

Optimal Dyck Reachability for Data-dependence and Alias Analysis

Krishnendu Chatterjee Bhavya Choudhary
Andreas Pavlogiannis

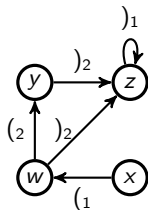


Static Analysis via Dyck Reachability

$$\Sigma = \{(1,)_1, \dots, (k,)_k\} \cup \{\epsilon\}$$

$$G = (V, E, \lambda : E \rightarrow \Sigma)$$

$$\mathcal{S} \rightarrow \mathcal{S} \mathcal{S} \mid (1 \mathcal{S})_1 \mid \dots \mid (k \mathcal{S})_k \mid \epsilon$$

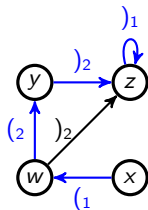


$$\Sigma = \{(1,)_1, \dots, (k,)_k\} \cup \{\epsilon\}$$

$$G = (V, E, \lambda : E \rightarrow \Sigma)$$

$$\mathcal{S} \rightarrow \mathcal{S} \mathcal{S} \mid (1 \mathcal{S})_1 \mid \dots \mid (k \mathcal{S})_k \mid \epsilon$$

$$P : x \rightsquigarrow z \quad \text{with} \quad \lambda(P) = (1(2)2)_1$$

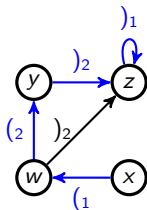


$$\Sigma = \{(1,)_1, \dots, (k,)_k\} \cup \{\epsilon\}$$

$$G = (V, E, \lambda : E \rightarrow \Sigma)$$

$$\mathcal{S} \rightarrow \mathcal{S} \mathcal{S} \mid (\mathcal{S})_1 \mid \dots \mid (\mathcal{S})_k \mid \epsilon$$

$$P : x \rightsquigarrow z \quad \text{with} \quad \lambda(P) = (1(2)2)_1$$



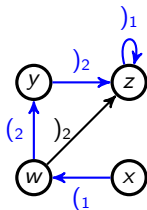
- Alias analysis
- Data-dependence analysis
- Data-flow analysis
- Shape analysis
- Impact analysis
- Bloat analysis
- Type-based flow analysis
- Program slicing

$$\Sigma = \{(1,)_1, \dots, (k,)_k\} \cup \{\epsilon\}$$

$$G = (V, E, \lambda : E \rightarrow \Sigma)$$

$$\mathcal{S} \rightarrow \mathcal{S} \mathcal{S} \mid (\mathcal{S})_1 \mid \dots \mid (\mathcal{S})_k \mid \epsilon$$

$$P : x \rightsquigarrow z \quad \text{with} \quad \lambda(P) = (1(2)2)_1$$



- **Alias analysis**
- **Data-dependence analysis**
- Data-flow analysis
- Shape analysis
- Impact analysis
- Bloat analysis
- Type-based flow analysis
- Program slicing

Q: Is v Dyck-reachable from u ?

- Dyck languages are a class of CFL
- Solution similar to CYK for CFL parsing
- $O(n^3)$
- $O(n^3/\log n)$ [Chaudhuri, POPL '08]
- Often prohibitive for lightweight static analysis

Q: Is v Dyck-reachable from u ?

- Dyck languages are a class of CFL
- Solution similar to CYK for CFL parsing
- $O(n^3)$
- $O(n^3/\log n)$ [Chaudhuri, POPL '08]
- Often prohibitive for lightweight static analysis

Can do better if the graph is simple!

Data-dependence analysis in the presence of libraries

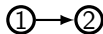
- Which variable depends on which others
- Identify Def-Use chains in a program
- LHS depends on RHS

Data-dependence analysis in the presence of libraries

- Which variable depends on which others
- Identify Def-Use chains in a program
- LHS depends on RHS

y depends on x if $x \rightsquigarrow y$

```
1 x ← 2
2 y ← x + 1
```



Analysis in the Presence of Libraries

- Large library
- Small client
- Goal:
 - Preprocess library once
 - Small cost of client analysis

Solution(?): Dyck reachability on the library once

Callbacks!

