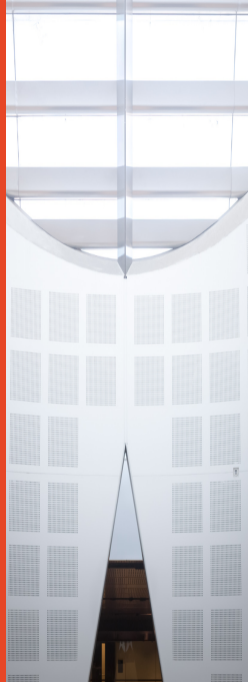


Incremental Datalog Prototype in Soufflé

David Zhao¹ Bernhard Scholz¹ Pavle Subotić²

¹University of Sydney

²Amazon



Introduction

Introduction

Datalog is a logic programming language, used for

- ▶ Program analysis - security property checking in OpenJDK
- ▶ Network security analysis - enforcing security rules in AWS networks
- ▶ Binary disassembly - Grammatech
- ▶ Smart contract analysis - checking for exploits on EVM

Introduction

Datalog is a logic programming language, used for

- ▶ Program analysis - security property checking in OpenJDK
- ▶ Network security analysis - enforcing security rules in AWS networks
- ▶ Binary disassembly - Grammatech
- ▶ Smart contract analysis - checking for exploits on EVM

Datalog is *succinct* and *easy to use*

Introduction

Datalog is a logic programming language, used for

- ▶ Program analysis - security property checking in OpenJDK
- ▶ Network security analysis - enforcing security rules in AWS networks
- ▶ Binary disassembly - Grammatech
- ▶ Smart contract analysis - checking for exploits on EVM

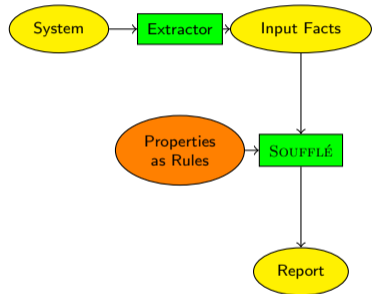
Datalog is *succinct* and *easy to use*

Parallel Datalog engines (e.g. SOUFFLÉ) allow high performance

- ▶ Comparable to optimized, hand-crafted tools
- ▶ Maintaining ease of use

Soufflé

- ▶ System as a set of facts
- ▶ System properties as a set of Datalog rules
- ▶ Datalog engine (SOUFFLÉ) evaluates rules using facts to check if properties are satisfied



Example - Pointer Analysis

System

```
1 a = new O();  
2 b = a;  
3  
4 a.f = c;  
5 d = a.f;
```

Example - Pointer Analysis

System

```
1 a = new 0();  
2 b = a;  
3  
4 a.f = c;  
5 d = a.f;
```

Input Facts

```
1 alloc(a,l1).  
2 assign(b,a).  
3  
4 store(a,f,c).  
5 load(d,a,f).
```

Example - Pointer Analysis

System

```
1 a = new 0();  
2 b = a;  
3  
4 a.f = c;  
5 d = a.f;
```

Input Facts

```
1 alloc(a, l1).  
2 assign(b, a).  
3  
4 store(a, f, c).  
5 load(d, a, f).
```

Output:

- ▶ vpt(a, l1).
- ▶ vpt(b, l1).

- ▶ alias(c, d).
- ▶ ...

Example

Pointer Analysis

```
1 // allocation sites (x = new 0())
2 vpt(Var, Obj) :- alloc(Var, Obj).
3
4 // assignments (x = y)
5 vpt(Var, Obj) :- assign(Var, Var2), vpt(Var2, Obj).
6
7 // load-store pairs (v = y.f; p.f = q; y and p alias)
8 vpt(Var, Obj) :- load(Var, Y, F), store(P, F, Q),
9                   vpt(Q, Obj), vpt(P, AliasObj), vpt(Y, AliasObj).
10
11 // generating the alias relation
12 alias(Var1, Var2) :- vpt(Var1, Obj), vpt(Var2, Obj), Var1 != Var2.
```

Incremental Evaluation

Motivation

Security analysis can be slow - up to 2 weeks for OpenJDK

Motivation

Security analysis can be slow - up to 2 weeks for OpenJDK

During software development:

- ▶ Make small change to codebase
- ▶ Run program analysis

Motivation

Security analysis can be slow - up to 2 weeks for OpenJDK

During software development:

- ▶ Make small change to codebase
- ▶ Run program analysis

Small change → fast re-run of analysis

Incremental Evaluation

Small change to system:

$$(input_1 \cup input^+) \setminus input^- = input_2$$

Incremental Evaluation

Small change to system:

$$(input_1 \cup input^+) \setminus input^- = input_2$$

Without incremental:

