

# 一、引入

前面聚集的点在算法思想上，但终归，我们要将其落实到具体的计算机系统上。

Goal: select adults from a set of persons

- Imperative: how to do (~implementation)

```
Set<Person> selectAdults(Set<Person> persons) {
    Set<Person> result = new HashSet<>();
    for (Person person : persons)
        if (person.getAge() >= 18)
            result.add(person);
    return result;
}
```

命令式

- Declarative: what to do (~specification)

```
SELECT * FROM Persons WHERE Age >= 18;
```

声明式

前者对程序设计的粒度更细，掌控的细节更多，但同时也意味着繁琐；后者的层次更高，面向逻辑，脏活都被系统规范给封装好了。

对于PTA等静态分析算法，两种方法的工作量是完全不同的。

## Pointer Analysis, Imperative Implementation

```
Solve(mentry)
WL=[], PFG={}, S={}, RM={}, CG={}
AddReachable(mentry)
while WL is not empty do
    remove (n, pts) from WL
    Δ = pts - p(n)
    Propagate(n, Δ)
    if n represents a variable x then
        foreach oi ∈ Δ do
            foreach x.f = y ∈ S do
                AddEdge(y, oi, f)
        foreach y = x.f ∈ S do
            ...
```

AddReachable(m)  
if m ∈ RM then  
 add m to PFG

- How to implement **worklist**?
  - Array list or linked list?
  - Which worklist entry should be processed first?
- How to implement **points-to set** ( $p()$ )?
  - Hash set or bit vector?
- How to connect **PFG nodes** and pointers?
- How to associate variables to the relevant statements?

```
AddEdge(s, t)
if s → t ∉ PFG then
    add s → t to PFG
    if pt(s) is not empty
        add (t, pt(s)) to
```

Propagate  
if pts is  
pt(n)  
...  
So many implementation details



## Pointer Analysis, Declarative Implementation (via Datalog)

```
VarPointsTo(x, o) :-  
    Reachable(o),  
    New(x, o, m).  
  
VarPointsTo(x, o) :-  
    Assign(x, y),  
    VarPointsTo(y, o).  
  
FieldPointsTo(o1, f, oj) :-  
    Store(x, f, y),  
    VarPointsTo(x, o1),  
    VarPointsTo(y, oj).  
  
VarPointsTo(y, oj) :-  
    Load(y, x, f),  
    VarPointsTo(x, o1),  
    FieldPointsTo(o1, f, oj).  
  
...  
  
• Succinct  
• Readable (logic-based specification)  
• Easy to implement
```



```
VarPointsTo(this, o),  
Reachable(o),  
CallGraph(1, m) :->  
    VCall(1, i, k),  
    VarPointsTo(x, o),  
    Dispatch(o, k, m),  
    ThisVar(m, this).
```

```
VarPointsTo(pi, o) :->  
    CallGraph(1, m),  
    Argument(1, i, ai),  
    Parameter(m, i, pi),  
    VarPointsTo(ai, o).
```

```
VarPointsTo(or, o) :->  
    CallGraph(1, m),  
    MethodReturn(m, ret),  
    VarPointsTo(ret, o),  
    CallReturn(1, r).
```

## 二. Datalog 介绍

### 1. 基础知识介绍

- Datalog is a **declarative logic** programming language that is a subset of **Prolog**.
- It emerged as a database language (mid-1980s)\*
- Now it has a variety of applications
  - Program analysis
  - Declarative networking
  - Big data
  - Cloud computing
  - ...

\*David Maier, K. Tuncay Tekle, Michael Kifer, and David S. Warren,  
“*Datalog: Concepts, History, and Outlook*”. Chapter, 2018.

