

# T1: Data models and query languages

## L5: Datalog continued

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp20)

<https://northeastern-datalab.github.io/cs7240/sp20/>

Version 1/21/2020

# Outline: Datalog

- Datalog
  - Datalog rules
  - Recursion
  - Semantics
  - Datalog<sup>¬</sup>: Negation, stratification
  - Datalog<sup>±</sup>
  - Stable model semantics (Answer set programming)
  - Datalog vs. RA
  - Naive and Semi-naive evaluation

# Syntactic Constraints



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

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

The rule stands for the following logical formula:

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

Restrictions:

- **Safety**: every variable in  $\mathbf{x}$  should occur in the body at least once

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

# Syntactic Constraints



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

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

The rule stands for the following logical formula:

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

Restrictions:

- **Safety**: every variable in  $\mathbf{x}$  should occur in the body at least once

~~$$R(\mathbf{x}, \mathbf{z}) \text{ :- } S(\mathbf{x}, \mathbf{y}), R(\mathbf{y}, \mathbf{x})$$~~

*forbidden rule:  $\mathbf{z}$  not in body*

- The head predicate  $R$  must be an **IDB** predicate  
(Body can include both EDBs and IDBs)

$$\text{Edge}(\mathbf{x}, \mathbf{y}) \text{ :- } \text{Edge}(\mathbf{x}, \mathbf{z}), \text{Edge}(\mathbf{z}, \mathbf{y})$$

# Syntactic Constraints



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

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

The rule stands for the following logical formula:

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

Restrictions:

- **Safety**: every variable in  $\mathbf{x}$  should occur in the body at least once

~~$$R(\mathbf{x}, \mathbf{z}) \text{ :- } S(\mathbf{x}, \mathbf{y}), R(\mathbf{y}, \mathbf{x})$$~~

*forbidden rule:  $\mathbf{z}$  not in body*

- The head predicate  $R$  must be an **IDB** predicate  
(Body can include both EDBs and IDBs)

~~$$\text{Edge}(\mathbf{x}, \mathbf{y}) \text{ :- } \text{Edge}(\mathbf{x}, \mathbf{z}), \text{Edge}(\mathbf{z}, \mathbf{y})$$~~

*assuming Edge is EDB*

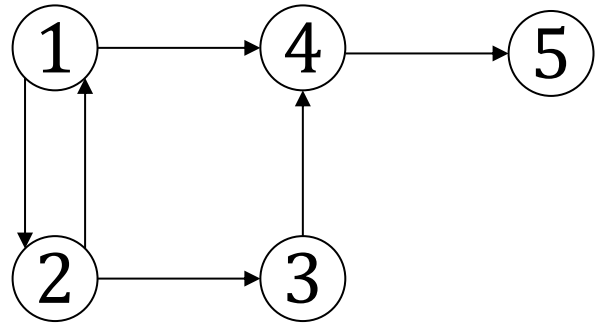
# Outline: Datalog

- Datalog
  - Datalog rules
  - **Recursion**
  - Semantics
  - Datalog<sup>¬</sup>: Negation, stratification
  - Datalog<sup>±</sup>
  - Stable model semantics (Answer set programming)
  - Datalog vs. RA
  - Naive and Semi-naive evaluation

# Datalog program

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

# Example



EDB

$P(x,y) :- E(x,y)$

IDB

$P(x,y) :- E(x,z), P(z,y)$

*recursion due to head in rule body*



*What does it compute?*

?

E

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



# Example

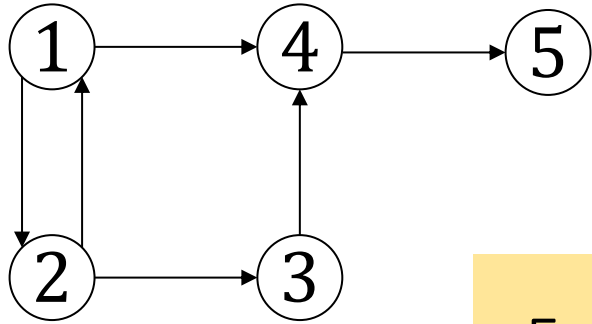


EDB  $P(x,y) :- E(x,y)$

*recursion due to head in rule body*

IDB  $P(x,y) :- E(x,z), P(z,y)$

*What does it compute?*



*Calculates all paths (transitive closure)*

For all nodes  $x$  and  $y$ :  
If there is an **edge** from  $x$  to  $y$ ,  
then there is a **path** from  $x$  to  $y$ .

For all nodes  $x$ ,  $y$ , and  $z$ :  
If there is an **edge** from  $x$  to  $z$ , and there is a **path** from  $z$  to  $y$   
then there is a **path** from  $x$  to  $y$ .