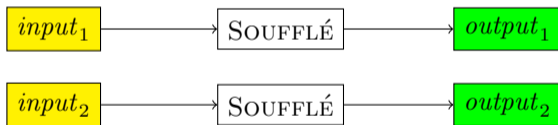


Incremental Evaluation

Small change to system:

$$(input_1 \cup input^+) \setminus input^- = input_2$$

Without incremental:

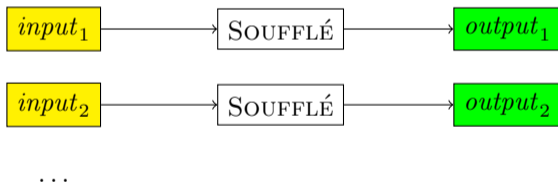


Incremental Evaluation

Small change to system:

$$(input_1 \cup input^+) \setminus input^- = input_2$$

Without incremental:



Incremental Evaluation

Small change to system:

$$(input_1 \cup input^+) \setminus input^- = input_2$$

With incremental:

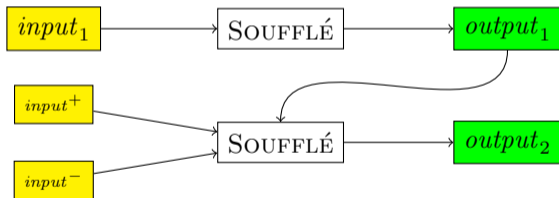


Incremental Evaluation

Small change to system:

$$(input_1 \cup input^+) \setminus input^- = input_2$$

With incremental:

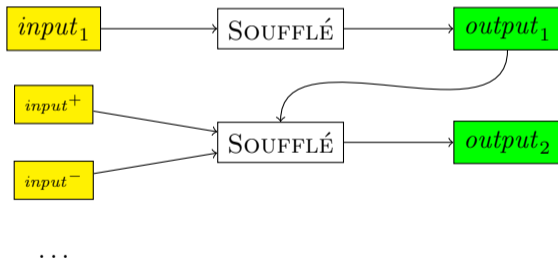


Incremental Evaluation

Small change to system:

$$(input_1 \cup input^+) \setminus input^- = input_2$$

With incremental:



Counting Method

For non-recursive, initially proposed by [Gupta et al., 1993]

Counting Method

For non-recursive, initially proposed by [Gupta et al., 1993]

- ▶ Maintain a *count* for number of ways to derive each tuple

Counting Method

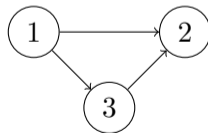
For non-recursive, initially proposed by [Gupta et al., 1993]

- ▶ Maintain a *count* for number of ways to derive each tuple

Simpler Example

-
- 1 `hop2(X, Y) :- edge(X, Y).`
 - 2 `hop2(X, Z) :- edge(X, Y), edge(Y, Z).`
-

`edge(1, 2), edge(1, 3), edge(3, 2)`



Counting Method

All Derivations

`hop2(1, 2) :- edge(1, 2)`

`hop2(1, 3) :- edge(1, 3)`

`hop2(3, 2) :- edge(3, 2)`

`hop2(1, 2) :- edge(1, 3), edge(3, 2)`

Counting Method

All Derivations

`hop2(1, 2) :- edge(1, 2)`

`hop2(1, 3) :- edge(1, 3)`

`hop2(3, 2) :- edge(3, 2)`

`hop2(1, 2) :- edge(1, 3), edge(3, 2)`

Relations are multisets instead of sets

$$\text{hop2} = \{(1, 2) \times 2, (1, 3) \times 1, (3, 2) \times 1\}$$

Counting Method

Deleting a tuple

$$\Delta_{\text{edge}} = \{(1, 2) \times -1\}$$

Counting Method

Deleting a tuple

$$\Delta\text{edge} = \{(1, 2) \times -1\}$$

-
- 1 $\Delta\text{hop2}(X, Y) :- \Delta\text{edge}(X, Y).$
 - 2 $\Delta\text{hop2}(X, Z) :- \Delta\text{edge}(X, Y), \text{edge}(Y, Z).$
 - 3 $\Delta\text{hop2}(X, Z) :- \text{edge}'(X, Y), \Delta\text{edge}(Y, Z). // \text{edge}' = \text{edge} + \Delta\text{edge}$
-

