

# Context-Sensitive Data-Dependence Analysis via Linear Conjunctive Language Reachability

Qirun Zhang, Zhendong Su  
POPL 2017

Presented by Ayaka Yorihiro  
PLDG 11/18/2020

# Static Program Analyses

## Context Sensitivity:

- Distinguish between multiple calls of the same function

```
1 def ident(i):  
2     return i  
3 ...  
4 a = 5  
5 b = "hi"  
6 c = ident(a)  
7 d = ident(b)
```

c and d have  
different  
values

## Data Sensitivity:

- Being able to follow pointers /information through indirection

```
1 def changeX(obj):  
2     obj.x = 5  
3 a = Thing()  
4 a.x = "hello"  
5 b = a  
6 changeX(b)  
7 print(a.x)
```

a.x is now 5

# Taint Analysis Example

```
1 class A {  
2     ...  
3     Text getT() {  
4         return this.t;  
5     }  
6 }  
7 A a; ...  
8 a.t = pw;  
9 msg = a.getT();
```

Want to prove that super important value `pw` gets passed to `msg`...

We want to be **sound**  
→ No false negatives

# Taint Analysis Example

```
1 class A {  
2   ...  
3   Text getT() {  
4     return this.t;  
5   }  
6 }  
7 A a; ...  
8 a.t = pw;  
9 msg = a.getT();
```

Return value gets t

method return  
to line 9

t set to pw

method call  
@line 9

Nodes: Storage  
Location  
Edges: Operations on  
Program  
(: call        ): return  
[: set        ]: get

Return value gets t

$\Pi_t$

method call  
@line 9

$\Pi_9$

pw

$\Pi_t$

t set to pw

this

ret

$\Pi_9$

msg

method  
return  
to line 9

# Taint Analysis Example

If there is a **valid path** from **pw** to **msg**, then it means that the value of **pw** would be propagated to **msg**

IS  $\llbracket_t \downarrow_9 \rrbracket_t \downarrow_9$  VALID?

# Dyck Language (Single Parenthesis)

- Matched Parenthesis
- Dyck language  $D_k$  of size  $k$ ...

$$S \rightarrow SS \mid \varepsilon \mid ((_0 S)_{0}) \mid \dots \mid ((_k S)_{k})$$

- ex)

○  $((_a)_{a}((_a)_{a})_{a})$  ✓

○  $((_a((b)_{b}))_{a})$  ✓

○  $((_a((b)_{a}))_{b})$  ✗

# Interleaved Dyck Languages (Paren + Bracket)

$\mathbb{D}_C$  : Context Sensitivity

ex)  $\langle \rangle_9 \langle \rangle_3 \langle \rangle_3 \rangle_9$

$\mathbb{D}_D$  : Data/Field Sensitivity

ex)  $\llbracket \rrbracket_a \rrbracket_a$

$\mathbb{D}_{CD}$  : Interleaved Dyck Languages – both Context Sensitive and Data/Field Sensitive

ex)  $\langle \rangle_9 \langle \rangle_3 \llbracket \rrbracket_a \rangle_3 \rrbracket_a \rangle_9$

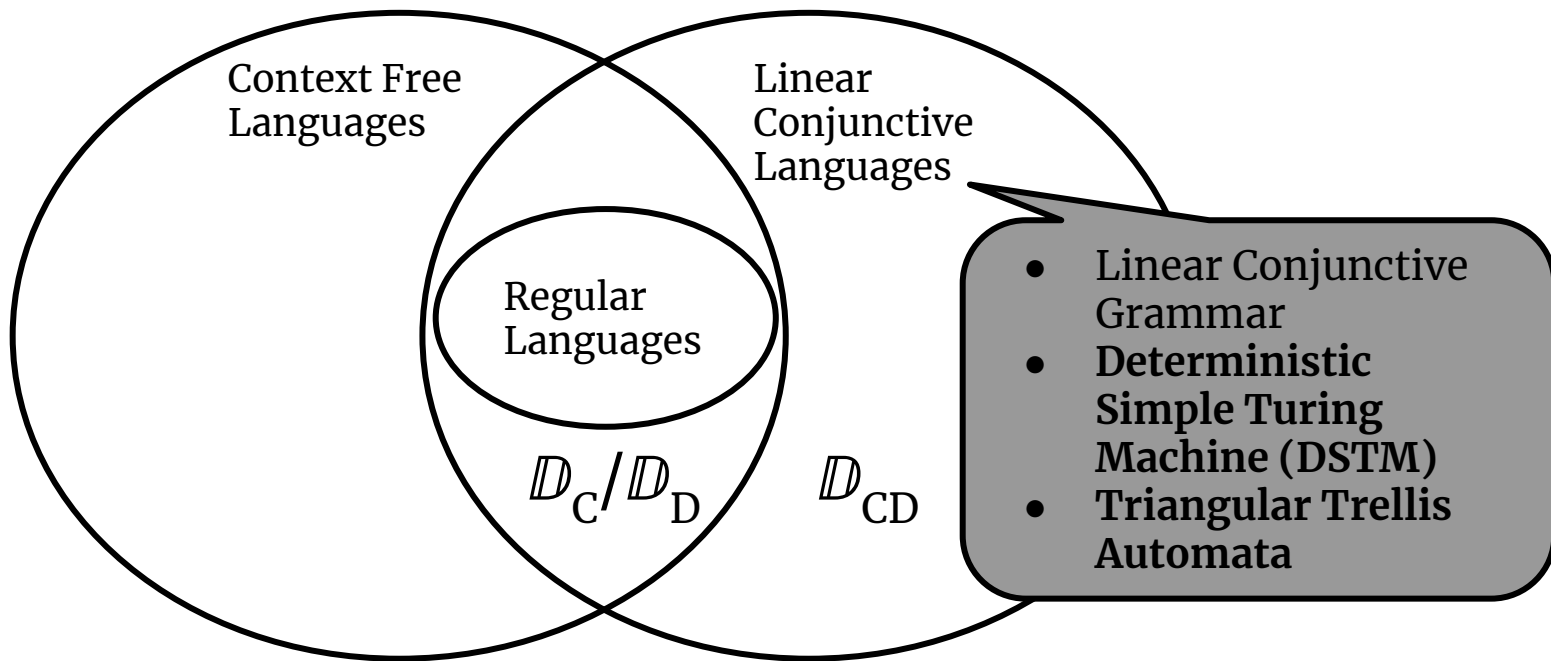
**Reasoning perfectly about context and data sensitivity simultaneously is UNDECIDABLE!**