E

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

# Example



EDB

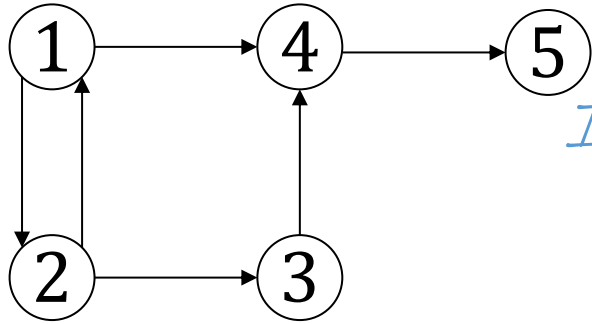
$P(x,y) :- E(x,y)$

IDB

$P(x,y) :- E(x,z), P(z,y)$

*recursion due to head in rule body*

*What does it compute?*



*Initially: P is empty*

E

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

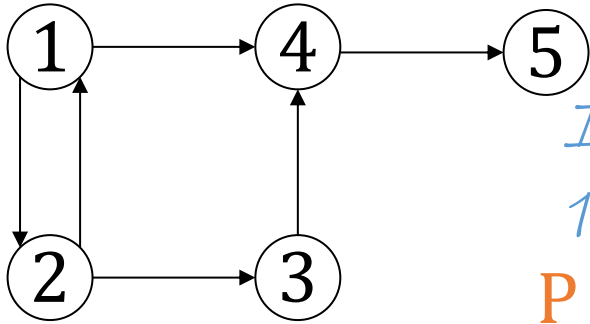
# Example



**EDB**  $P(x,y) :- E(x,y)$   
**IDB**  $P(x,y) :- E(x,z), P(z,y)$

*recursion due to head in rule body*

*What does it compute?*



*Initially: P is empty*

*1<sup>st</sup> iteration*

**P**

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

*P=E from 1<sup>st</sup> rule*

**E**

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

*2<sup>nd</sup> rule generates nothing (because P is empty)*

# Example



EDB

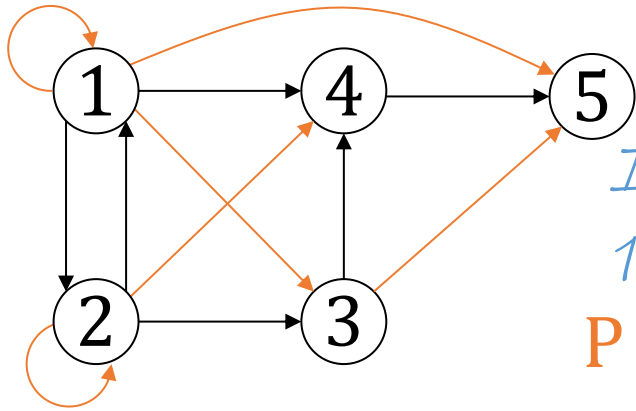
$P(x,y) :- E(x,y)$

*recursion due to head in rule body*

IDB

$P(x,y) :- E(x,z), P(z,y)$

*What does it compute?*



*Initially: P is empty*

*1<sup>st</sup> iteration*

**P**

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

*P=E from 1<sup>st</sup> rule*

*2<sup>nd</sup> iteration*

**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

*1<sup>st</sup> rule*

*2<sup>nd</sup> rule*

*2<sup>nd</sup> rule generates nothing (because P is empty)*

*New facts from 2<sup>nd</sup> rule*

**E**

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

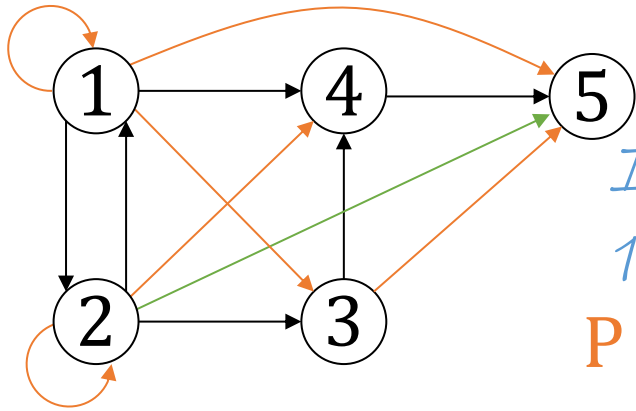
# Example



EDB  $P(x,y) :- E(x,y)$   
 IDB  $P(x,y) :- E(x,z), P(z,y)$

*recursion due to head in rule body*

*What does it compute?*



*Initially: P is empty*

*1<sup>st</sup> iteration*

**P**

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

*P=E from 1<sup>st</sup> rule*

*2<sup>nd</sup> iteration*

**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

*1<sup>st</sup> rule*

*2<sup>nd</sup> rule*

*3<sup>rd</sup> iteration*

**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
2	5

*1<sup>st</sup> + 2<sup>nd</sup> rule*

*1<sup>st</sup> rule*

*2<sup>nd</sup> rule*

*2<sup>nd</sup> rule generates nothing (because P is empty)*

*New facts from 2<sup>nd</sup> rule*

**E**

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

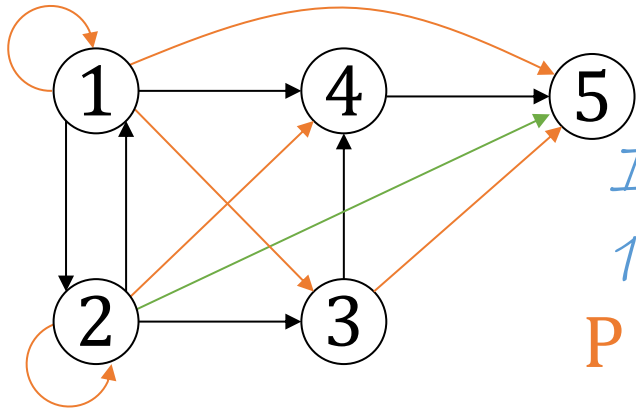
# Example



EDB  $P(x,y) :- E(x,y)$   
 IDB  $P(x,y) :- E(x,z), P(z,y)$

*recursion due to head in rule body*

*What does it compute?*



*Initially: P is empty*

*1<sup>st</sup> iteration*

**P**

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

*P=E from 1<sup>st</sup> rule*

*2<sup>nd</sup> iteration*

**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

*1<sup>st</sup> rule*

*2<sup>nd</sup> rule*

*3<sup>rd</sup> iteration = 4<sup>th</sup> iteration*

**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
2	5

*1<sup>st</sup> + 2<sup>nd</sup> rule*

*1<sup>st</sup> rule*

*2<sup>nd</sup> rule*

**E**

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

*2<sup>nd</sup> rule generates nothing (because P is empty)*

*New facts from 2<sup>nd</sup> rule*

# 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 edge  $R \rightarrow S$  whenever there is a rule with  $S$  in the head and  $R$  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