Counting Method

Deleting a tuple

$$\Delta\text{edge} = \{(1, 2) \times -1\}$$

-
- 1 $\Delta\text{hop2}(X, Y) :- \Delta\text{edge}(X, Y).$
 - 2 $\Delta\text{hop2}(X, Z) :- \Delta\text{edge}(X, Y), \text{edge}(Y, Z).$
 - 3 $\Delta\text{hop2}(X, Z) :- \text{edge}'(X, Y), \Delta\text{edge}(Y, Z). // \text{edge}' = \text{edge} + \Delta\text{edge}$
-

From rule 1: $\Delta\text{hop2}(1, 2)^{-1} :- \Delta\text{edge}(1, 2)^{-1}$

$$\begin{aligned}\Delta\text{hop2} &= \{(1, 2) \times -1\} \\ \implies \text{hop2} &= \{(1, 2) \times 1, (1, 3) \times 1, (3, 2) \times 1\}\end{aligned}$$

Problems with Recursion

Transitive Closure

```
1 path(X, Y) :- edge(X, Y).  
2 path(X, Z) :- edge(X, Y), path(Y, Z).
```

Problems with Recursion

Transitive Closure

-
- 1 `path(X, Y) :- edge(X, Y).`
 - 2 `path(X, Z) :- edge(X, Y), path(Y, Z).`
-

A cycle:

`edge(a,a), edge(a,b), edge(b,a)`



Problems with Recursion

Transitive Closure

-
- 1 `path(X, Y) :- edge(X, Y).`
 - 2 `path(X, Z) :- edge(X, Y), path(Y, Z).`
-

A cycle:

`edge(a,a)`, `edge(a,b)`, `edge(b,a)`



How many times is `path(a,a)` derived?

How many times to decrement count when deleting?

Extending to Recursive Datalog

An additional dimension to the counting

Maintain counts *per iteration* of Datalog evaluation

Inspired by Differential Dataflow [McSherry et al., 2013]

Extending to Recursive Datalog

An additional dimension to the counting

Maintain counts *per iteration* of Datalog evaluation

Inspired by Differential Dataflow [McSherry et al., 2013]

$$\text{Iteration}_0 : \{ \text{edge}(a, a)^{0,1}, \text{edge}(a, b)^{0,1}, \text{edge}(b, a)^{0,1} \}$$

Extending to Recursive Datalog

An additional dimension to the counting

Maintain counts *per iteration* of Datalog evaluation

Inspired by Differential Dataflow [McSherry et al., 2013]

Iteration₀ : {edge(*a*, *a*)^{0,1}, edge(*a*, *b*)^{0,1}, edge(*b*, *a*)^{0,1}}

Iteration₁ : {path(*a*, *a*)^{1,1}, path(*a*, *b*)^{1,1}, path(*b*, *a*)^{1,1}}

Extending to Recursive Datalog

An additional dimension to the counting

Maintain counts *per iteration* of Datalog evaluation

Inspired by Differential Dataflow [McSherry et al., 2013]

Iteration₀ : {edge(*a*, *a*)^{0,1}, edge(*a*, *b*)^{0,1}, edge(*b*, *a*)^{0,1}}

Iteration₁ : {path(*a*, *a*)^{1,1}, path(*a*, *b*)^{1,1}, path(*b*, *a*)^{1,1}}

Iteration₂ : {path(*a*, *a*)^{2,1}, path(*b*, *b*)^{2,1}}

Extending to Recursive Datalog

An additional dimension to the counting

Maintain counts *per iteration* of Datalog evaluation

Inspired by Differential Dataflow [McSherry et al., 2013]

Iteration₀ : {edge(*a*, *a*)^{0,1}, edge(*a*, *b*)^{0,1}, edge(*b*, *a*)^{0,1}}

Iteration₁ : {path(*a*, *a*)^{1,1}, path(*a*, *b*)^{1,1}, path(*b*, *a*)^{1,1}}

Iteration₂ : {path(*a*, *a*)^{2,1}, path(*b*, *b*)^{2,1}}

path(*a*, *a*) derived:

- ▶ Once in iteration 1
- ▶ Once in iteration 2

Unrolling the Recursion

Storing per-iteration counts is like unrolling the recursion

```
1 path(X, Y)_1 :- edge(X, Y).
2 path(X, Y)_2 :- edge(X, Y1), edge(Y1, Y).
3 path(X, Y)_3 :- edge(X, Y1), edge(Y1, Y2), edge(Y2, Y).
4 path(X, Y)_4 :- edge(X, Y1), edge(Y1, Y2), edge(Y2, Y3), edge(Y3, Y).
5 ...
```