# Problem Space and Contributions

- Static Program Analyses aim for high Context- and Data-Sensitivity in order to reason about programs
- Context-Sensitive Data-Sensitive Analysis can be expressed as a graph reachability problem over the Interleaved Dyck Language  $\mathbb{D}_{CD}$
- This problem is undecidable
- $\mathbb{D}_{CD}$  is a subset of the Linear Conjunctive Languages
- Authors devise approximation algorithms for Linear Conjunctive Language Reachability



# Language Classes wrt $\mathbb{D}_{CD}$ ...



# Deterministic Simple Turing Machine (DSTM)

6-tuple  $M = (Q, \Sigma, \Gamma, \delta, q_0, F)$

$Q$ : set of states

$q_0$ : start state

$\Sigma$ : input alphabet

$F \subseteq Q$ : accepting states

$\Gamma$ : work tape alphabet

$\delta$ : transition function

\$ - end of tape

$Q \times \Gamma \times (\Sigma \cup \{\varepsilon\}) \rightarrow$

$\sqcup$  - blank symbol

$Q \times (\Gamma \setminus \{\sqcup\}) \times \{L, R\}$

# DSTM Transition Function Rules

For all  $\delta(q, Z, \sigma) = (q', Z', d)$  where  $q, q' \in Q$ ,  $Z, Z' \in \Gamma$ ,  $\sigma \in \Sigma$ ,  $d \in \{L, R\}$ ...

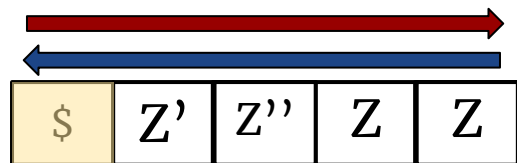
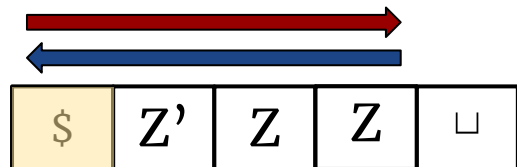
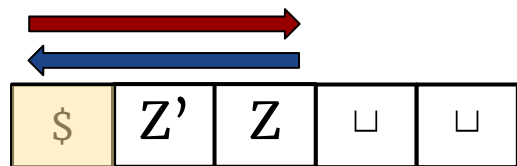
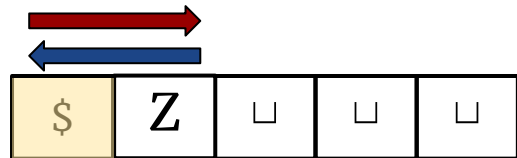
1. If  $Z = \$$ , then  $Z' = \$$
2. If  $q = \text{"L} \rightarrow \text{R sweep"}$  and  $Z \neq \sqcup$ , then  $\sigma = \varepsilon$ ,  $q' = \text{"L} \rightarrow \text{R sweep"}$ ,  $Z' = Z$  and  $d = R$
3. If  $q = \text{"L} \rightarrow \text{R sweep"}$  and  $Z = \sqcup$ , then  $\sigma \neq \varepsilon$ ,  $q' \neq \text{"L} \rightarrow \text{R sweep"}$ , and  $d = L$

(L  $\rightarrow$  R sweep... keep going right, process input at the rightmost end of tape)

4. If  $q \neq \text{"L} \rightarrow \text{R sweep"}$  and  $Z \neq \$$ , then  $\sigma = \varepsilon$ ,  $q' \neq \text{"L} \rightarrow \text{R sweep"}$ , and  $d = L$
5. If  $q \neq \text{"L} \rightarrow \text{R sweep"}$  and  $Z = \$$ , then  $\sigma = \varepsilon$ ,  $q' = \text{"L} \rightarrow \text{R sweep"}$ , and  $d = R$

(L  $\leftarrow$  R sweep... we can write to tape at any point, but we can't read input and we have to keep going left until we get to the leftmost end)

# DSTM Transition Function Rules



## Left to Right Sweep:

- Stay in designated state for L→R sweep
- Keep going Right
- Can't read input
- Can't change tape content

## Read a \$:

- End of Right to Left
- If no more input, computation ends here
- Start of Left to Right otherwise

## Read a □:

- End of Left to Right
- Can read input
- Can write to tape
- Start Right to Left

## Right to Left Sweep:

- Can't go to L→R Sweep state, can go to any other state
- Keep going Left
- Can't read input
- Can change tape content

# Example... $\langle 1 \rangle \langle 2 \rangle \langle 2 \rangle \langle 1 \rangle$

q: L → R 😞

\$	$\langle 1 \rangle$	□	□	□
----	---------------------	---	---	---

q: L → R 😞 😞

\$	$\langle 1 \rangle$	$\langle 2 \rangle$	□	□
----	---------------------	---------------------	---	---

q: L → R 😞 😊  $\langle 2 \rangle$

\$	$\langle 1 \rangle$	#	#	□
----	---------------------	---	---	---

q: DONE 😊  $\langle 1 \rangle$   $\langle 1 \rangle$   $\langle 1 \rangle$

\$	#	#	#	#
----	---	---	---	---

Input  $\langle 1 \rangle$ :

- Want to remember which  $\langle 1 \rangle$  we read  
→ Write  $\langle 1 \rangle$  to tape
- We need to resolve  $\langle 1 \rangle$  in the future  
→ State 😞

Input  $\langle 2 \rangle$ :