Datalog = **Data + Logic**

(and, or, not)

- No side-effects 赋值, 改变.
- No control flows
- No functions
- Not Turing-complete 非图灵完备

## 2. 谓词

用来描述若干项之间的联系,  $R: D^n \rightarrow \{0, 1\}$ 。  
由于论域  $D$  的离散性, 其在某种程度下可以  
视为全组合论域  $D^n$  中所有变元组针对某个  
agenda 的陈述集(有真有假)

Ps: 为真的陈述被称为 fact

### 3. 原子公式

谓词 + 项构成原子公式  $\rightarrow P(t_1, t_2, \dots, t_n)$

在形式化中，将使用以下四类符号。

(1) 常量符号： $a, b, c, \dots, a_i, b_i, c_i, \dots, i \geq 1$ ，当论域D给出时，它可以是D中的某个元素；

(2) 变量符号： $x, y, z, \dots, x_i, y_i, z_i, \dots, i \geq 1$ ，当论域D给出时，它可以是D中的任何一个元素；

(3) 函数符号： $f, g, h, \dots, f_i, g_i, h_i, \dots, i \geq 1$ ，当论域D给出时， $n$ 元函数符号  $f(x_1, \dots, x_n)$  可以是  $D^n$  到 D 的任意一个映射；

(4) 谓词符号： $P, Q, R, \dots, P_i, Q_i, R_i, \dots, i \geq 1$ ，当论域D给出时， $n$ 元谓词符号  $P(x_1, \dots, x_n)$  可以是  $D^n$  到 {1, 0} 的任意一个谓词。

定义1 一阶逻辑中的项(item) 被递归定义如下：

(1) 常量符号是项；

(2) 变量符号是项；

(3) 若  $f(x_1, \dots, x_n)$  是  $n$  元函数符号， $t_1, \dots, t_n$  是项，则  $f(t_1, \dots, t_n)$  是项；

(4) 只有有限次地使用(1)、(2)、(3)所生成的符号串才是项。

例如  $a, b, x, y$  是项， $f(x, y) = x + y, g(x, y) = x \cdot y$  是项，  
 $f(a, g(x, y)) = a + x \cdot y$  也是项。

定义2 设  $P(x_1, \dots, x_n)$  是  $n$  元谓词， $t_1, \dots, t_n$  是项，则称  $P(t_1, \dots, t_n)$  为原子公式，或简称原子。<sup>[2]</sup>

① 使用算术谓词来描述项之间的关系，在 Datalog 中被称为 arithmetic atom

- In addition to relational atoms, Datalog also has arithmetic atoms

- E.g., age >= 18

② 使用其它特殊的谓词来描述项之间的关系，被称为 relational atom

• Age("Xiaoming", 18) is true

• Age("Alan", 23) is false

Age	person	age
	Xiaoming	18
	Xiaohong	23
	Alan	16
	Abao	31

本质上都是  $D^n$  谓词  $\rightarrow f_0, 13$ ，是 Datalog 中的基本单元

## 4. 逻辑推断

我们前面讲到的一些概念中，论域作为承载所有谓词的域，取决于需求和真实情况，谓词作为定义在其上的函数，其产生的不同事实分布描述了论域的一种关系性结构。推断实际上就是谓词在论域上的推导。

$H \leftarrow B_1, B_2, \dots, B_n$

Head (consequent)

$H$  is an atom

Body (antecedent)

$B_i$  is a (possibly negated) atom

Each  $B_i$  is called a subgoal

$H(X_1, X_2) \leftarrow B_1(X_1, X_3), B_2(X_2, X_4), \dots, B_n(X_m)$ .