$: f_1(x, y)$

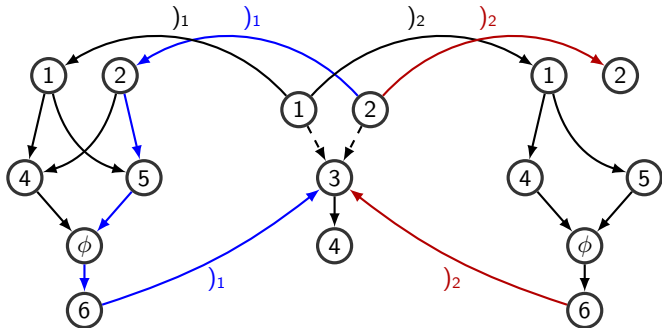
```
1 if  $y \% 2 = 1$  then
2    $z \leftarrow x + y$ 
3 else
4    $z \leftarrow x \cdot y$ 
5 end
6 return  $z$ 
```

$: g()$

```
1  $x \leftarrow 2$ 
2  $y \leftarrow 2$ 
3  $p \leftarrow f(x, y)$ 
4 return  $p$ 
```

$: f_2(x, y)$

```
1 if  $x \% 2 = 1$  then
2    $z \leftarrow 2 \cdot x$ 
3 else
4    $z \leftarrow 2 \cdot x + 1$ 
5 end
6 return  $z$ 
```



Results

- Library size n_1
- Library call sites k_1
- Client size n_2
- Client call sites k_2

$$n_1 \gg n_2 \quad n_i \gg k_i$$

Approach	Time	
	Library	Client
CFL	$O(n_1^3)$	$O((n_1 + n_2)^3)$
TAL*	$O(n_1^6)$	$O((s + n_2)^6)$

- * [TWZXZM, POPL '15]

- Library size n_1
- Library call sites k_1

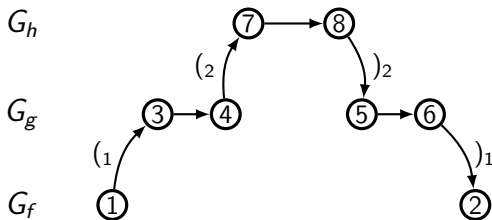
- Client size n_2
- Client call sites k_2

$$n_1 \gg n_2$$

$$n_i \gg k_i$$

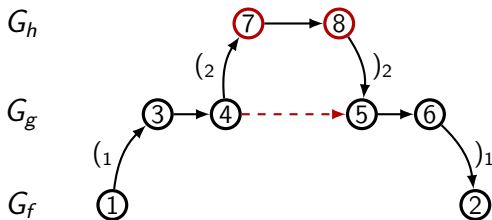
Approach	Time	
	Library	Client
CFL	$O(n_1^3)$	$O((n_1 + n_2)^3)$
TAL*	$O(n_1^6)$	$O((s + n_2)^6)$
This Paper	$O(n_1 + k_1 \cdot \log n_1)$	$O(n_2 + k_1 \cdot \log n_1 + k_2 \cdot \log n_2)$

- * [TWZXZM, POPL '15]



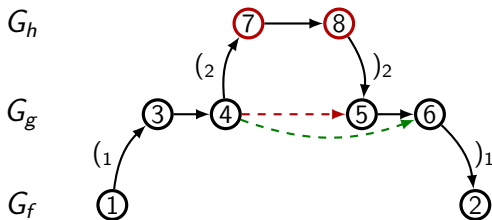
- 1 Reachability ignoring parenthesis edges
- 2 If Entry \rightsquigarrow Exit
 - Insert summary edge

Algorithmic Principle



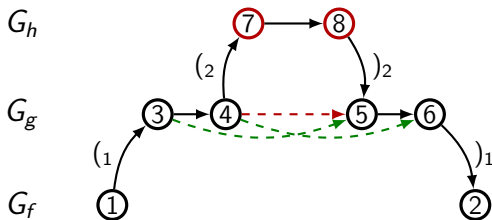
- 1 Reachability ignoring parenthesis edges
- 2 If Entry \rightsquigarrow Exit
 - Insert summary edge