- Don't need to record this input  
→ Write # to tape
- “Looking for” tape symbol  $\langle 1 \rangle$   
→ State  $\langle 1 \rangle$

State 😞:

There is a  $\langle 1 \rangle$

Always write  
same symbol

State  $\langle 1 \rangle$ :

Ignores #

Tape symbol  $\langle 1 \rangle$  →  
write # & go to 😊

State 😊:

Ignores #

Tape symbol  $\langle 1 \rangle$  →  
write  $\langle 1 \rangle$  & go to 😞

# DSTM for $\mathbb{D}_C$

Goals:

Want to match  $q_i, q_i$  together

$$Q = \{\text{😊}, \text{😞}\} \cup \{q_0, \dots\}$$

$$\Sigma = \{q_0, q_0, \dots\}$$

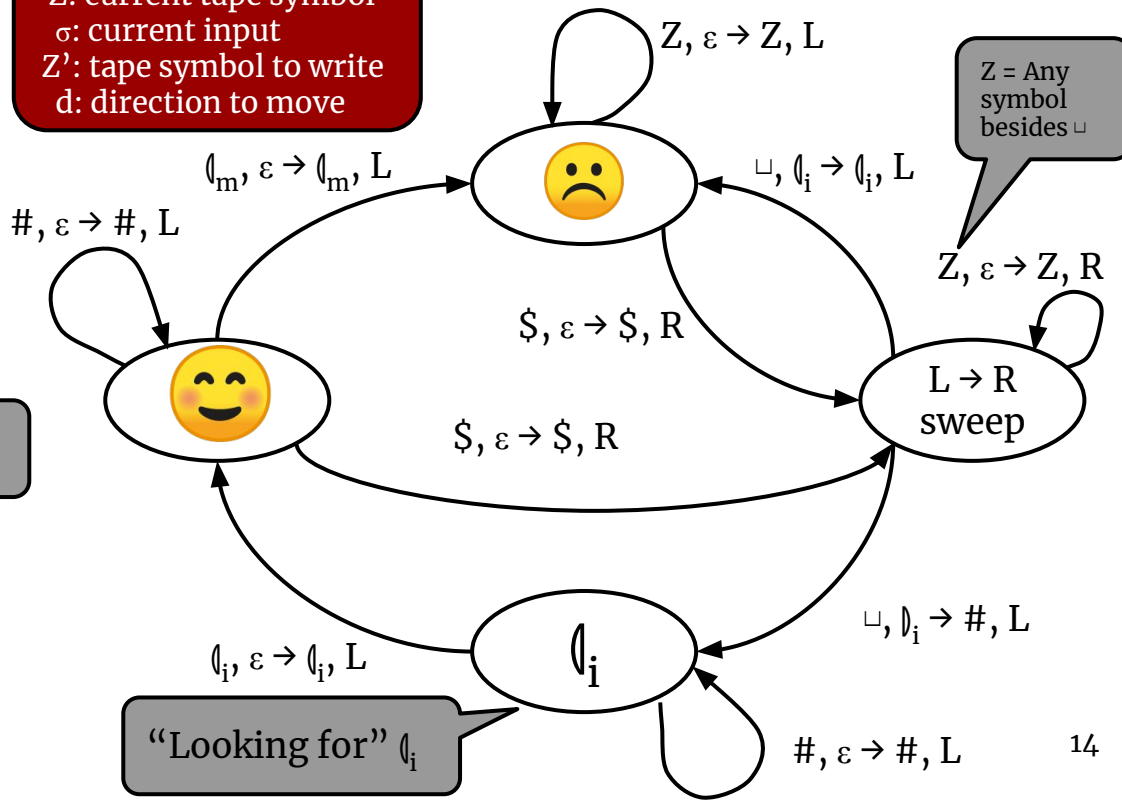
“Looking for”  $q_0$

$$\Gamma = \{\$, \sqcup, \#\} \cup \{q_0, \dots\}$$

$$F = \{\text{😊}\}$$

“Don't care”

**Notation:**  $Z, \sigma \rightarrow Z', d$   
 $Z$ : current tape symbol  
 $\sigma$ : current input  
 $Z'$ : tape symbol to write  
 $d$ : direction to move



# Example... $(q_1 \sqcup_a)_1 \sqcup_a$

q: L → R 😞

\$	$q_1$	$\sqcup$	$\sqcup$	$\sqcup$
----	-------	----------	----------	----------

z:

q: L → R 😞 😞

\$	$q_1$	$\sqcup_a$	$\sqcup$	$\sqcup$
----	-------	------------	----------	----------

z:

q: L → R 😞 😞  $q_1$   $q_1$

\$	#	$\sqcup_a$	#	$\sqcup$
----	---	------------	---	----------

z:

q: DONE 😊 😊  $\sqcup_a$   $\sqcup_a$

\$	#	#	#	#
----	---	---	---	---

z:

**Input  $q_i/\sqcup_i$ :**  
 → Write  $q_i/\sqcup_i$  to tape  
 → State 😞

**State 😞:**  
 There is a  $q/\sqcup$   
 Always write same symbol

**State  $q_i$ :**  
 Ignores #  
 Tape symbol  $q_i$   
 → write # & go to 😊  
**NEW:** Tape Symbol  $\sqcup_m$   
 → write  $\sqcup_m$  & go to state 😞  $q_i$

**State  $\sqcup_i$ :**  
 Ignores #  
 Tape symbol  $\sqcup_i$   
 → write # & go to 😊  
**NEW:** Tape Symbol  $q_m$   
 → write  $q_m$  & go to 😞  $\sqcup_i$