抽象地来看，其意味着当且仅当所有  $B_i$  为真时， $H$  为真。但由于  $B_i$  与  $H$  都是原子公式，所以这种过程实际上是由具体的事实在描述构建的。结果是，我们可以得到一种全新的、对论域的结构性认识。

Datalog program = Facts + Rules

Age	person	age
	Xiaoming	18
	Xiaohong	23
	Alan	16
	Abao	31



$\text{Adult}(\text{person}) \leftarrow \text{Age}(\text{person}, \text{age}), \text{age} \geq 18.$



Adult	person
	Xiaoming
	Xiaohong
	Abao

一些新的Facts

因此我们发现 Datalog 实际上就是在论域上，  
基于已有的谓词和 Rules 进行推导，得到一些  
新的谓词（以及其用在论域上的结果）

Conventionally, predicates in Datalog are divided  
into two kinds:

### 1. EDB (extensional database)

- The predicates that are defined in a priori
- Relations are immutable → 不可变
- Can be seen as input relations

### 2. IDB (intensional database)

- The predicates that are established only by rules
- Relations are inferred by rules
- Can be seen as output relations

$$H \leftarrow B_1, B_2, \dots, B_n.$$

- H can only be IDB
- B<sub>i</sub> can be EDB or IDB

## ★ Datalog By Rules

(1) 与 ,

(2) 或

SportFan(person)  $\leftarrow$   
Hobby(person, "jogging");  
Hobby(person, "swimming").

Ps: 与的运算优先级要

高于或

分成两行去写

SportFan(person)  $\leftarrow$  Hobby(person, "jogging").  
SportFan(person)  $\leftarrow$  Hobby(person, "swimming").

(3) 非 !

MakeupExamStd(student)  $\leftarrow$   
Student(student),  
!PassedStd(student).

Where Student stores all students,  
and PassedStd stores the students  
who passed the exam.

(4) 递归 Recursion

也即 IDB 可由其自身直接/间接的推断得出。  
其使得 DataLog 超越了操作简单代数的查  
询语言 (SQL)，而能表达更加复杂的语义特性。

- For example, we can compute the reachability information (i.e., transitive closure) of a graph with recursive rules:

Reach(from, to)  $\leftarrow$   
Edge(from, to).

Reach(from, to)  $\leftarrow$   
Reach(from, node),  
Edge(node, to).

Where Edge(a, b) means that the graph  
has an edge from node a to node b,  
and Reach(a, b) means that b is  
reachable from a.

Ps: ① 变量域约束

如果论域无穷 + 推断结构不闭合  
 $\Downarrow$

可能需要枚举无穷个项来判断结果

A(x)  $\leftarrow$  B(y),  $x > y$ .

For both rules, infinite values of x can satisfy the rule,  
which makes A an *infinite relation*.

A(x)  $\leftarrow$  B(y), !C(x, y).

因此我们规定一个 rule 是 safe 的当且仅当其  
所涉及到的变量至少被一个 non-negated relational  
atom 所约束。

即使一个变量的论域是无穷的，施加一个  
non-negated 的关系型原子公式后也变成有穷的了

## ② 循环消元

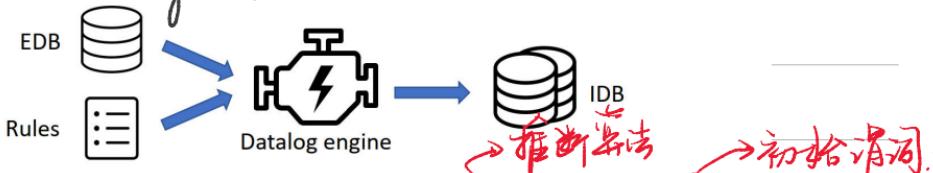
$A(x) \leftarrow B(x), \neg A(x)$ .