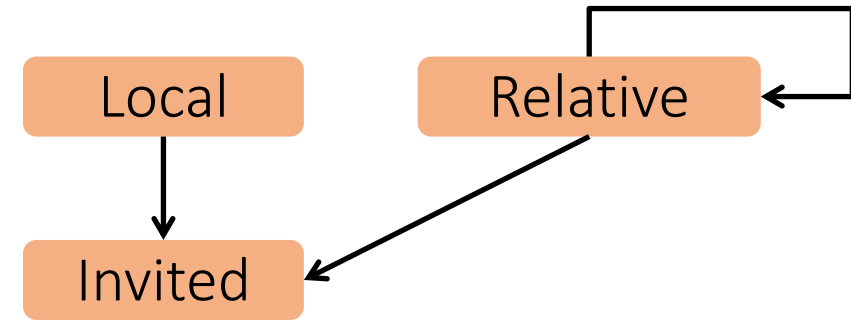
```
Local(x) :- Person(x,y,'MA')
Relative(x,x) :- Person(x,y,z)
Relative(x,y) :- Relative(x,z),Parent(z,y)
Relative(x,y) :- Relative(x,z),Parent(y,z)
Relative(x,y) :- Relative(x,z),Spouse(z,y)
Invited(y) :- Relative('myself',y),Local(y)
```

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

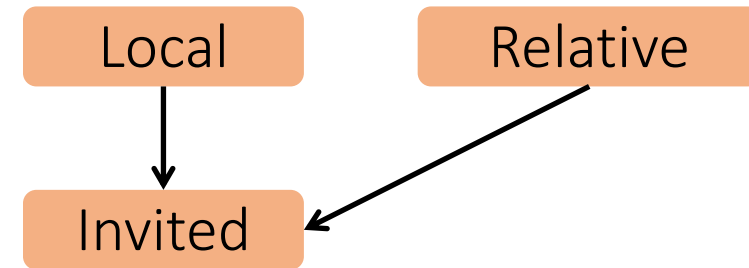
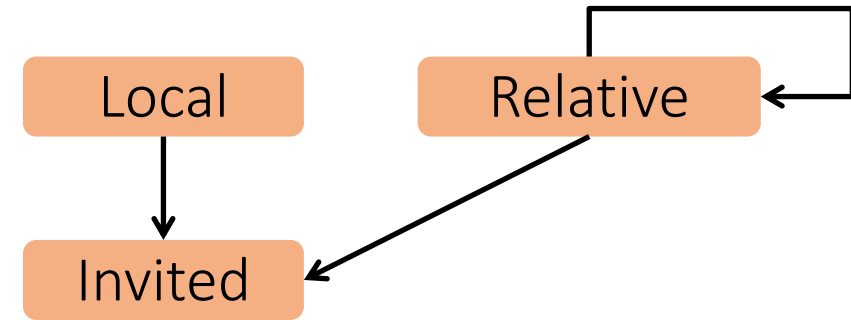
```
Local(x) :- Person(x,y,'MA')
Relative(x,x) :- Person(x,y,z)
Relative(x,y) :- Relative(x,z),Parent(z,y)
Relative(x,y) :- Relative(x,z),Parent(y,z)
Relative(x,y) :- Relative(x,z),Spouse(z,y)
Invited(y) :- Relative('myself',y),Local(y)
```

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

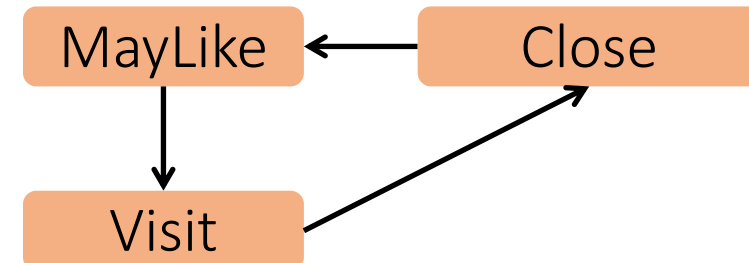
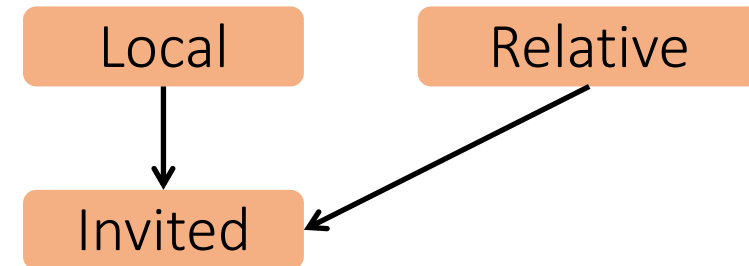
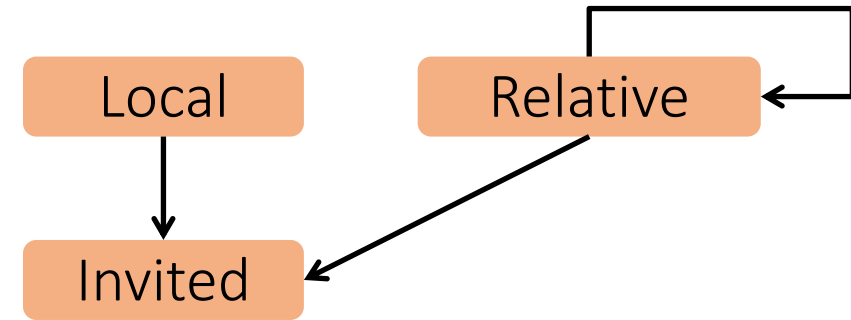
```
Local(x) :- Person(x,y,'MA')
Relative(x,x) :- Person(x,y,z)
Relative(x,y) :- Relative(x,z),Parent(z,y)
Relative(x,y) :- Relative(x,z),Parent(y,z)
Relative(x,y) :- Relative(x,z),Spouse(z,y)
Invited(y) :- Relative('myself',y),Local(y)
```