**Input  $q_i/\sqcup_i$ :**  
 → Write # to tape  
 → State  $q_i/\sqcup_i$

**State 😊:**  
 Ignores #  
 Symbol  $q_m/\sqcup_m$  → write  $q_m/\sqcup_m$  & go to 😞

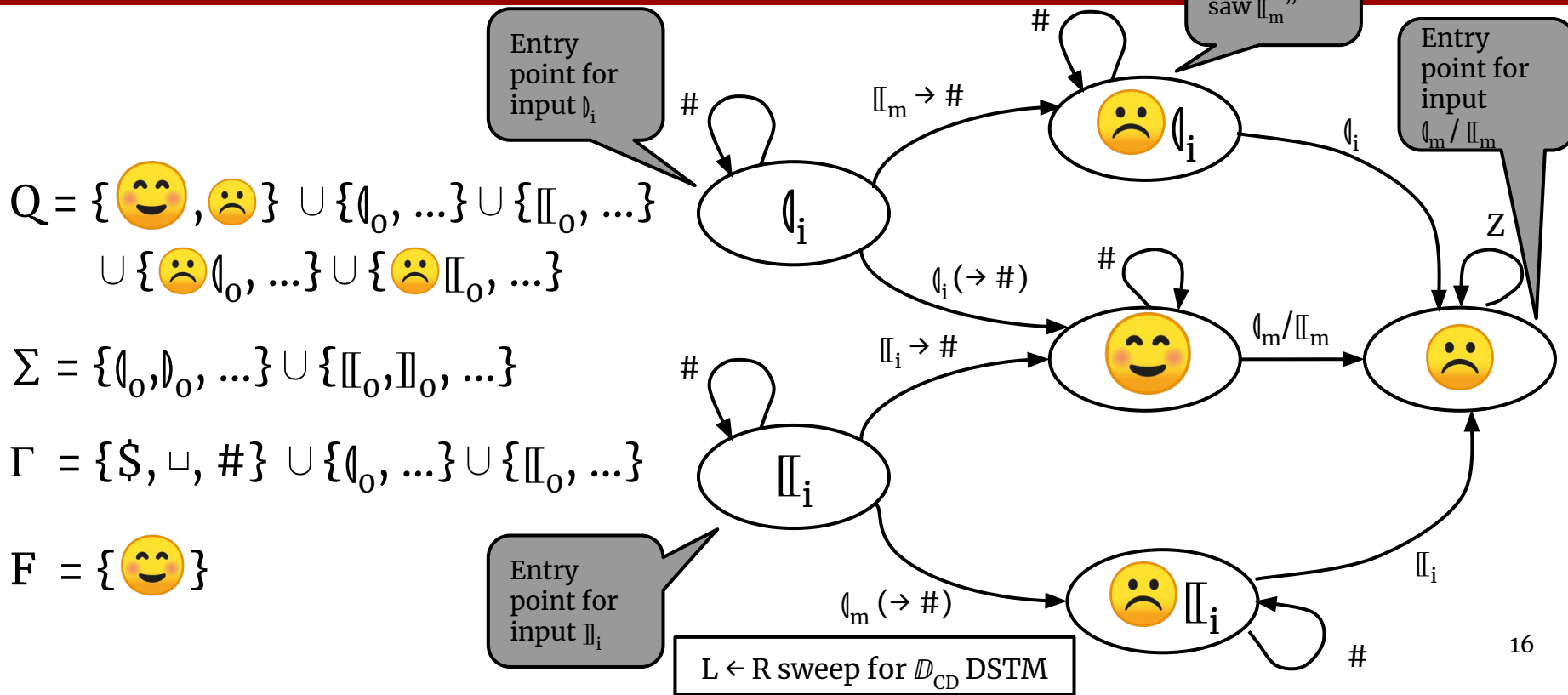
**NEW! State 😞  $q_i$ :**  
 “Looking for  $q_i$  but seen  $\sqcup_m$ ”  
 Ignores #,  $\sqcup_m$

Tape symbol  $q_i$   
 → write # & go to 😞

**NEW! State 😞  $\sqcup_i$ :**  
 “Looking for  $\sqcup_i$  but seen  $q_m$ ”  
 Ignores #,  $q_m$

Tape symbol  $\sqcup_i$   
 → write # & go to 😞

# Extend to build DSTM for $\mathbb{D}_{CD}$





# DSTM for $\mathbb{D}_{CD}$

**Notation:**  $Z, \sigma \rightarrow Z', d$   
 $Z$ : current tape symbol  
 $\sigma$ : current input  
 $Z'$ : tape symbol to write  
 $d$ : direction to move

