

Program Analysis with Datalog

Sergey Mehtaev

s.mehtaev@pku.edu.cn

Peking University

GitHub CodeQL – Language For Defining Custom Program Analyses

JavaConverter.java

```
public static Object deserialize (InputStream is)
    throws IOException {
    ObjectInputStream ois = new ObjectInputStream(is);
    return ois.readObject();
}
```

UnsafeDeserialization.ql

```
from PathNode source, PathNode sink
where flowPath(source, sink)
select sink.getNode().(UnsafeDeserializationSink)
    .getMethodAccess(),
source, sink, "Unsafe deserialization of $@",
source.getNode(), "user input"
```

QL Query Results

alerts ▾

> ☰ Unsafe deserialization of [user input](#).

▾ ☰ Unsafe deserialization of [user input](#).

▾ Path

1 [getContent\(...\) : InputStream](#)

2 [getContentAsStream\(...\) : InputStream](#)

3 [toBufferedInputStream\(...\) : InputStream](#)

4 [getInputStream\(...\) : InputStream](#)

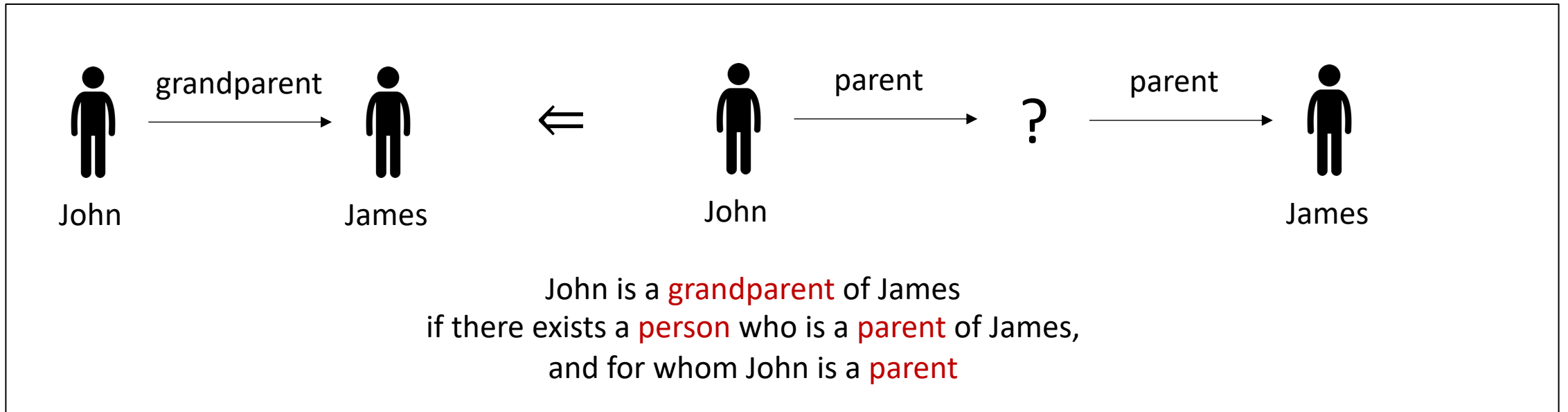
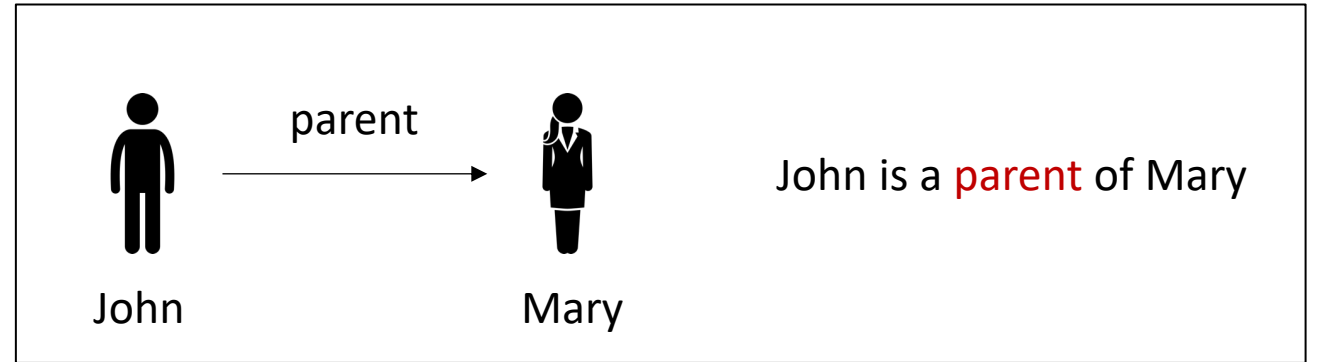
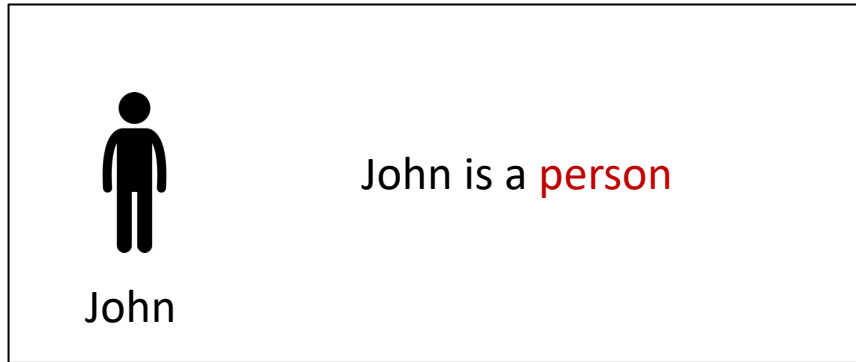
5 [is : InputStream](#)