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: Datalog

- Datalog
  - Datalog rules
  - Recursion
  - **Semantics**
  - Datalog<sup>¬</sup>: Negation, stratification
  - Datalog<sup>±</sup>
  - Stable model semantics (Answer set programming)
  - Datalog vs. RA
  - Naive and Semi-naive evaluation

# 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
  2. model-theoretic

# 1. Fixpoint semantics via the chase (operative definition)

A **chase** procedure is a program of the following form:

Chase(**P**,**D**)

- **I** := empty
- while(true) {
  - if(**D**∪**I** satisfies all the rules of **P**), then **return I**
  - **Find a rule** head(**x**) :- body(**x**,**y**) and tuples **a**,**b** s.t. that **D**∪**I** contains body(**a**,**b**) but not head(**a**)
  - **I** := **I** ∪ {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**



# Example



Person(id, gender, country)  
Parent(parent, child)

(To redo with Parent(child, parent) and Descendent..)

*redo*

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

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

Invited(y) :- Relative('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

Relative      Invited

# Example



Person(id, gender, country)  
Parent(parent, child)

→ **Relative**(x, x) :- **Person**(x, y, z)  
**Relative**(x, y) :- **Relative**(x, z), **Parent**(z, y)  
**Invited**(y) :- **Relative**('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

Relative      Invited

# Example



Person(id, gender, country)  
Parent(parent, child)

→ **Relative**(x, x) :- **Person**(x, y, z)  
**Relative**(x, y) :- **Relative**(x, z), **Parent**(z, y)  
**Invited**(y) :- **Relative**('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

1	1
---	---

---

*Invited*

# Example



Person(id, gender, country)  
Parent(parent, child)

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

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

→ Invited(y) :- Relative('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

1	1
---	---

---

Invited

# Example



Person(id, gender, country)  
Parent(parent, child)

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

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

→ Invited(y) :- Relative('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

1	1
---	---

1
---

# Example



Person(id, gender, country)  
Parent(parent, child)

Relative(x, x) :- Person(x, y, z)  
→ Relative(x, y) :- Relative(x, z), Parent(z, y)  
Invited(y) :- Relative('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

1	1
---	---

1
---

# Example



Person(id, gender, country)  
Parent(parent, child)

Relative(x, x) :- Person(x, y, z)  
→ Relative(x, y) :- Relative(x, z), Parent(z, y)  
Invited(y) :- Relative('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

1	1
1	2

1
---

# Example



Person(id, gender, country)  
Parent(parent, child)

Relative(x, x) :- Person(x, y, z)  
→ Relative(x, y) :- Relative(x, z), Parent(z, y)  
Invited(y) :- Relative('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

1	1
1	2

1
---



# Example



Person(id, gender, country)  
Parent(parent, child)

Relative(x, x) :- Person(x, y, z)  
→ Relative(x, y) :- Relative(x, z), Parent(z, y)  
Invited(y) :- Relative('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

1	1
1	2
1	3

1
---

# Example



Person(id, gender, country)  
Parent(parent, child)

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

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

→ Invited(y) :- Relative('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

1	1
1	2
1	3

1
---

# Example



Person(id, gender, country)  
Parent(parent, child)

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

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

→ Invited(y) :- Relative('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

1	1
1	2
1	3

1
2

# Example



Person(id, gender, country)  
Parent(parent, child)

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

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

→ Invited(y) :- Relative('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

1	1
1	2
1	3

1
2

# Example



Person(id, gender, country)  
Parent(parent, child)

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

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

→ Invited(y) :- Relative('1', y)

1	F	IL
2	M	AU
3	F	NL
4	M	IL
5	F	FR

1	2
2	3

1	1
1	2
1	3

1
2
3

## 2. Minimal model semantics

- We say that **IDB I** is a **model** of **P** (w.r.t. **EDB D**) if **DUI** 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

$\text{Path}(x,y) \text{ :- Edge}(x,y)$   
 $\text{Path}(x,z) \text{ :- Edge}(x,y), \text{Path}(y,z)$

## 1. Fixpoint semantics



## 2. Minimal model semantics: smallest $\text{Path}$ s.t.

# Illustration with our example

$\text{Path}(x,y) \text{ :- Edge}(x,y)$   
 $\text{Path}(x,z) \text{ :- Edge}(x,y), \text{Path}(y,z)$

## 1. Fixpoint semantics

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

## 2. Minimal model semantics: smallest $\text{Path}$ s.t.



# Illustration with our example

$\text{Path}(x,y) :- \text{Edge}(x,y)$   
 $\text{Path}(x,z) :- \text{Edge}(x,y), \text{Path}(y,z)$