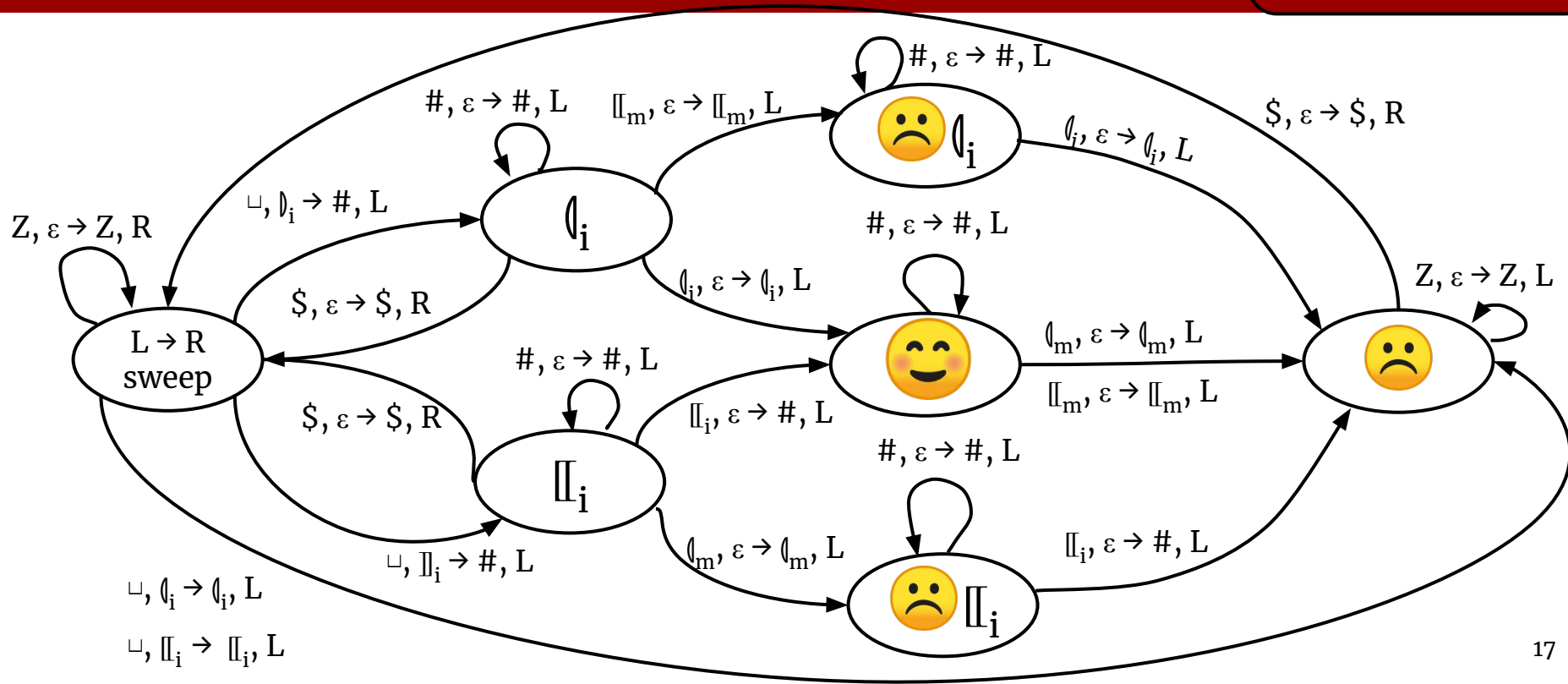


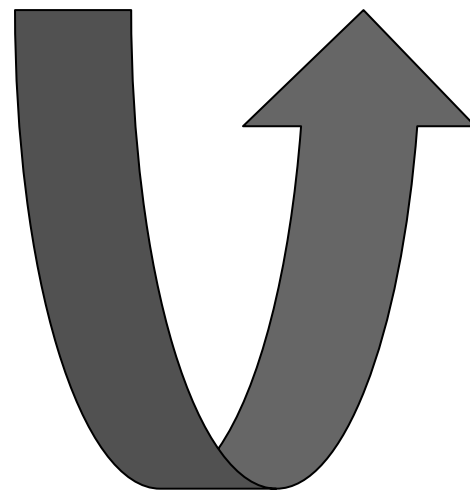
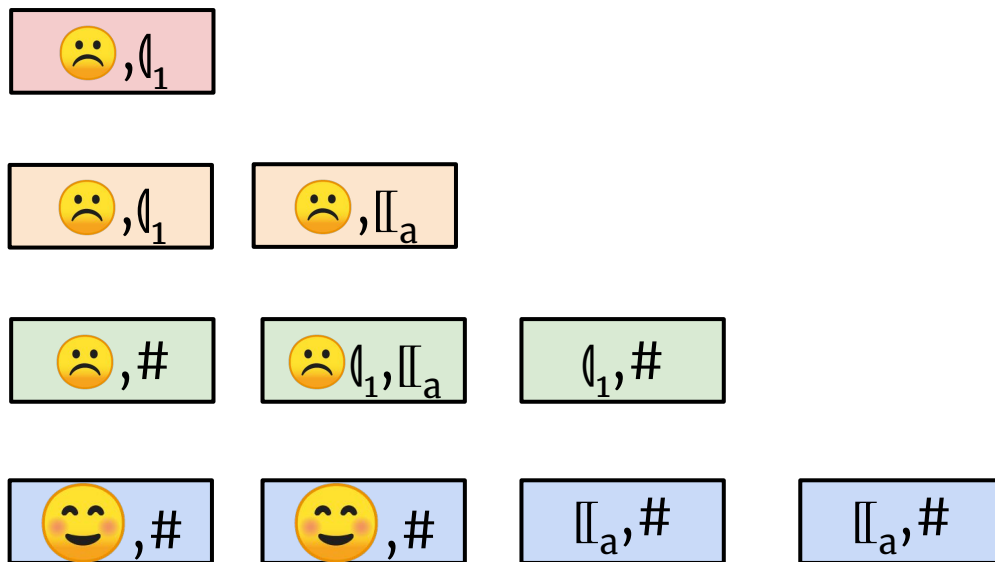
Diagram illustrating the execution of a Turing Machine (TM) on the input string "a". The diagram shows five rows of the TM's state and tape configuration.

- Row 1: State is  $L \rightarrow R$  (sad face), tape is  $\$ \square \square \square$ .
- Row 2: State is  $L \rightarrow R$  (sad face), tape is  $\$ \square_1 \square \square$ .
- Row 3: State is  $L \rightarrow R$  (sad face), tape is  $\$ \square_1 \square_a \square$ . The head is at the first  $\square_a$ .
- Row 4: State is  $L \rightarrow R$  (sad face), tape is  $\$ \# \square_1 \square$ . The head is at the first  $\square_1$ .
- Row 5: State is **DONE** (happy face), tape is  $\$ \# \# \# \#$ . The head is at the first  $\#$ .

Arrows indicate the head's movement: from the first  $\square_a$  in Row 3 to the first  $\square_1$  in Row 4, and from the first  $\square_1$  in Row 4 to the first  $\#$  in Row 5.