某个原子公式成立当且仅当其

取反成立本身就很滑稽。

因此规则二就是避免对一个原子公式同时进行取反与递归。

## 3. Datalog 程序的执行



- Datalog engine deduces facts by given rules and EDB predicates until no new facts can be deduced. Some modern Datalog engines  
LogicBlox, Soufflé, XSB, Datomic, Flora-2, ...
  - Monotonicity: Datalog is monotone as facts cannot be deleted
  - Termination: A Datalog program always terminates as
    - 1) Datalog is monotone 单调
    - 2) Possible values of IDB predicates are finite (rule safety) 可得到的IDB谓词有限?
- 其可终止性某种程度上就揭示了 Datalog 并不是一种图灵完备的语言

### 三. DataLog 的应用

#### 1. 指针分析

首先确定初始的、已想穿的谓词。

Kind	Statement
New	$i : x = \text{new } T()$
Assign	$x = y$
Store	$x.f = y$
Load	$y = x.f$

EDB

New( $x : V, o : O$ )

Assign( $x : V, y : V$ )

Store( $x : V, f : F, y : V$ )

Load( $y : V, x : V, f : F$ )

Variables: V  
Fields: F  
Objects: O

IDB

VarPointsTo( $v : V, o : O$ )

e.g., fact VarPointsTo( $x, o_i$ ) represents  $o_i \in pt(x)$

FieldPointsTo( $oi : O, f : F, oj : O$ )

e.g., fact FieldPointsTo( $o_i, f, o_j$ ) represents  $o_j \in pt(o_i.f)$

然后基于程序流域，得到 EDB 初始谓词的事实域。

#### An Example

```

1 b = new C();
2 a = b;
3 c = new C();
4 c.f = a;
5 d = c;
6 c.f = d;
7 e = d.f;
  
```

Variables: V  
Fields: F  
Objects: O

New( $x : V, o : O$ )

New

b	$o_1$
c	$o_3$

域。

Assign( $x : V, y : V$ )

Assign

a	b
d	c

Store( $x : V, f : F, y : V$ )

Store

c	f	a
c	f	d

Load( $x : V, y : V, f : F$ )

Load

e	d	f
---	---	---

# Datalog Rules for Pointer Analysis

Kind	Statement	Rule
New	$i : x = \text{new } T()$	$\frac{}{o_i \in pt(x)}$
Assign	$x = y$	$\frac{o_i \in pt(y)}{o_i \in pt(x)}$
Store	$x.f = y$	$\frac{o_i \in pt(x) \quad o_j \in pt(y)}{o_j \in pt(o_i.f)}$
Load	$y = x.f$	$\frac{o_i \in pt(x) \quad o_j \in pt(o_i.f)}{o_j \in pt(y)}$

$\text{VarPointsTo}(x, o) \leftarrow \text{New}(x, o).$   
 $\text{VarPointsTo}(x, o) \leftarrow \text{Assign}(x, y), \text{VarPointsTo}(y, o).$   
 $\text{FieldPointsTo}(oi, f, oj) \leftarrow \text{Store}(x, f, y), \text{VarPointsTo}(x, oi), \text{VarPointsTo}(y, oj).$   
 $\text{VarPointsTo}(y, oj) \leftarrow \text{Load}(y, x, f), \text{VarPointsTo}(x, oi), \text{FieldPointsTo}(oi, f, oj).$

之后根据推断规则集，在EDB的事实域的基础上得到IDB的事实域。

## An Example

```

1 b = new C();
2 a = b;
3 c = new C();
4 c.f = a;
5 d = c;
6 c.f = d;
7 e = d.f;

```

```

VarPointsTo(x, o) <- New(x, o).

VarPointsTo(x, o) <- Assign(x, y),
VarPointsTo(y, o).

FieldPointsTo(oi, f, oj) <- Store(x, f, y),
VarPointsTo(x, oi),
VarPointsTo(y, oj).

VarPointsTo(y, oj) <- Load(y, x, f),
VarPointsTo(x, oi),
FieldPointsTo(oi, f, oj).