## 1. Fixpoint semantics

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

## 2. Minimal model semantics: smallest $\text{Path}$ s.t.

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

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

# Datalog Semantics & equivalence b/w the definitions

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
3. Pick a chase and prove by induction: If  $I'$  is a model, then every intermediate  $I$  is contained in  $I'$

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 has the same expressive power as the algebra  $\{\sigma_=\, \pi, \times, \rho, \cup\}$  (where  $\sigma_ =$  means selection with a single equality)

This fragment is often called “**positive RA**” or **USPJ**  
(union-select-project-join)

# 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'** ( $\mathbf{D} \subseteq \mathbf{D}'$ ), then  $\mathbf{P}(\mathbf{D}) \subseteq \mathbf{P}(\mathbf{D}')$

$$\mathbf{D} \subseteq \mathbf{D}' \Rightarrow \mathbf{P}(\mathbf{D}) \subseteq \mathbf{P}(\mathbf{D}')$$

# Outline: Datalog

- Datalog
  - Datalog rules
  - Recursion
  - Semantics
  - Datalog<sup>¬</sup>: Negation, stratification
  - Datalog<sup>±</sup>
  - Stable model semantics (Answer set programming)
  - Datalog vs. RA
  - Naive and Semi-naive evaluation

# What is the Semantics?



`Buddy(x,y) :- Likes(x,y), ¬Parent(y,x)`





# What is the Semantics?



$\text{Buddy}(x,y) :- \text{Likes}(x,y), \neg \text{Parent}(y,x)$

~~$\text{LeftBox}(x) :- \neg \text{LeftBox}(x), \text{Item}(x)$~~



# What is the Semantics?



$\text{Buddy}(x,y) \text{ :- Likes}(x,y), \neg \text{Parent}(y,x)$

$\text{LeftBox}(x) \text{ :- } \neg \text{LeftBox}(x), \text{Item}(x)$

$\text{Left}(x) \text{ :- } \neg \text{Right}(x), \text{Item}(x)$

$\text{Right}(x) \text{ :- } \neg \text{Left}(x)$

$\text{Item}(\text{'ball'}) \longrightarrow ? \text{LeftBox}(\text{'ball'})$

?

# What is the Semantics?



$\text{Buddy}(x,y) \text{ :- Likes}(x,y), \neg \text{Parent}(y,x)$

$\text{LeftBox}(x) \text{ :- } \neg \text{LeftBox}(x), \text{Item}(x)$

$\text{Left}(x) \text{ :- } \neg \text{Right}(x), \text{Item}(x)$

$\text{Right}(x) \text{ :- } \neg \text{Left}(x)$

$\text{Left}(x) \text{ :- } \neg \text{Right}(x), \text{Item}(x)$

$\text{Right}(x) \text{ :- } \neg \text{Left}(x), \text{Item}(x)$

$\text{Item}(\text{'ball'}) \longrightarrow ? \text{LeftBox}(\text{'ball'})$

*unsafe*

?

# What is the Semantics?



$\text{Buddy}(x,y) \text{ :- Likes}(x,y), \neg \text{Parent}(y,x)$

$\text{LeftBox}(x) \text{ :- } \neg \text{LeftBox}(x), \text{Item}(x)$

$\text{Left}(x) \text{ :- } \neg \text{Right}(x), \text{Item}(x)$

$\text{Right}(x) \text{ :- } \neg \text{Left}(x)$

$\text{Left}(x) \text{ :- } \neg \text{Right}(x), \text{Item}(x)$

$\text{Right}(x) \text{ :- } \neg \text{Left}(x), \text{Item}(x)$

$\text{Item}(\text{'ball'}) \longrightarrow ? \text{LeftBox}(\text{'ball'})$

*unsafe*

$\text{Item}(\text{'ball'}) \begin{cases} \longrightarrow \text{LeftBox}(\text{'ball'}) \\ \longrightarrow \text{RightBox}(\text{'ball'}) \end{cases}$

*Adding negation to Datalog is not straightforward!*

# Negation in Datalog

- Various semantics have been proposed for supporting negation in Datalog
  - In a way that makes sense
- We will look at two:
  - *Semipositive* programs (restricted)
  - *Stratified* programs (standard)

# Semipositive Programs

```
Buddy(x,y) :- Likes(x,y), ¬Parent(y,x)
```

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 relational atom
  - Guarantees domain independence
  - Are following rules safe?

```
S(x) :- T(y), Edge(z,y), !Edge(x,y)
```

?

```
S(x) :- T(y), !T(x)
```

?

# Semipositive Programs

```
Buddy(x,y) :- Likes(x,y), ¬Parent(y,x)
```

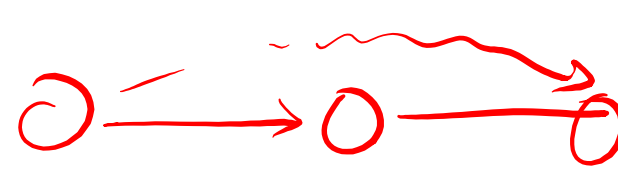
A **semipositive** program is a program where **only EDBs** may be negated

- Semantics: same as ordinary Datalog programs
- **Safety**: rule is safe if every variable occurs in a positive relational atom
  - Guarantees domain independence
  - Are following rules safe?