Algorithmic Principle



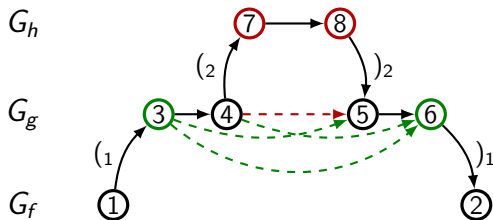
- 1 Reachability ignoring parenthesis edges
- 2 If Entry \rightsquigarrow Exit
 - Insert summary edge

Algorithmic Principle



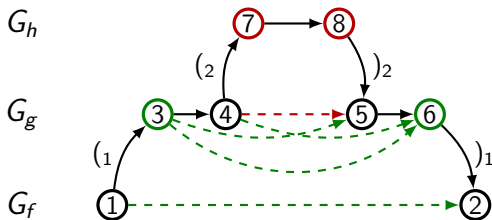
- 1 Reachability ignoring parenthesis edges
- 2 If Entry \rightsquigarrow Exit
 - Insert summary edge

Algorithmic Principle



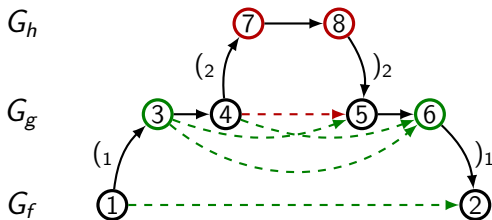
- 1 Reachability ignoring parenthesis edges
- 2 If Entry \rightsquigarrow Exit
 - Insert summary edge

Algorithmic Principle



- 1 Reachability ignoring parenthesis edges
- 2 If Entry \rightsquigarrow Exit
 - Insert summary edge

Algorithmic Principle



1 Reachability ignoring parenthesis edges

2 If Entry \rightsquigarrow Exit

- Insert summary edge

- Most time spent in (1)

- Each G_i is a “simple graph”

- Low treewidth

Given a graph G

- An integer parameter $t_G \geq 1$
- Measures how “similar” G is to a tree
- G is a tree iff $t_G = 1$
- Trees are very simple to handle
- Low-treewidth graphs are somewhat simple to handle

Given a graph G

- An integer parameter $t_G \geq 1$
- Measures how “similar” G is to a tree
- G is a tree iff $t_G = 1$
- Trees are very simple to handle
- Low-treewidth graphs are somewhat simple to handle

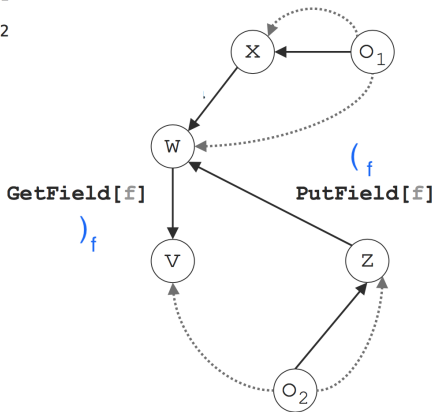
Theorem

For a low-treewidth graph we can construct a dynamic reachability oracle that supports

- *edge insertions in $O(\log n)$ time*
- *reachability queries in $O(\log n)$ time*

May Alias Analysis

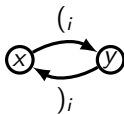
```
x = new Obj(); // o1  
z = new Obj(); // o2  
w = x;  
w.f = z;  
v = w.f;
```



$o_1 \xrightarrow{\dots} x$: object o_1 flows to pointer x
[YE, ESOP '09] \Leftrightarrow pointer x points to object o_1

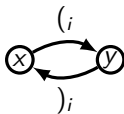
Bidirected graphs

$$\forall u, v \in V : \lambda(u, v) = (i \iff \lambda(v, u) =)i$$



Bidirected graphs

$$\forall u, v \in V : \lambda(u, v) = (i \iff \lambda(v, u) =)_i$$



- Demand-driven field-sensitive alias analysis
- [YXR, ISSTA '11] [ZLYS, PLDI '13]

