
FiveStar(X), notBeach(X, Y) → ∃Z Pool(X, Z),
Beach(X, Y) → ∃Z SwimOpp(X, Z),
Pool(X, Y) → ∃Z SwimOpp(X, Z),*

The above example corresponds to the following set of DL axioms, expressed in an extension of \mathcal{ELHI} by nonmonotonic negation:

<i>FiveStar</i>	\sqsubseteq	Hotel,
<i>FiveStar</i> \sqcap not \exists Pool	\sqsubseteq	\exists Beach,
<i>FiveStar</i> \sqcap not \exists Beach	\sqsubseteq	\exists Pool,
\exists Beach	\sqsubseteq	\exists SwimOpp,
\exists Pool	\sqsubseteq	\exists SwimOpp,

Interesting Observations

- Exploiting schema knowledge in query answering is not trivial
- Languages and algorithms exist that allow for tractable query answering
- Applications in real-world scenarios are possible
 - Industrial applications in data integration, Semantic Web, ontological reasoning
 - Companies and Products: RelationalAI, Deepreason.ai, Oracle Semantic Technologies, IBM IODT, OntoDLV (Vienna)