~~$S(x) :- T(y), \text{Edge}(z,y), \neg \text{Edge}(x,y)$~~

~~$S(x) :- T(y), \neg T(x)$~~

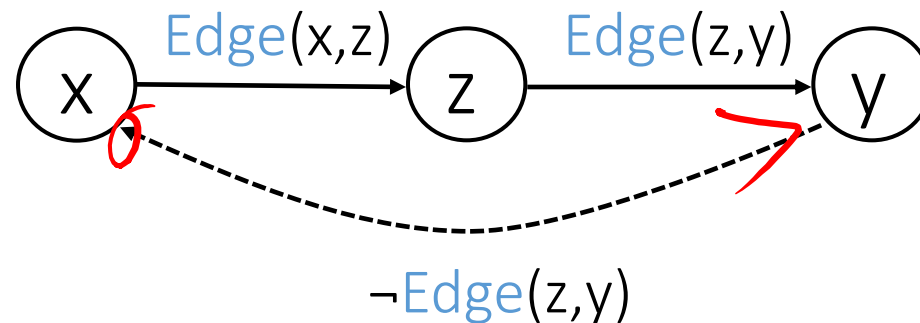
# Negated Atoms



$E(z', yz)$

- We may put ! (NOT) in front of an atom, to negate its meaning.
- Example: For any given node  $x$ , return all nodes  $y$  that are two hops away, where  $y$  is not an immediate neighbor of  $x$ .

$\text{TwoHop}(x,y) :- \text{Edge}(x,z), \text{Edge}(z,y), \neg \text{Edge}(x,y)$



Alternative notations:  $\neg R(\text{args})$  instead of  $!R(\text{args})$



# Stratified Programs

- Let  $P$  be a Datalog program and  $E_0$  be set of EDB predicates
- A *stratification* of  $P$  is a partitioning of the IDB predicates into disjoint sets  $E_1, \dots, E_k$  where:
  - For  $i=1, \dots, k$ , every rule with head in  $E_i$  has body predicates only from  $E_0, \dots, E_i$
  - For  $i=1, \dots, k$ , every rule with head in  $E_i$  can have *negated* body predicates only from  $E_0, \dots, E_{i-1}$

# Practice: unconnected nodes



Compute all pairs of disconnected nodes in a graph.

# Practice: unconnected nodes



Compute all pairs of disconnected nodes in a graph.

```
Reachable(x,y) :- Edge(x,y)
```

```
Reachable(x,y) :- Edge(x,z), Reachable(z,y)
```

# Practice: unconnected nodes

*Node is basically ADom:  
Node(x) :- Edge(x,y)  
Node(y) :- Edge(x,y)*

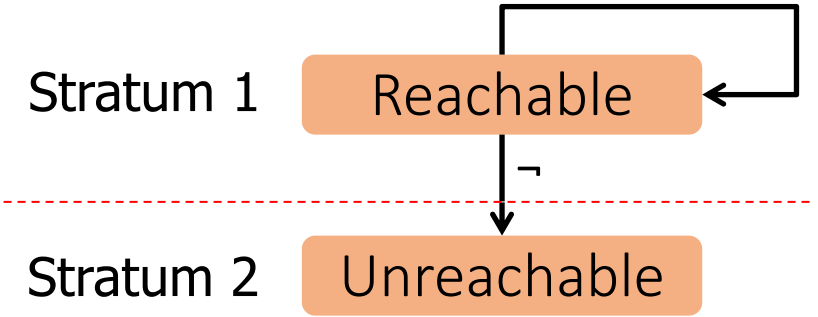


Compute all pairs of disconnected nodes in a graph.

Reachable(x,y) :- Edge(x,y)

Reachable(x,y) :- Edge(x,z), Reachable(z,y)

Unreachable(x,y) :- Node(x), Node(y), !Reachable(x,y)



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

*think: "topological sort"*

*Precedence graph :*

- Nodes = IDB predicates.
- Edge (p,q) if predicate q depends on p.
- Label this arc "!" if the predicate p is negated.

# Practice: unconnected nodes

*Node is basically ADom:  
Node(x) :- Edge(x,y)  
Node(y) :- Edge(x,y)*

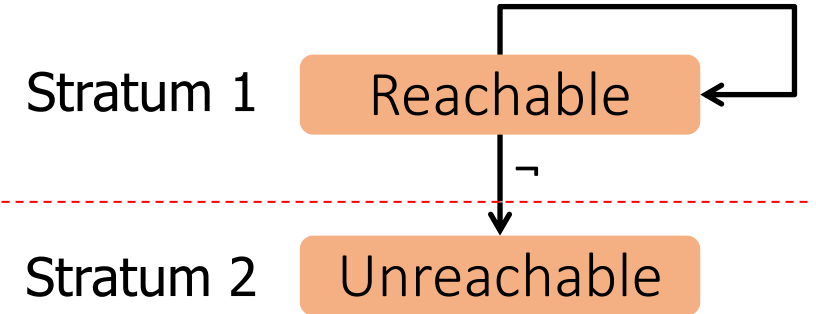


Compute all pairs of disconnected nodes in a graph.

Reachable(x,y) :- Edge(x,y)

Reachable(x,y) :- Edge(x,z), Reachable(z,y)

Unreachable(x,y) :- Node(x), Node(y), !Reachable(x,y)



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

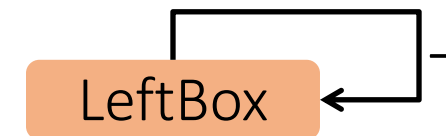
*think: "topological sort"*

- Non-stratified example:

LeftBox(x) :- ¬LeftBox(x), Item(x)

*Precedence graph :*

- Nodes = IDB predicates.
- Edge (p,q) if predicate q depends on p.
- Label this arc "!" if the predicate p is negated.