6 [ois](#)

> Path

> ☰ Unsafe deserialization of [user input](#).

Introduction to Datalog



Introduction to Datalog

- A declarative, non-Turing complete logic programming languages
- Some notable applications:
 - Database queries, e.g. Datomic database
 - Program analysis, e.g. Doop, Semmle
 - Network protocol specification
- Formally, a subset of Horn clauses, logical formulas in the form

$$L_0 \leftarrow L_1, \dots L_n$$

Syntax

- A Datalog program is a set of Horn clauses

$$\underbrace{L_0}_{\text{Head}} \leftarrow \underbrace{L_1, \dots, L_n}_{\text{Body}}$$

- L_i is a literal symbol in the form $p_i(t_1, \dots, t_{k_i})$, where p_i is a predicate symbol and t_j are terms
- A term t_j is either a constant or a variable
- Clauses with empty bodies are facts.

Safety Conditions

Safety conditions ensure program termination:

1. Each fact in the program should not contain variables (a clause without variables is called a ground clause)
2. Each variable occurring in the head of a rule should also occur in its body

Examples

- Father-son relation

father(bob, john)

- Grandparent relation:

*grandparent(X, Y) ←
parent(X, Z),
parent(Z, Y)*

EDB & IDB

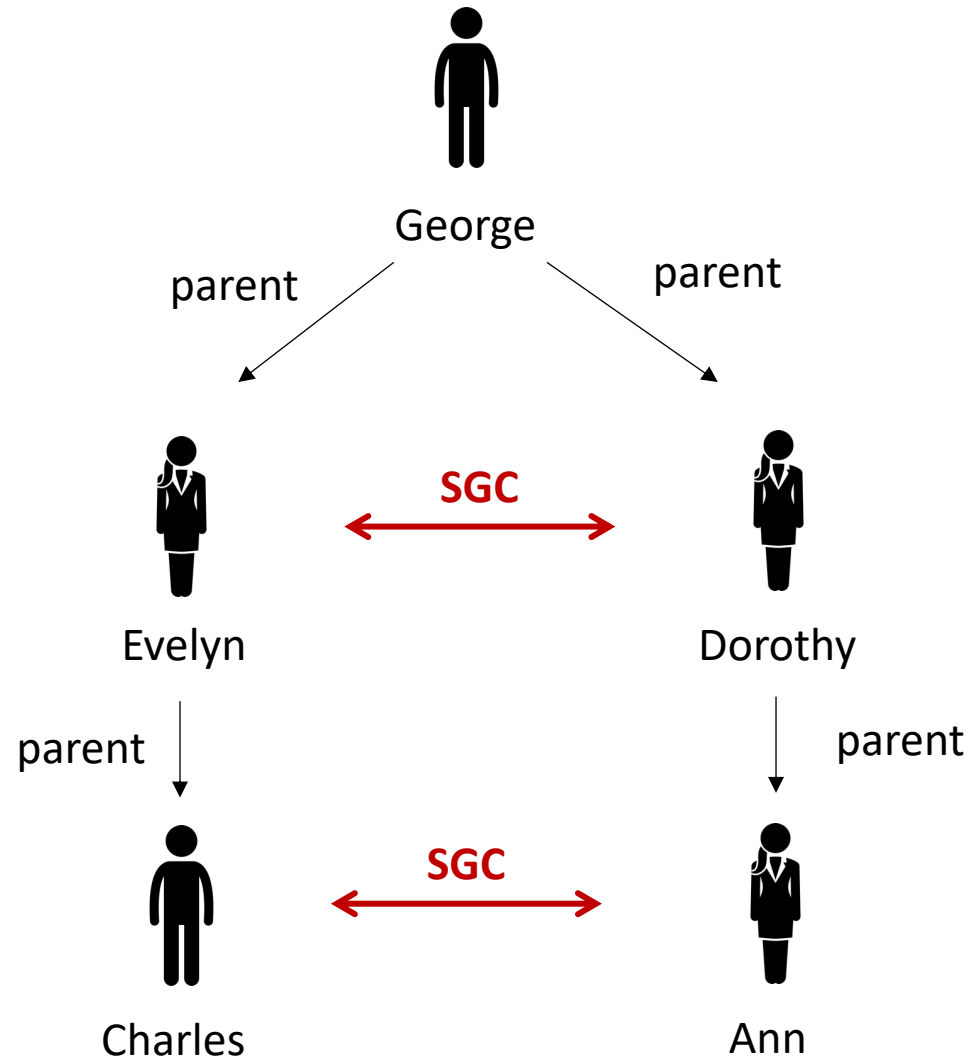
- EDB-predicates are those occurring in EDB