```

<b>New(x:V, o:O)</b>	<b>Store(x:V, f:F, y:V)</b>																		
<b>New</b>	<b>Store</b>																		
<table border="1"> <tr> <td>b</td><td><math>o_1</math></td></tr> <tr> <td>c</td><td><math>o_3</math></td></tr> </table>	b	$o_1$	c	$o_3$	<table border="1"> <tr> <td>c</td><td>f</td><td>a</td></tr> <tr> <td>c</td><td>f</td><td>d</td></tr> </table>	c	f	a	c	f	d								
b	$o_1$																		
c	$o_3$																		
c	f	a																	
c	f	d																	
<b>Assign(x:V, y:V)</b>	<b>Load(x:V, y:V, f:F)</b>																		
<b>Assign</b>	<b>Load</b>																		
<table border="1"> <tr> <td>a</td><td>b</td></tr> <tr> <td>d</td><td>c</td></tr> </table>	a	b	d	c	<table border="1"> <tr> <td>e</td><td>d</td><td>f</td></tr> </table>	e	d	f											
a	b																		
d	c																		
e	d	f																	
<b>VarPointsTo(v:V, o:O)</b>	<b>FieldPointsTo</b>																		
<b>VarPointsTo</b>	<b>FieldPointsTo</b>																		
<table border="1"> <tr> <td>b</td><td><math>o_1</math></td></tr> <tr> <td>c</td><td><math>o_3</math></td></tr> <tr> <td>a</td><td><math>o_1</math></td></tr> <tr> <td>d</td><td><math>o_3</math></td></tr> <tr> <td>e</td><td><math>o_1</math></td></tr> <tr> <td>e</td><td><math>o_3</math></td></tr> </table>	b	$o_1$	c	$o_3$	a	$o_1$	d	$o_3$	e	$o_1$	e	$o_3$	<table border="1"> <tr> <td><math>o_3</math></td><td>f</td><td><math>o_1</math></td></tr> <tr> <td><math>o_3</math></td><td>f</td><td><math>o_3</math></td></tr> </table>	$o_3$	f	$o_1$	$o_3$	f	$o_3$
b	$o_1$																		
c	$o_3$																		
a	$o_1$																		
d	$o_3$																		
e	$o_1$																		
e	$o_3$																		
$o_3$	f	$o_1$																	
$o_3$	f	$o_3$																	
<b>FieldPointsTo(oi:O, f:F, oj:O)</b>																			

上面是单方法，不含调用的例子。下面我们将讨论如何处理方法调用。

Kind	Statement	Rule
Call	$L: r = x.k(a_1, \dots, a_n)$	$o_i \in pt(x), m = \text{Dispatch}(o_i, k) \rightarrow o_u \in pt(aj), 1 \leq j \leq n$ $\frac{o_v \in pt(m_{ret})}{\rightarrow o_i \in pt(m_{this})}$ $\rightarrow o_u \in pt(m_{pj}), 1 \leq j \leq n$ $o_v \in pt(r)$

EDB

- VCall(l:S, x:V, k:M)
- Dispatch(o:O, k:M, m:M)
- ThisVar(m:M, this:V)

IDB

- Reachable(m:M)
- CallGraph(l:S, m:M)

VarPointsTo(this, o),  
 Reachable(m),  
 CallGraph(l, m)  $\leftarrow$   
 VCall(l, x, k),  
 VarPointsTo(x, o),  
 Dispatch(o, k, m),  
 ThisVar(m, this).

Statements S  
 (Labels):

Methods: M

Tian Tan @ Nanjing University

Kind	Statement	Rule
Call	$L: r = x.k(a_1, \dots, a_n)$	$o_i \in pt(x), m = \text{Dispatch}(o_i, k) \rightarrow o_u \in pt(aj), 1 \leq j \leq n$ $\frac{o_v \in pt(m_{ret})}{\rightarrow o_i \in pt(m_{this})}$ $\rightarrow o_u \in pt(m_{pj}), 1 \leq j \leq n$ $o_v \in pt(r)$

EDB

- Argument(l:S, i:N, ai:V)
- Parameter(m:M, i:N, pi:V)

VarPointsTo(pi, o)  $\leftarrow$   
 CallGraph(l, m),  
 Argument(l, i, ai),  
 Parameter(m, i, pi),  
 VarPointsTo(ai, o).