Graphs of n nodes

Source	Worst-case Time	Average-case Time
[ZLYS, PLDI '13]	$O(n^2)$	$O(n \cdot \log n)$

Graphs of n nodes, $\alpha(n)$ the inverse Ackermann

Source	Worst-case Time	Average-case Time
[ZLYS, PLDI '13]	$O(n^2)$	$O(m \cdot \log n)$
This paper	$O(n \cdot \alpha(n))$	$O(n)$

Graphs of n nodes, $\alpha(n)$ the inverse Ackermann

Source	Worst-case Time	Average-case Time
[ZLYS, PLDI '13]	$O(n^2)$	$O(m \cdot \log n)$
This paper	$O(n \cdot \alpha(n))$	$O(n)$

optimal!

optimal!

Reachability is Equivalence

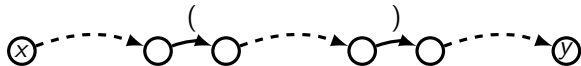
Theorem (ZLYS, PLDI '13)

Dyck reachability on bidirected graphs is an equivalence relation.

Reachability is Equivalence

Theorem (ZLYS, PLDI '13)

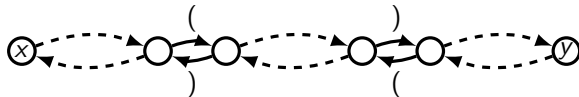
Dyck reachability on bidirected graphs is an equivalence relation.



Reachability is Equivalence

Theorem (ZLYS, PLDI '13)

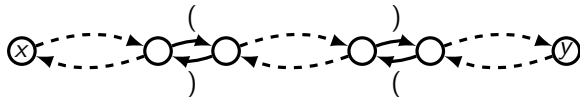
Dyck reachability on bidirected graphs is an equivalence relation.



Reachability is Equivalence

Theorem (ZLYS, PLDI '13)

Dyck reachability on bidirected graphs is an equivalence relation.

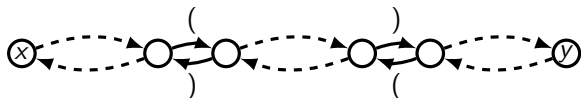


- Compute Strongly Connected Components

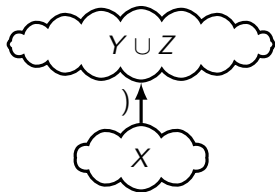
Reachability is Equivalence

Theorem (ZLYS, PLDI '13)

Dyck reachability on bidirected graphs is an equivalence relation.



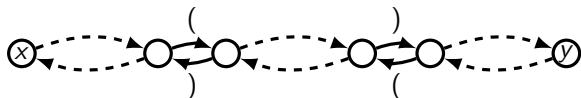
- Compute Strongly Connected Components



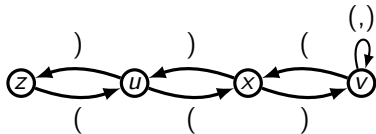
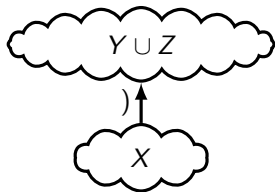
Reachability is Equivalence

Theorem (ZLYS, PLDI '13)

Dyck reachability on bidirected graphs is an equivalence relation.



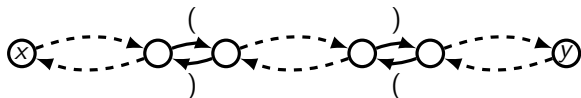
- Compute Strongly Connected Components



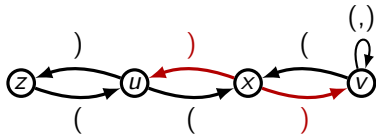
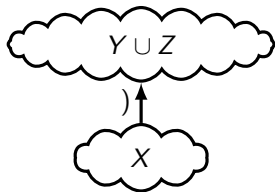
Reachability is Equivalence

Theorem (ZLYS, PLDI '13)

Dyck reachability on bidirected graphs is an equivalence relation.