- IDB-predicates are those occurring in IDB, but not EDB
- **EDB** (extensional database) is a set of instances over the EDB-predicates (the input)
- **IDB** (intensional database) is a set of instances over the IDB-predicates (the output)

Example (EDB)

person(ann)
person(bertrand)
person(charles)
person(dorothy)
person(evelyn)
person(fred)
person(george)
person(hilary)

parent(george, dorothy)
parent(george, evelyn)
parent(dorothy, bertrand)
parent(dorothy, ann)
parent(hilary, ann)
parent(evelyn, charles)

Example: Same Generation Cousins (SGC)



Example (IDB)

Same generation cousins (*sgc*):

Rule 1: $sgc(X, X) \leftarrow person(X)$
Rule 2: $sgc(X, Y) \leftarrow parent(X1, X),$
 $parent(Y1, Y),$
 $sgc(X1, Y1).$

Some facts belonging to *sgc*:

$sgc(dorothy, evelyn)$
 $sgc(charles, ann)$

Elementary production principle

The meaning of Datalog program is defined through fact inference.

Elementary production principle (EPP): for a rule $L_0 \leftarrow L_1, \dots, L_n$ and a list of facts F_1, \dots, F_n , if there is a substitution θ (of variables with values) such that $L_i\theta = F_i$, then from the facts F_1, \dots, F_n we can infer the fact $L_0\theta$.

Example: from $sgc(X, X) \leftarrow person(X)$ and $person(ann)$, we can infer $sgc(ann, ann)$.