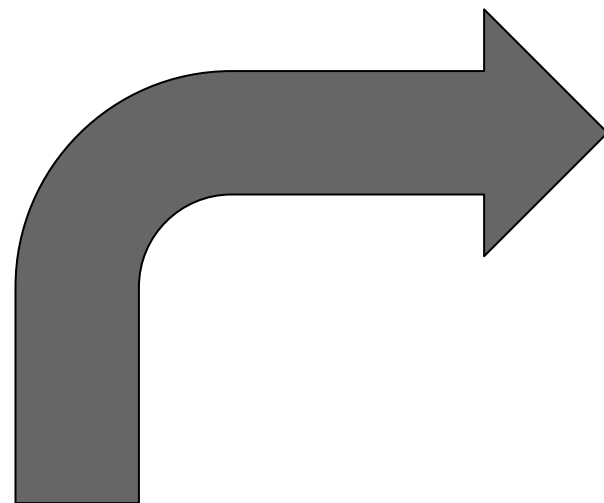
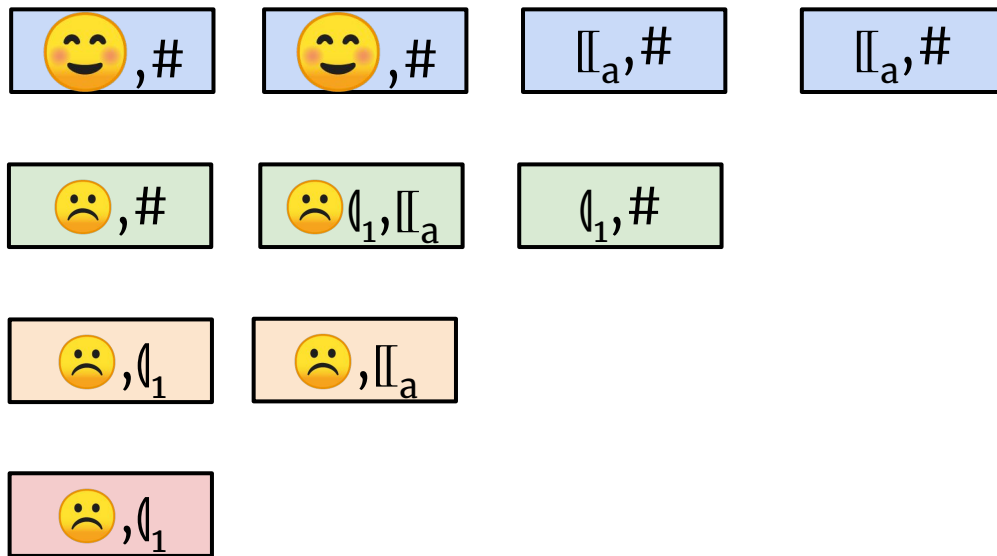


# Let's put together the DSTM state and tape



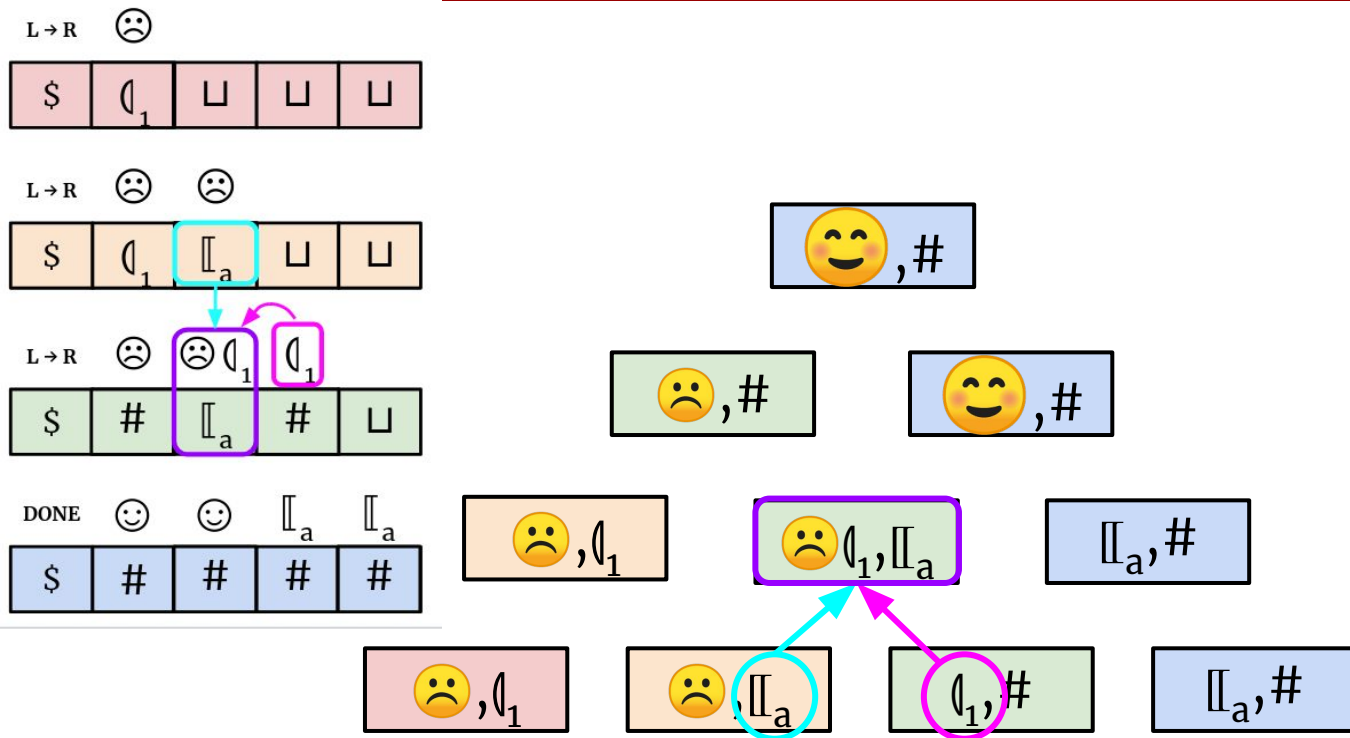
Flip Vertically...

# Flip it...



Small Clockwise Rotation...