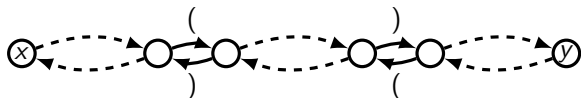
- Compute Strongly Connected Components



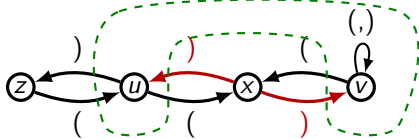
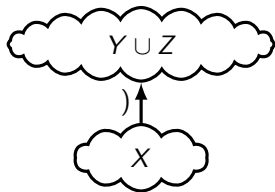
Reachability is Equivalence

Theorem (ZLYS, PLDI '13)

Dyck reachability on bidirected graphs is an equivalence relation.



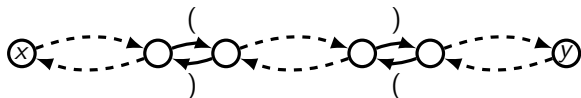
- Compute Strongly Connected Components



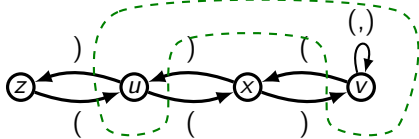
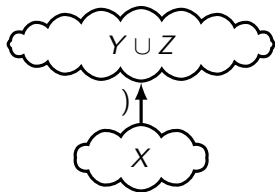
Reachability is Equivalence

Theorem (ZLYS, PLDI '13)

Dyck reachability on bidirected graphs is an equivalence relation.



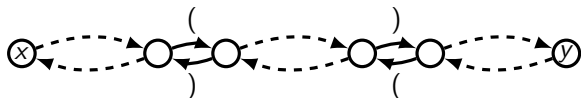
- Compute Strongly Connected Components



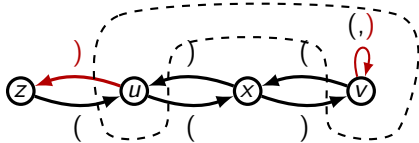
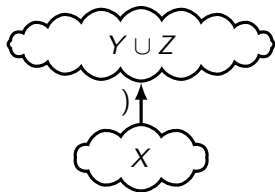
Reachability is Equivalence

Theorem (ZLYS, PLDI '13)

Dyck reachability on bidirected graphs is an equivalence relation.



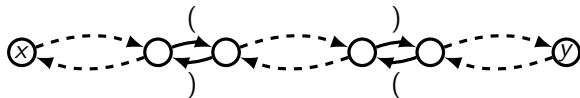
- Compute Strongly Connected Components



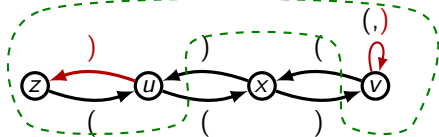
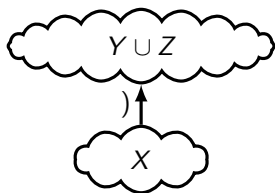
Reachability is Equivalence

Theorem (ZLYS, PLDI '13)

Dyck reachability on bidirected graphs is an equivalence relation.



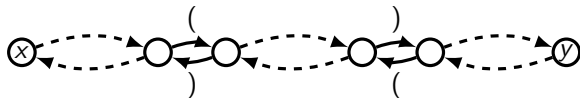
- Compute Strongly Connected Components



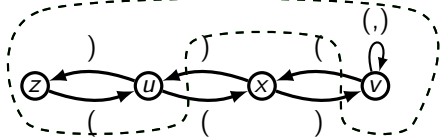
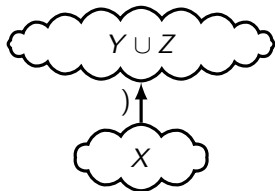
Reachability is Equivalence

Theorem (ZLYS, PLDI '13)

Dyck reachability on bidirected graphs is an equivalence relation.