Proof-theoretic Semantics of Datalog

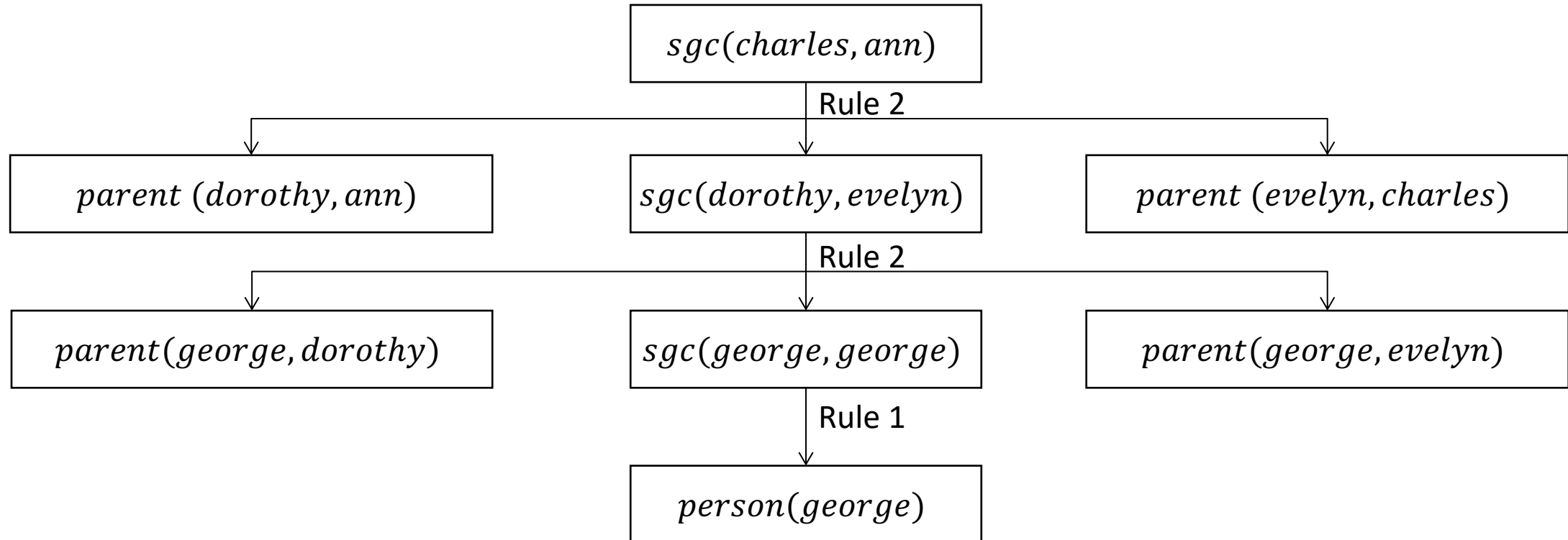
Given a set of clauses S , we define the inference relation \vdash for fact F :

1. $S \vdash F$ if $F \in S$
2. $S \vdash F$ if a rule $R \in S$ and facts F_1, \dots, F_n exists such that $S \vdash F_i$ for all $i \in [1..n]$ and F is inferred by applying EPP to R and F_1, \dots, F_n .

Theorem. All facts that can be inferred from S are exactly the facts that logically follow from S .

Proof Tree

A **proof tree** for $sgc(charles, ann)$ visualises \vdash relationship:



Negation

$\neg sgc(bertrand, X)$ – X is not a same generation cousin of Bertrand

Negation is incorporated based on **Closed World Assumption (CWA)**: if a fact does not logically follow from a set of Datalog clauses, then we conclude that the negation of this fact is true.

Safety: every variable appearing in a negated literal must also appear in a positive literal of the same rule. Negation in the head of a rule is not allowed.