# Rotate... gives a Triangular Trellis Automaton!!



In DSTM for  $\mathbb{D}_{CD}$

$$\delta(q_1, \llbracket_a, \epsilon) =$$

$$(\text{😞 } q_1, \llbracket_a)$$

# Triangular Trellis Automata

6-tuple  $K = (\Sigma, L, Q_1, I_1, \delta_1, F)$

$\Sigma$ : input alphabet

$L$ : node labels

$Q_1$ : set of states for  $L$ -labeled nodes

$I_1: \Sigma \rightarrow Q_1$  (function generating bottom row states from input)

$\delta: Q_{11} \times Q_{12} \rightarrow Q_1$  (transition function where  $Q_{11}/Q_{12}$  is left/right child)

$F \subseteq Q_1$ : set of accepting states

# Quick Trellis Example

Input: “b b a”

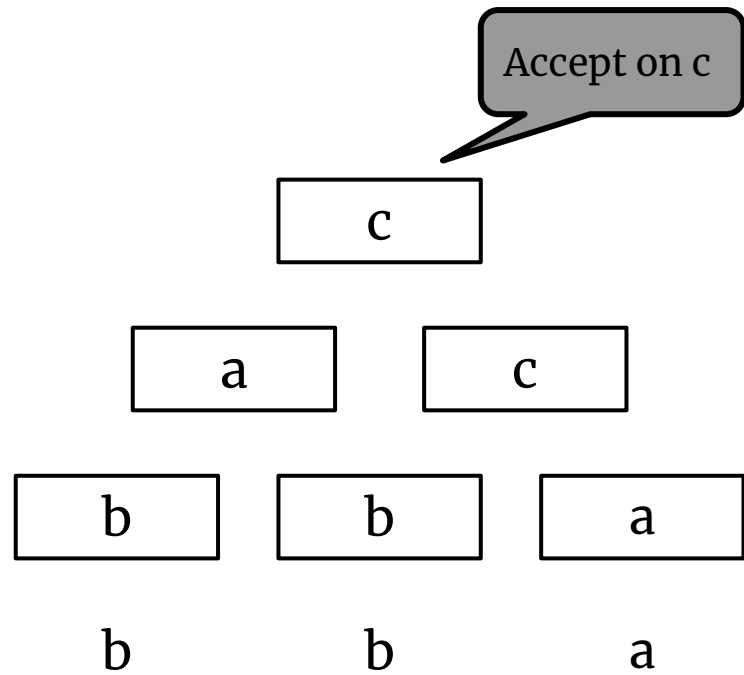
$F: \{c\}$

$I_1$  : identity

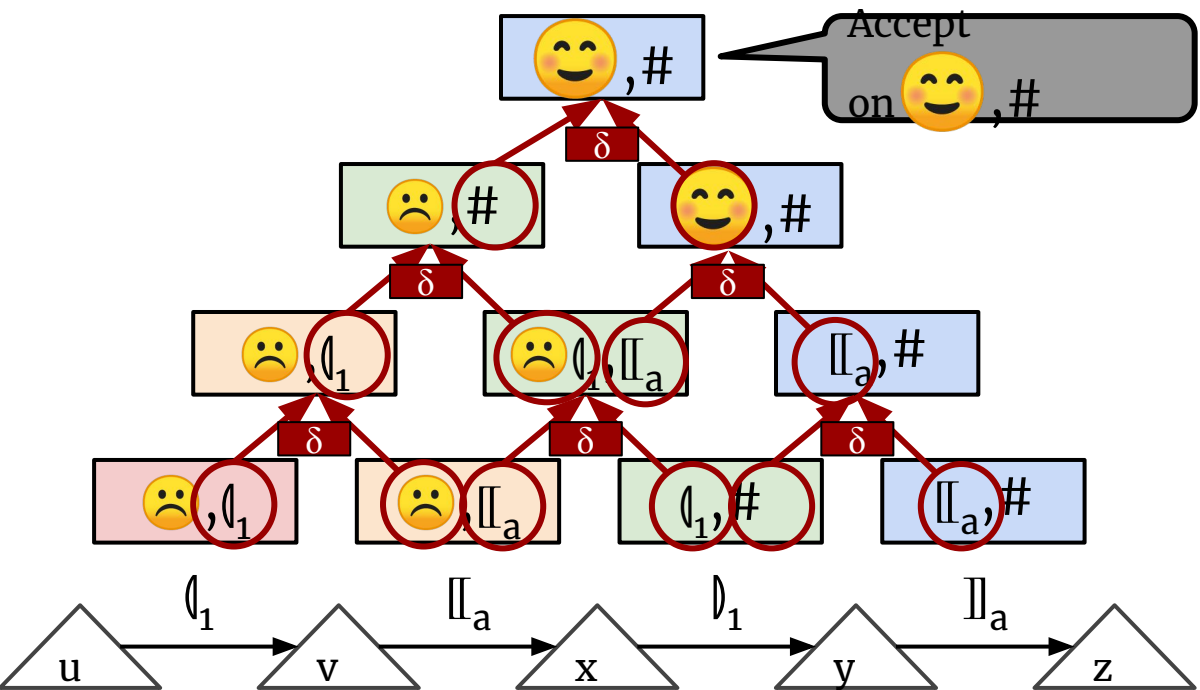
$\delta_1$ :

Right Child

	a	b	c
a	a	a	c
b	c	a	b
c	c	b	a



# Trellis for Computing membership of $(q_1 \parallel_a) q_1 \parallel_a$



## LCL Rules Part 2:

Trellis Transition:  
 $\delta_1((q_1, Z_1), (q_2, Z_2)) = q, Z$

If  $\delta(q_2, Z_1, \epsilon) = (q, Z, L)$   
 in DSTM



# Algorithm 1 (Baseline) Overview

- Input: Graph  $G = (V, E)$  and LCL rules for  $\mathbb{D}_{\text{CD}}$
- Output: Set  $\{S(n_1, n_2) \mid n_1, n_2 \in V\}$
- Worklist algorithm
  - Worklist consists of **Summary Edges**  $S(n_1, n_2) = (q, Z)$
  - $(q, Z)$  is result of applying Trellis transition function on a path  $n_1 \rightarrow \dots \rightarrow n_2$
  - Brute-force applies Triangular Trellis Transition function to corresponding values
- Approximation in  $O(|M|mn)$  time with  $O(n^2)$  space
  - $|V| = n, |E| = m$
  - $|M| = O(q^2Z + qZ^2)$  where  $q/Z$  are number of states/tape symbols in DSTM  $M$