- Compute Strongly Connected Components



Alias Analysis

- Implementation in C++
- Compared with [ZLYS, PLDI '13]
- DaCapo-2006 benchmarks
- SPGs from [YXR, ISSTA '11]

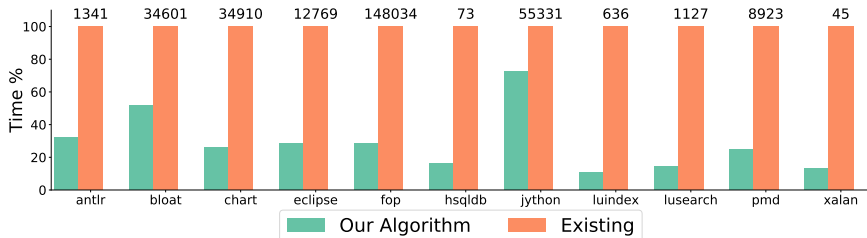
Alias Analysis

- Implementation in C++
- Compared with [ZLYS, PLDI '13]
- DaCapo-2006 benchmarks
- SPGs from [YXR, ISSTA '11]

Data-dependence Analysis

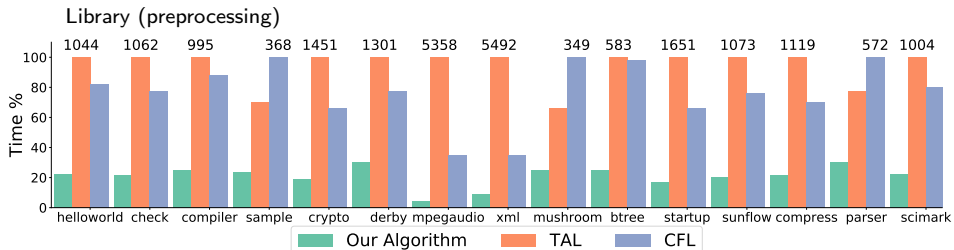
- Implementation in Java
- Compared with TAL reachability [TWZXZM, POPL '15]
- Benchmarks:
 - SPECjvm2008
 - 4 randomly chosen GitHub projects

Time-usage comparison (ms)



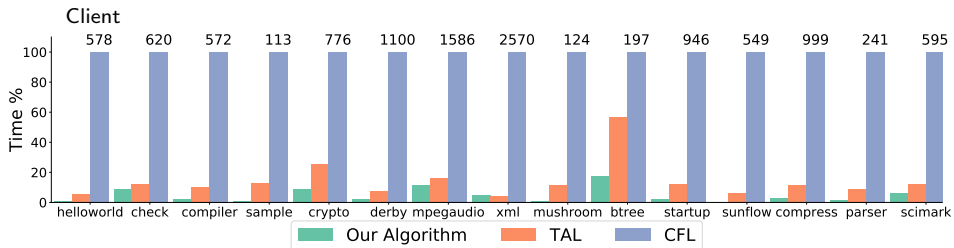
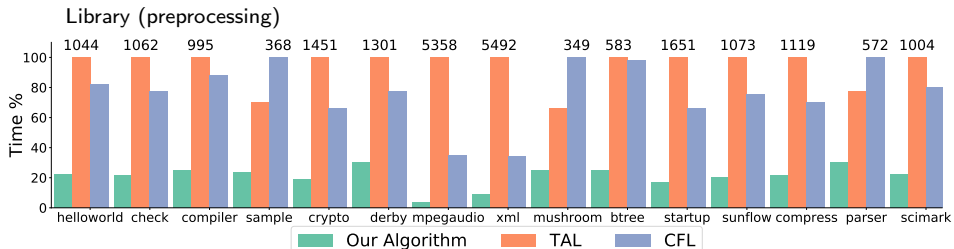
Data-dependence Analysis

Time-usage comparison (ms)