Problem of Negation with Recursion

IDB:

$$\begin{aligned}r(X) &\leftarrow s(X), \neg t(X) \\t(X) &\leftarrow s(X), \neg r(X)\end{aligned}$$

EDB:

$$s(1)$$

In this case, we have two solutions according to the fixpoint semantics:

1. $cons(S) := \{r(1), \neg t(1)\}$
2. $cons(S) := \{\neg r(1), t(1)\}$

Stratified Datalog

A stratified Datalog program can be partitioned into disjoint sets of clauses $P = P^1 \cup \dots \cup P^n$ called strata, so that

1. Each IDB predicate has its defining clauses within one stratum
2. P^1 can contains only negation of EDB predicates
3. P^i can only contain negation of IDB predication from $P^j, j < i$

Ensures that the semantics with negation is well-defined.

Stratified Datalog Example

$r_1: p(X, Y) \leftarrow \neg q(X, Y), s(X, Y)$

$r_2: q(X, Y) \leftarrow q(X, Z), q(Z, Y)$

$r_3: q(X, Y) \leftarrow d(X, Y), \neg r(X, Y)$

$r_4: r(X, Y) \leftarrow d(Y, X)$

$r_5: s(X, Y) \leftarrow q(X, Y), q(Y, T), X \neq T$

Strata:

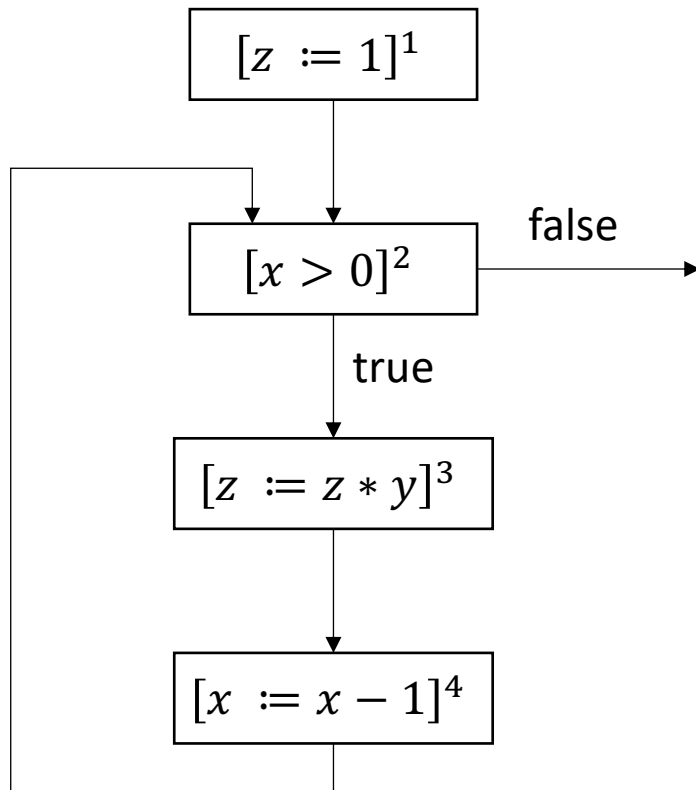
$P^1 = \{r_4\}$

$P^2 = \{r_2, r_3, r_5\}$

$P^3 = \{r_1\}$

EDB predicates: d

CFG As EDB



label(1)
label(2)
label(3)
label(4)
init(1)
final(2)
flow(1,2)
flow(2,3)
flow(3,4)
flow(4,2)