Statements S  
 (Labels):

Methods: M

Nature numbers N  
 (indexes)

Tian Tan @ Nanjing University

Kind	Statement	Rule
Call	$L: r = x.k(a_1, \dots, a_n)$	$o_i \in pt(x), m = \text{Dispatch}(o_i, k) \rightarrow o_u \in pt(aj), 1 \leq j \leq n$ $\frac{o_v \in pt(m_{ret})}{\rightarrow o_i \in pt(m_{this})}$ $\rightarrow o_u \in pt(m_{pj}), 1 \leq j \leq n$ $\rightarrow o_v \in pt(r)$

EDB

- MethodReturn(m:M, ret:V)
- CallReturn(l:S, r:V)

VarPointsTo(r, o)  $\leftarrow$   
 CallGraph(l, m),  
 MethodReturn(m, ret),  
 VarPointsTo(ret, o),  
 CallReturn(l, r).

```
Reachable(m) <-
EntryMethod(m).
```

```
VarPointsTo(x, o) <-
Reachable(m),
New(x, o, m).
```

```
VarPointsTo(x, o) <-
Assign(x, y),
VarPointsTo(y, o).
```

```
FieldPointsTo(oi, f, oj) <-
Store(x, f, y),
VarPointsTo(x, oi),
VarPointsTo(y, oj).
```

```
VarPointsTo(y, oj) <-
Load(y, x, f),
VarPointsTo(x, oi),
FieldPointsTo(oi, f, oj).
```

```
VarPointsTo(this, o),
Reachable(m),
CallGraph(l, m) <-
VCall(l, x, k),
VarPointsTo(x, o),
Dispatch(o, k, m),
ThisVar(m, this).
```

```
VarPointsTo(pi, o) <-
CallGraph(l, m),
Argument(l, i, ai),
Parameter(m, i, pi),
VarPointsTo(ai, o).
```

```
VarPointsTo(r, o) <-
CallGraph(l, m),
MethodReturn(m, ret),
VarPointsTo(ret, o),
CallReturn(l, r).
```

## 2. 污染分析

- EDB predicates

- Source(m: M) // source methods
- Sink(m: M, i: N) // sink methods
- Taint(l: S, t: T) // associates each call site to the tainted data from the call site

- IDB predicate

- TaintFlow(sr: S, sn: S, i: N) // detected taint flows, e.g., TaintFlow(sr, sn, i) denotes that tainted data from source call *sr* may flow to *i*-th argument of sink call *sn*

- Handles sources (generates tainted data)

Kind	Statement	Rule
Call	$l: r = x.k(a_1, \dots, a_n)$	$l \rightarrow m \in CG$ $\underline{m \in \text{Sources}}$ $t_l \in pt(r)$

```
VarPointsTo(r, t) <-
CallGraph(l, m),
Source(m),
CallReturn(l, r),
Taint(l, t).
```

Call-site 与对应接收变量为 t

- Handles sinks (generates taint flow information)

Kind	Statement	Rule
Call	$l: r = x.k(a_1, \dots, a_n)$	$l \rightarrow m \in CG$ $\langle m, i \rangle \in \text{Sinks}$ $t_j \in pt(a_i)$ $\langle j, l, t \rangle \in \text{TaintFlows}$

```
TaintFlow(j, l, i) <-
CallGraph(l, m),
Sink(m, i),
Argument(l, i, ai),
VarPointsTo(ai, t),
Taint(j, t).
```

### 3. Datalog 的优缺点对比

- Pros

- Succinct and readable 可读性
- Easy to implement 易使用
- Benefit from off-the-shelf optimized Datalog engines 可单独优化引擎模块

- Cons

- Restricted expressiveness, i.e., it is impossible or inconvenient to express some logics 表达能力受限
- Cannot fully control performance 能控制的粒度较粗