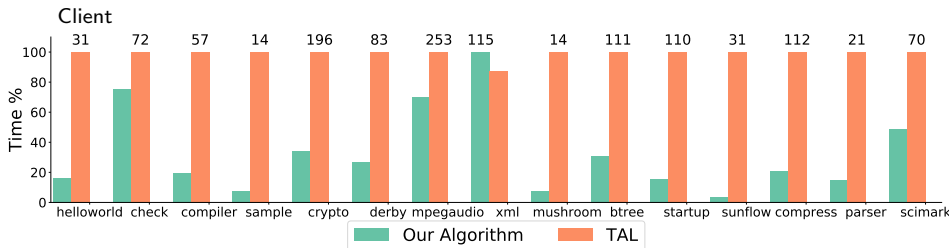
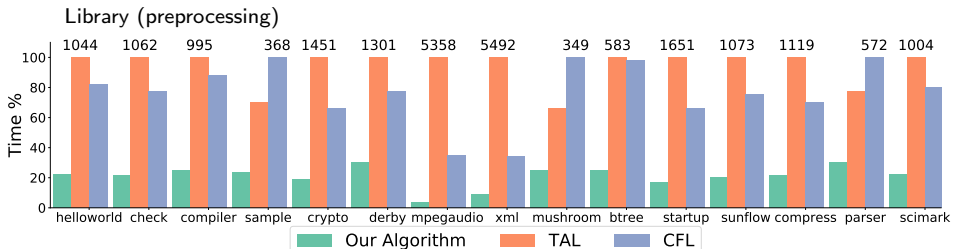
Data-dependence Analysis

Time-usage comparison (ms)



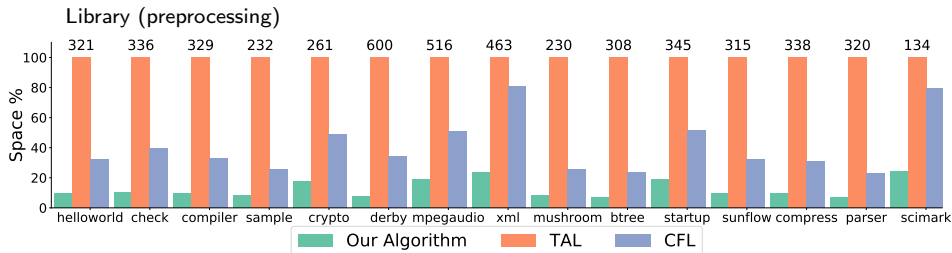
Data-dependence Analysis

Time-usage comparison (ms)



Data-dependence Analysis

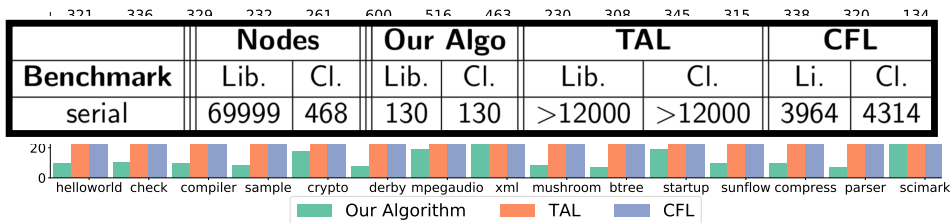
Space-usage comparison (MB)



Data-dependence Analysis

Space-usage comparison (MB)

Library (preprocessing)



- Graphs in Alias and Data-dependence analysis have special structure
- Faster, better algorithms for Dyck reachability
- Optimality guarantees

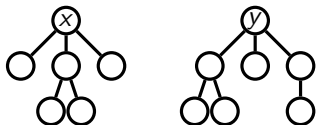
Thank you!
Questions?

Need to make SCC merging fast

Need to make SCC merging fast

Merging nodes

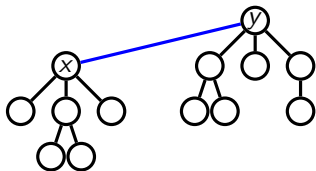
- Disjoint-sets data structure
- $O(\alpha(n))$ time for merging and querying



Need to make SCC merging fast

Merging nodes

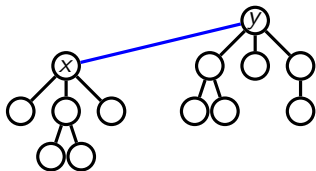
- Disjoint-sets data structure
- $O(\alpha(n))$ time for merging and querying



Need to make SCC merging fast

Merging nodes

- Disjoint-sets data structure
- $O(\alpha(n))$ time for merging and querying



Merging edges

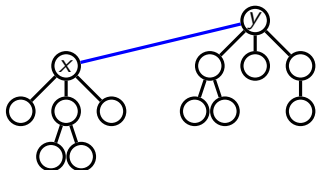
- Only outgoing!
- Linked list representation
- $O(1)$ time for merging