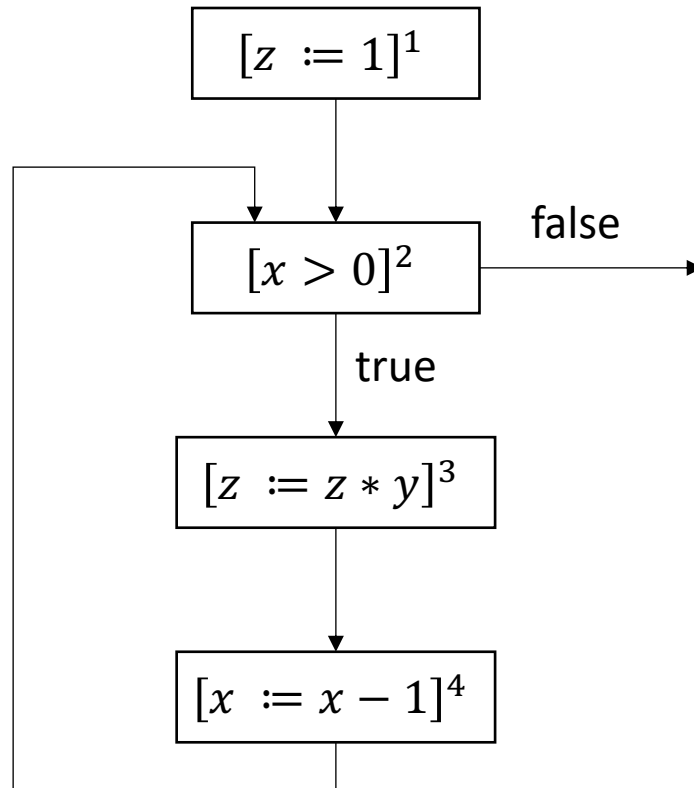
CFG Reachability

The relation $reach(L_1, L_2)$ states that there is a path from L_1 to L_2 in CFG:

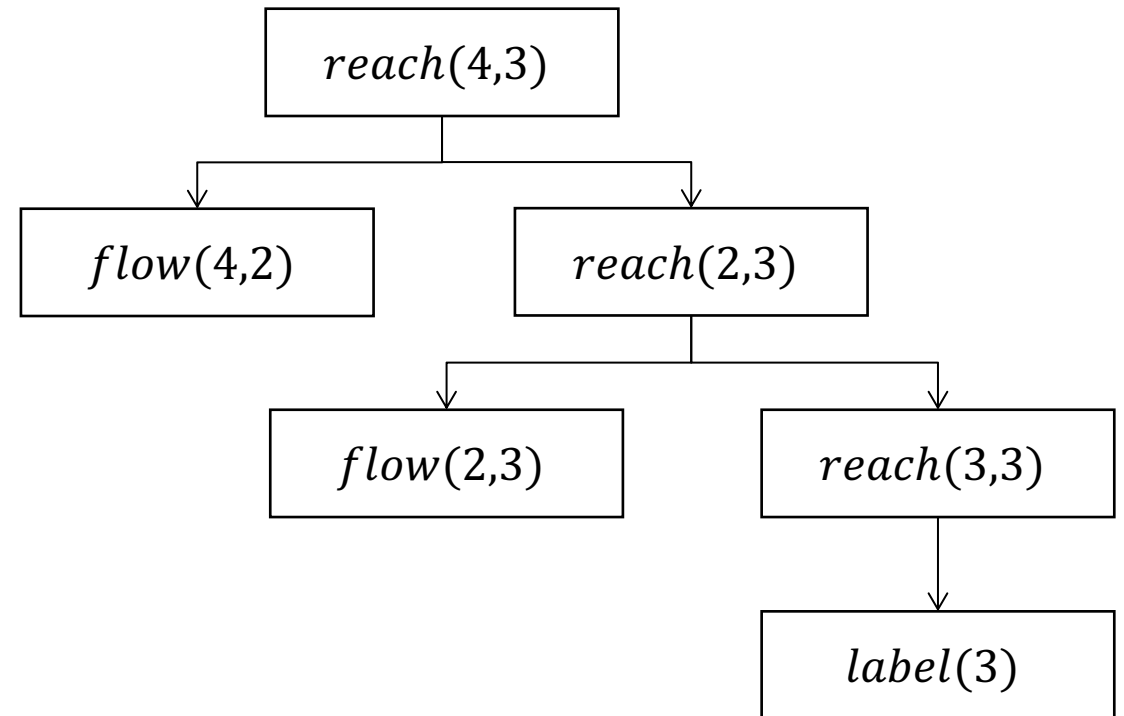
$$\begin{aligned} reach(X, X) &\leftarrow label(X) \\ reach(X, Y) &\leftarrow flow(X, Z), reach(Z, Y) \end{aligned}$$

CFG Reachability

CFG



Proof Tree of $reach(4, 3)$



Dominators (Part 1)

dominate(D, N): A node D **dominates** a node N if every path from the entry node to N must go through D .