# Example

```
Person(id, gender, country)
Fake(id)
Parent(parent, child)
Spouse(person1, person2)
Likes(person1, person2)
```

```
RealPerson(x) :- Person(x,y,z),!Fake(x)
```

```
Relative(x,x) :- RealPerson(x)
```

```
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)
```

```
Buddy(x,y) :- !Relative(x,y),Likes(x,y)
```

```
Buddy(x,y) :- !Relative(x,y),Buddy(x,z),Buddy(z,y)
```

# A possible stratification

```
Person(id, gender, country)
Fake(id)
Parent(parent, child)
Spouse(person1, person2)
Likes(person1, person2)
```

1 **RealPerson**(x) :- **Person**(x,y,z),!**Fake**(x)

**Relative**(x,x) :- **RealPerson**(x)

**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)

2 **Buddy**(x,y) :- !**Relative**(x,y),**Likes**(x,y)

**Buddy**(x,y) :- !**Relative**(x,y),**Buddy**(x,z),**Buddy**(z,y)

# Another possible stratification

```
Person(id, gender, country)
Fake(id)
Parent(parent, child)
Spouse(person1, person2)
Likes(person1, person2)
```

1 `RealPerson(x) :- Person(x,y,z), !Fake(x)`

2 `Relative(x,x) :- RealPerson(x)`

`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)`

3 `Buddy(x,y) :- !Relative(x,y), Likes(x,y)`

`Buddy(x,y) :- !Relative(x,y), Buddy(x,z), Buddy(z,y)`



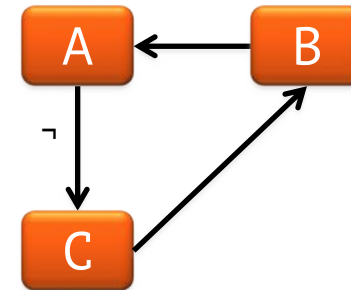
# Semantics of Stratified Programs

- For  $i=1, \dots, k$ :
  - Compute the IDBs of the stratum  $E_i$
  - Add computed IDBs to the EDBs
- Then, due to the definition of stratification, each  $E_i$  can be viewed as **semipositive**
- Does the result depend on the specific stratification of choice?
  - Answer on the next slide