Our Implementation

Our Ideas

Motivation

Integrate counting inside semi-naïve

- ▶ Minimal changes to evaluation strategy

Our Ideas

Motivation

Integrate counting inside semi-naïve

- ▶ Minimal changes to evaluation strategy

Maintaining the Δ version of the relation

1. Maintaining Δ explicitly - requires extra indexes and rules (new evaluation strategy)
2. Simulating Δ with extra counts - requires more storage per tuple

Simulating Δ Relations

Introduce a third dimension for each tuple - *epoch*

Simulating Δ Relations

Introduce a third dimension for each tuple - *epoch*

1. Initial evaluation is epoch 0
2. Each subsequent insertion/deletion and Datalog re-evaluation is a new epoch

Epoch value for a tuple \iff latest epoch in which it was updated

E.g. $\text{path}(a, a)^3 \iff \text{path}(a, a)$ was last updated in epoch 3

Simulating Δ Relations

Introduce a third dimension for each tuple - *epoch*

1. Initial evaluation is epoch 0
2. Each subsequent insertion/deletion and Datalog re-evaluation is a new epoch

Epoch value for a tuple \iff latest epoch in which it was updated

E.g. $\text{path}(a, a)^3 \iff \text{path}(a, a)$ was last updated in epoch 3

Simulating Δ

Epoch value = *current* epoch \iff tuple is in Δ

Challenges

How do we know if an updated tuple is *actually* new?

Challenges

How do we know if an updated tuple is *actually* new?

Reverse Path

```
1 path(X, Y) :- edge(X, Y).
2 path(X, Z) :- edge(X, Y), path(Y, Z).
3 reversepath(X, Y) :- path(Y, X).
```

Challenges

How do we know if an updated tuple is *actually* new?

Reverse Path

- 1 `path(X, Y) :- edge(X, Y).`
 - 2 `path(X, Z) :- edge(X, Y), path(Y, Z).`
 - 3 `reversepath(X, Y) :- path(Y, X).`
-

$$\text{edge} = \{(a, b)^{0,1}, (b, c)^{0,1}\}$$

$$\text{path} = \{(a, b)^{1,1}, (b, c)^{1,1}, (a, c)^{2,1}\}$$



Challenges

How do we know if an updated tuple is *actually* new?

Reverse Path

-
- 1 `path(X, Y) :- edge(X, Y).`
 - 2 `path(X, Z) :- edge(X, Y), path(Y, Z).`
 - 3 `reversepath(X, Y) :- path(Y, X).`
-

$$\text{edge} = \{(a, b)^{0,1}, (b, c)^{0,1}\}$$

$$\text{path} = \{(a, b)^{1,1}, (b, c)^{1,1}, (a, c)^{2,1}\}$$

$$\text{reversepath} = \{(b, a)^{2,1}, (c, b)^{2,1}, (c, a)^{3,1}\}$$



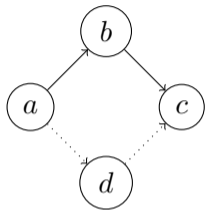
Challenges

$$\text{edge} = \{(a, b)^{0,1}, (b, c)^{0,1}\}$$

$$\text{path} = \{(a, b)^{1,1}, (b, c)^{1,1}, (a, c)^{2,1}\}$$

$$\text{reversepath} = \{(b, a)^{2,1}, (c, b)^{2,1}, (c, a)^{3,1}\}$$

Consider: $\Delta\text{edge} = \{(a, d)^{0,1}, (d, c)^{0,1}\}$

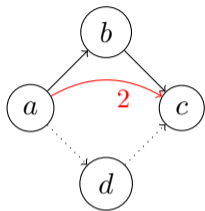


Challenges

$$\text{edge} = \{(a, b)^{0,1}, (b, c)^{0,1}\}$$

$$\text{path} = \{(a, b)^{1,1}, (b, c)^{1,1}, (a, c)^{2,1}\}$$

$$\text{reversepath} = \{(b, a)^{2,1}, (c, b)^{2,1}, (c, a)^{3,1}\}$$



Consider: $\Delta\text{edge} = \{(a, d)^{0,1}, (d, c)^{0,1}\}$

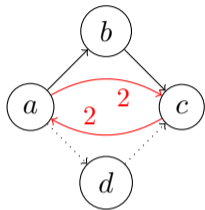
$$\Delta\text{path} = \{(a, d)^{1,1}, (d, c)^{1,1}, (a, c)^{2,1}\} \implies \text{path}(a, c)^{2,2}$$