$reach_without(L, L, D) \leftarrow$
 $label(L),$
 $label(D),$
 $L \neq D$

$reach_without(L_1, L_2, D) \leftarrow$
 $L_1 \neq D,$
 $flow(L_1, L_3),$
 $reach_without(L_3, L_2)$

Dominators (Part 2)

dominate(D, N): A node D **dominates** a node N if every path from the entry node to N must go through D .

```
not_dominates( $D, N$ ) ←  
  init( $L$ ),  
  reach_without( $L, N, D$ )  
dominate( $D, N$ ) ←  
   $D \neq N$ ,  
  init( $L$ ),  
  reach( $L, N$ ),  
  label( $D$ ),  
  !not_dominates( $D, N$ )
```

The role of negation

Negation helps to express a property that hold for a set of paths.

$$\begin{aligned} \textit{dominate}(D, N) \leftarrow & \\ & D \neq N, \\ & \textit{init}(L), \\ & \textit{reach}(L, N), \\ & \textit{label}(D), \\ & \textit{!not_dominate}(D, N) \end{aligned}$$

D dominates N if there **exists** a path from an init node to N, and it is **not** true that there **exists** a path from an init node to N that does **not** visit D.

Reaching Definitions

Reaching definitions analysis determines for each program point, which assignments may have been made and not overwritten, when program execution reaches this point along some path.

```
[x := 5]1;  
[y := 1]2;  
while [x > 1]3 do  
    [y := x * y]4;  
    [x := x - 1]5;
```

All assignments reach the entry of 4; only the assignments 1,4,5 reach the entry of 5.

EDB for Reaching Definitions

Domains

- program variables (V_i)
- labels (L_i)

Relations

label(L)
flow(L₁, L₂)
init(L)
final(L)
variable(V)
defined(V, L)

Killed Assignments

$$\mathit{kill}_{RD}([x := a]^l) = \{(x, ?)\} \cup \{(x, l') \mid B^{l'} \text{ is an assignment to } x\}$$

$$\mathit{kill}_{RD}([\text{skip}]^l) = \emptyset$$

$$\mathit{kill}_{RD}([b]^l) = \emptyset$$

where ? is the label for uninitialised variables.

$$\begin{aligned} \mathit{kill}(L, V, D) \leftarrow \\ & \mathit{defined}(V, L), \\ & \mathit{defined}(V, D), \\ & L \neq D \end{aligned}$$

$$\begin{aligned} \mathit{kill}(L, V, ?) \leftarrow \\ & \mathit{defined}(V, L) \end{aligned}$$

Generated Assignments

$$gen_{RD}([x := a]^l) = \{(x, l)\}$$

$$gen_{RD}([skip]^l) = \emptyset$$

$$gen_{RD}([b]^l) = \emptyset$$

$$gen(L, V, L) \leftarrow defined(V, L)$$

Analysis

$$RD_{entry}(l) = \begin{cases} \{(x, ?) \mid x \in Vars\} & \text{if } l = \text{init}(\text{program}) \\ \cup \{RD_{exit}(l') \mid (l', l) \in \text{flow}(\text{program})\} & \text{otherwise} \end{cases}$$

$$RD_{exit}(l) = \left(RD_{entry}(l) \setminus \text{kill}_{RD}(B^l) \right) \cup \text{gen}_{RD}(B^l)$$

rd_entry(L, V, ?) ← init(L)

rd_entry(L, V, D) ← flow(L', L), rd_exit(L', V, D)

rd_exit(L, V, D) ← rd_entry(L, V, D), !kill(L, V, D)

rd_exit(L, V, D) ← gen(L, V, D)