# Optimizations (Algorithm 2)

- Sanity checks Heuristics for “Feasibility”
- Using Gray-White Trellis Automata (GWTA)
  - Gray nodes:  $(q, Z)$  such that  $q \in \{\text{😊}, \text{😞}\}$
  - Allows filtering of “spurious” nodes
- Approximation in  $O(|\mathbb{D}|mn)$  time with  $O(n^2)$  space
  - $|\mathbb{D}| = O((c+d)cd)$  where  $c/d$  is size of  $\mathbb{D}_C/\mathbb{D}_D$

# Benchmarks & Evaluation

- Apply Algorithm 2 to two different analysis
  - Java Alias Analysis
    - field sensitivity edge  $(u, \llbracket_f, v)$ : field  $f$  of object  $u$  may point to  $v$
  - Android app Taint analysis
    - field sensitivity edge  $(u, \llbracket_f, v)$ : value of  $u$  flows to the field  $f$  of variable  $v$
- Comparison with traditional CFL-reachability algorithm
  - CIFS: Context-insensitive field-sensitive
  - CSFI: Context-sensitive field-insensitive
  - CRFS: Regular-approximating context-sensitive field-sensitive
  - CSFR: Context-sensitive Regular-approximating field-sensitive

# Evaluation on Alias Analysis

Program	Precision (#S-pair)			Time (s)			Space (MB)		
	CIFS	CSFI	Lin	CIFS	CSFI	Lin	CIFS	CSFI	Lin
antlr	2,153,091	1,827,705	1,519,527	165.46	768.11	7.76	674.42	913.20	3822.73
bloat	2,305,351	1,904,837	1,588,287	195.15	1161.30	11.38	800.62	1142.45	5632.24
chart	6,181,087	5,746,258	4,868,772	1170.75	7190.24	70.96	1588.62	2405.58	18024.00
eclipse	2,300,305	1,826,423	1,515,782	190.66	794.00	9.20	709.95	948.41	4218.52
fop	4,733,519	4,069,404	3,575,835	824.71	4393.99	45.07	1443.70	2166.72	14220.40
hsqldb	2,094,016	1,814,974	1,513,500	149.07	675.71	7.46	630.21	834.00	3534.08
jython	2,786,450	1,878,085	1,557,348	243.20	1225.47	13.33	840.79	1271.31	5420.17
luindex	2,151,119	1,820,530	1,513,628	161.44	698.43	7.33	657.75	866.41	3484.07
lusearch	2,199,729	1,824,220	1,517,576	171.55	718.47	9.75	644.55	895.12	3909.12
pmd	2,353,255	1,879,584	1,570,281	191.28	796.97	10.22	709.25	950.11	4181.22
xalan	2,078,462	1,814,712	1,511,005	147.38	665.81	7.47	622.53	828.00	3390.27

# Evaluation on Taint Analysis

Program	Precision (#S-pairs)					Time (s)					Space (MB)				
	CIFS	CRFS	CSFI	CSFR	Lin	CIFS	CRFS	CSFI	CSFR	Lin	CIFS	CRFS	CSFI	CSFR	Lin
Backflash	32,081	11,679	7,115	6,819	597	0.09	996.11	0.44	7.10	0.02	11.4	2,054.8	14.0	56.7	16.1
BatteryDoctor	109,662	-	15,978	-	3,071	1.86	-	2.25	-	0.12	43.6	-	47.6	-	97.0
DroidKungFu	41,072	-	11,813	9,523	1,510	0.44	-	0.89	178.34	0.03	18.6	-	19.9	686.3	23.8
Fakebanker	12,098	4,439	2,463	2,363	542	0.08	180.40	0.08	8.90	0.01	9.4	1,038.5	10.4	138.1	9.6
Fakedaum	59,104	-	6,480	6,237	1,560	0.50	-	0.59	206.43	0.04	24.3	-	27.4	1,246.2	42.9
Faketaobao	3,196	1,179	732	676	274	0.02	5.28	0.02	0.34	0.00*	4.7	163.9	5.5	21.5	3.4
Jollyserv	22,960	12,449	1,463	1,182	801	0.19	937.41	0.06	5.91	0.02	10.9	1,519.1	11.2	200.3	11.9
Loozfon	3,044	1,865	646	618	218	0.01	2.15	0.02	0.06	0.00*	3.4	37.5	4.2	5.0	2.6
Phospy	-	-	-	-	1,158K	-	-	-	-	35.57	-	-	-	-	3,383.3
Roidsec	81,485	-	18,598	16,592	611	0.38	-	1.88	55.49	0.02	-	16.7	18.1	167.8	16.8
Scipix	976,255	-	71,759	-	146,913	333.50	-	21.69	-	2.74	476.2	-	93.3	-	478.6
simhosy	1,711,493	-	171,052	-	30,736	833.30	-	177.33	-	1.46	1,242.8	-	262.7	-	670.9
Skullkey	-	-	-	-	2,703K	-	-	-	-	190.40	-	-	-	-	23,290.8
Uranai	24,802	11,379	1,062	874	587	0.07	554.68	0.07	0.90	0.02	9.9	1,210.1	13.0	47.6	14.0
Zertsecurity	24,534	9,361	2,512	1,705	479	0.08	131.64	0.06	1.06	0.01	7.5	453.5	7.0	26.6	6.0

# Summary

- Context-Sensitive Data-Dependent Static Analysis can be formulated as solving Interleaved Dyck Language  $\mathbb{D}_{\text{CD}}$  on paths of graphs
  - Reachability on graphs wrt  $\mathbb{D}_{\text{CD}}$  is undecidable!
- DSTM/Triangular Trellis lend to sound approximation algorithms
  - $\mathbb{D}_{\text{CD}}$  is subset of Linear Conjunctive Languages (LCL)
  - Reachability on graphs wrt LCL is undecidable!
  - $O(mn)$  time approximation algorithm using ideas from Triangular Trellis