Challenges

$$\text{edge} = \{(a, b)^{0,1}, (b, c)^{0,1}\}$$

$$\text{path} = \{(a, b)^{1,1}, (b, c)^{1,1}, (a, c)^{2,1}\}$$

$$\text{reversepath} = \{(b, a)^{2,1}, (c, b)^{2,1}, (c, a)^{3,1}\}$$



Consider: $\Delta\text{edge} = \{(a, d)^{0,1}, (d, c)^{0,1}\}$

$$\Delta\text{path} = \{(a, d)^{1,1}, (d, c)^{1,1}, (a, c)^{2,1}\} \implies \text{path}(a, c)^{2,2}$$

$$\Delta\text{reversepath} = \{(d, a)^{2,1}, (c, d)^{2,1}, (c, a)^{3,1}\} \implies \text{reversepath}(c, a)^{3,2}$$

Challenges

$$\Delta \text{reversepath}(c, a)^{3,1} \implies \text{reversepath}(c, a)^{3,2}$$

But actually

The only derivation is:

$$\text{reversepath}(c, a) \text{ :- path}(a, c).$$

Prototype in Soufflé

Implemented in the Soufflé
Datalog compiler

- ▶ Per-iteration count
- ▶ Simulating Δ relations with epoch

```
y = x; a.f = x; y = a.f;
```

```
-----  
vpt  
x      y      @iteration      @epoch  @count  
=====
```

x	11	1	0	1
y	11	2	0	2
a	12	1	0	1

```
=====
```

x	y	@iteration	@epoch	@count
x	y	3	0	1
y	x	3	0	1

```
=====
```

Incremental is invoked.
Enter command >

Prototype in Soufflé

Implemented in the Soufflé
Datalog compiler

- ▶ Per-iteration count
- ▶ Simulating Δ relations with epoch

```
y = x; a.f = x; y = a.f;
```

```
Enter command > remove assign("y","x")
vpt
x      y      @iteration      @epoch  @count
=====
x      11      1      0      1
y      11      2      -1     1
a      12      1      0      1
=====
alias
x      y      @iteration      @epoch  @count
=====
x      y      3      0      1
y      x      3      0      1
=====
Incremental is invoked.
Enter command >
```

Prototype in Soufflé

Implemented in the Soufflé
Datalog compiler

- ▶ Per-iteration count
- ▶ Simulating Δ relations with epoch

```
y = x; a.f = x; y = a.f;
```

```
Enter command > remove load("y","a","f")
```

```
vpt
```

x	y	@iteration	@epoch	@count
=====				
x	11	1	0	1
y	11	2	-2	0
a	12	1	0	1

```
=====
```

```
alias
```

x	y	@iteration	@epoch	@count
=====				
x	y	3	-2	0
y	x	3	-2	0

```
=====
```

```
Incremental is invoked.
```

```
Enter command >
```

Prototype in Soufflé

Implemented in the Soufflé
Datalog compiler

- ▶ Per-iteration count
- ▶ Simulating Δ relations with epoch

```
y = x; a.f = x; y = a.f;
```

```
Enter command > insert assign("y","x")
```

```
vpt
```

x	y	@iteration	@epoch	@count
=====				
x	11	1	0	1
y	11	2	3	1
a	12	1	0	1

```
=====
```

```
alias
```

x	y	@iteration	@epoch	@count
=====				
x	y	3	3	1
y	x	3	3	1

```
=====
```

```
Incremental is invoked.
```

```
Enter command >
```


Preliminary Results

Super early prototype:

- ▶ The reversepath issue
- ▶ B-Tree deletion not yet implemented

Preliminary Results

Super early prototype:

- ▶ The reversepath issue
- ▶ B-Tree deletion not yet implemented

Benchmark: initial evaluation of transitive closure, path graph of 1000 nodes

Soufflé (no inc.)	Soufflé (inc.)	DDlog (Differential Dataflow)
0.3 sec	0.5 sec	18 sec

Preliminary Results

Super early prototype:

- ▶ The reversepath issue
- ▶ B-Tree deletion not yet implemented

Benchmark: initial evaluation of transitive closure, path graph of 1000 nodes

Soufflé (no inc.)	Soufflé (inc.)	DDlog (Differential Dataflow)
0.3 sec	0.5 sec	18 sec

Incremental update: removing edge(499, 500)

Soufflé (no inc.)	Soufflé (inc.)	DDlog (Differential Dataflow)
-	0.2 sec	14 sec

Outlook

- ▶ We propose some new ideas for incremental evaluation of recursive Datalog by counting, with promising preliminary results
- ▶ We need to resolve the `reversepath` issue
- ▶ We need to extend to negations and aggregates