# Theorems on Stratification (1)

- **THEOREM 1:** All stratifications are equivalent
  - That is, they give the same result on every input
- **THEOREM 2:** 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) :- !A(x)
```



# Theorems on Stratification (2)

- **THEOREM 3:** Non-recursive Datalog programs with negation are stratifiable
  - Via the topological order
- **THEOREM 4:** Nonrecursive Datalog with negation has the same expressive power as the algebra  $\{\sigma_{=}, \pi, \times, \rho, \cup, \setminus\}$ 
  - Extendable to RA if we add the predicates  $>, <$

# Practice



ParentChild(p,c)

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

?

# Practice



ParentChild(p,c)

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

*first compute for each  
person their descendants*

*then use negation*

?

# Practice



ParentChild(p,c)

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

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

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

*first compute for each  
person their descendants*

*then use negation*

?

# Practice



ParentChild(p,c)

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

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

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

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

*first compute for each  
person their descendants*

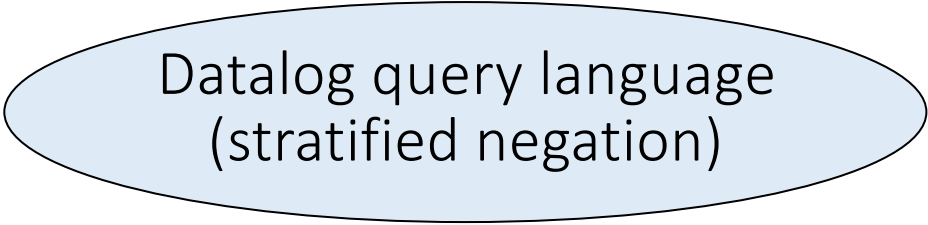
*then use negation*

# Outline: Datalog

- Datalog
  - Datalog rules
  - Recursion
  - Semantics
  - Datalog<sup>¬</sup>: Negation, stratification
  - Datalog<sup>±</sup>
  - Stable model semantics (Answer set programming)
  - Datalog vs. RA
  - Naive and Semi-naive evaluation



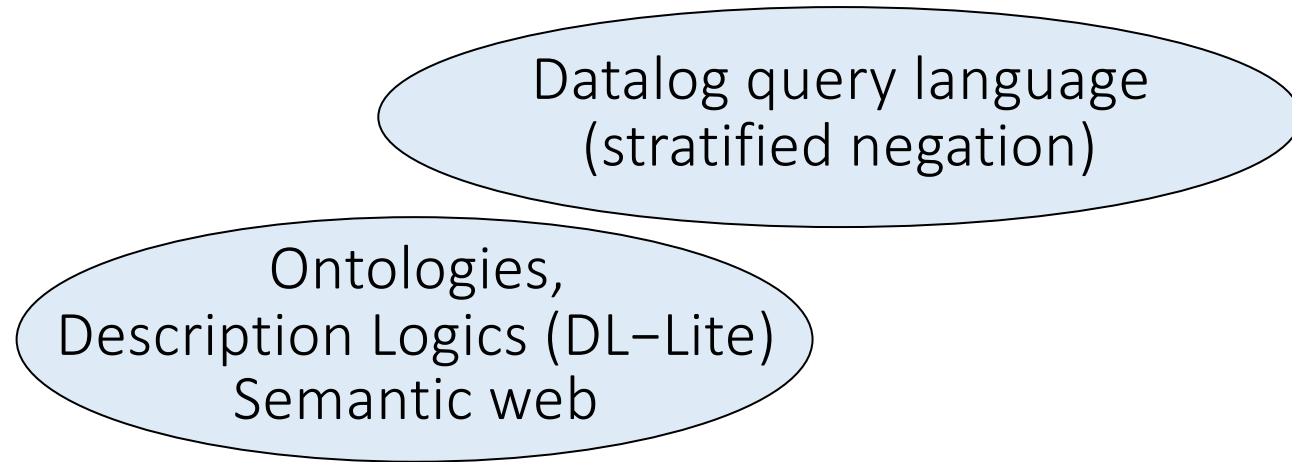
# Datalog<sup>±</sup>: background



Datalog query language  
(stratified negation)

- Much is possible with Datalog

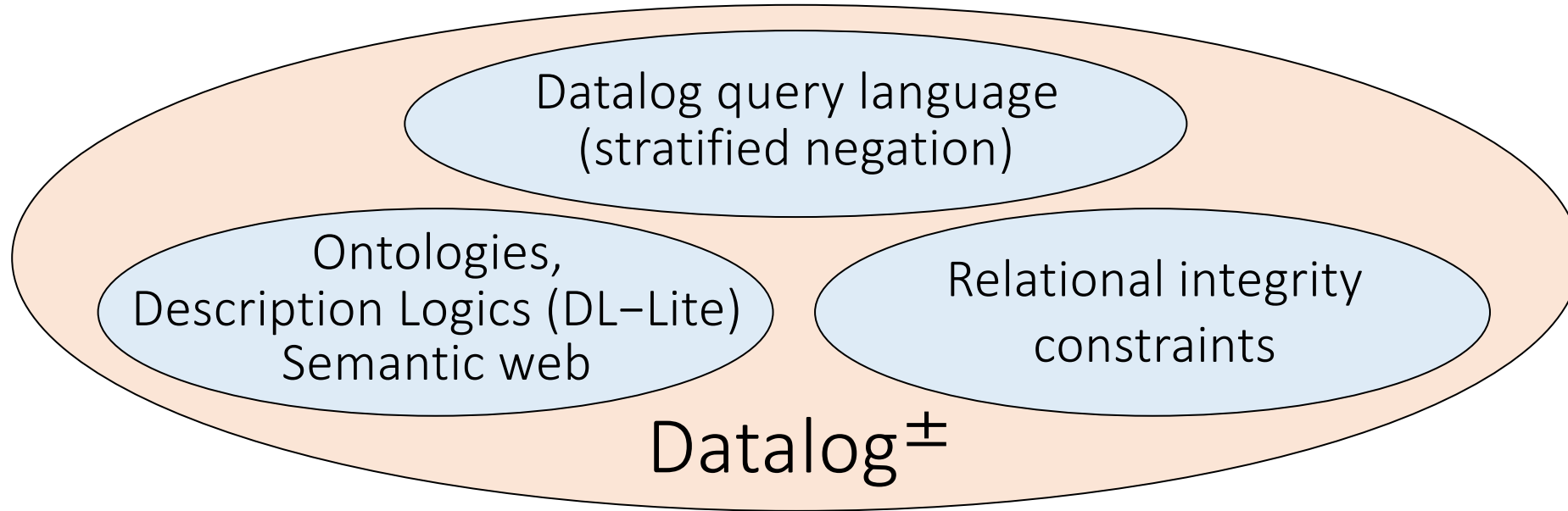
# Datalog<sup>±</sup>: background



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

Patel-Schneider, Horrocks. A comparison of two modelling paradigms in the Semantic Web. WWW (Semantic Web track). 2006.

# Datalog<sup>±</sup>: goal



- Much is possible with Datalog
- Much is not (observed e.g. by Patel-Schneider, Horrocks 2006)
- **Datalog<sup>±</sup>** 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)

Patel-Schneider, Horrocks. A comparison of two modelling paradigms in the Semantic Web. WWW (Semantic Web track). 2006.

Cali, Gottlob, Lukasiewicz, Marnette, Pieris. Datalog<sup>+/-</sup>: A Family of Logical Knowledge Representation and Query Languages for New Applications. LICS 2011.

Slide adopted from Andrea Cali

# 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 <sup>±</sup> rule
Participation	$\text{boss}(X) \rightarrow \exists Y \text{ reports}(Y, X)$
Disjointness	$\text{customer}(X), \text{boss}(X) \rightarrow \perp$
Functionality	$\text{reports}(X, Y1), \text{reports}(X, Y2) \rightarrow Y1 = Y2$

# Datalog<sup>±</sup> vs. DL

$FiveStar(X) \rightarrow Hotel(X),$   
 $FiveStar(X), \text{not} Pool(X, Y) \rightarrow \exists Z Beach(X, Z),$   
 $FiveStar(X), \text{not} Beach(X, Y) \rightarrow \exists Z Pool(X, Z),$   
 $Beach(X, Y) \rightarrow \exists Z SwimOpp(X, Z),$   
 $Pool(X, Y) \rightarrow \exists Z SwimOpp(X, Z),$

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

$FiveStar \sqsubseteq Hotel,$   
 $FiveStar \sqcap \text{not} \exists Pool \sqsubseteq \exists Beach,$   
 $FiveStar \sqcap \text{not} \exists Beach \sqsubseteq \exists Pool,$   
 $\exists Beach \sqsubseteq \exists SwimOpp,$   
 $\exists Pool \sqsubseteq \exists SwimOpp,$

# 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)

# Outline: Datalog

- Datalog
  - Datalog rules
  - Recursion
  - Semantics
  - Datalog<sup>¬</sup>: Negation, stratification
  - Datalog<sup>±</sup>
  - Stable model semantics (Answer set programming)
  - Datalog vs. RA
  - Naive and Semi-naive evaluation