Need to make SCC merging fast

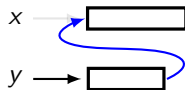
Merging nodes

- Disjoint-sets data structure
- $O(\alpha(n))$ time for merging and querying



Merging edges

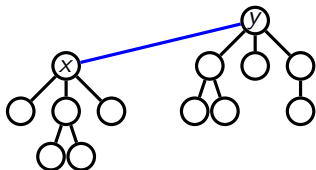
- Only outgoing!
- Linked list representation
- $O(1)$ time for merging



Need to make SCC merging fast

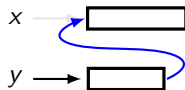
Merging nodes

- Disjoint-sets data structure
- $O(\alpha(n))$ time for merging and querying



Merging edges

- Only outgoing!
- Linked list representation
- $O(1)$ time for merging



Theorem

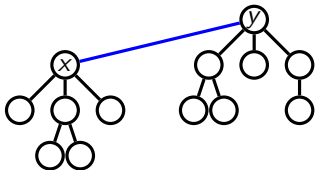
Dyck reachability on bidirected graphs solvable in

- $O(m + n \cdot \alpha(n))$ worst-case time
- $O(m)$ average-case time

Need to make SCC merging fast

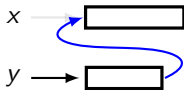
Merging nodes

- Disjoint-sets data structure
- $O(\alpha(n))$ time for merging and querying



Merging edges

- Only outgoing!
- Linked list representation
- $O(1)$ time for merging



Theorem

Dyck reachability on bidirected graphs solvable in

- $O(m + n \cdot \alpha(n))$ worst-case time, *optimally!*
- $O(m)$ average-case time, *optimally!*

Treewidth for Bidirected Graphs

Theorem

Dyck reachability on constant-treewidth, bidirected graphs requires $\Omega(m + n \cdot \alpha(n))$ time.

Treewidth for Bidirected Graphs

Theorem

Dyck reachability on constant-treewidth, bidirected graphs requires $\Omega(m + n \cdot \alpha(n))$ time.

Compare with

Theorem (ZLYZ, PLDI '13)

Dyck reachability on bidirected trees solvable in $O(n)$ time.

Typically complexity on low-treewidth graphs = complexity on trees

Complexity of Dyck reachability

- $O(n^3)$
- 2NPDA-hard
 - Conditional cubic lower bound

Complexity of Dyck reachability

- $O(n^3)$
- 2NPDA-hard
 - Conditional cubic lower bound

Theorem

Dyck reachability is Boolean Matrix Multiplication (BMM) - hard.

Complexity of Dyck reachability

- $O(n^3)$
- 2NPDA-hard
 - Conditional cubic lower bound

Theorem

Dyck reachability is Boolean Matrix Multiplication (BMM) - hard.

Corollary

An $O(n^{3-\delta})$ algorithm for Dyck reachability requires an $O(n^{3-\epsilon})$ combinatorial algorithm for BMM.

- $\delta, \epsilon > 0$

Existing Dyck-reachability algorithms conditionally optimal
(s.t. BMM-hardness)

Complexity of Dyck reachability

- $O(n^3)$
- 2NPDA-hard
 - Conditional cubic lower bound

Theorem

Dyck reachability *on low-treewidth graphs* is Boolean Matrix Multiplication (BMM) - hard.

Corollary

An $O(n^{3-\delta})$ algorithm for Dyck reachability *on low-treewidth graphs* requires an $O(n^{3-\epsilon})$ combinatorial algorithm for BMM.

- $\delta, \epsilon > 0$

Existing Dyck-reachability algorithms conditionally optimal
(s.t. BMM-hardness)

Time-usage comparison (s)

Benchmark	Fields	Nodes	Edges	Our Algo	Existing Algo
antlr	172	13708	23547	0.428783	1.34152
bloat	316	43671	103361	17.7888	34.6012
chart	711	53500	91869	8.99378	34.9101
eclipse	439	34594	52011	3.62835	12.7697
fop	1064	101507	178338	42.5447	148.034
hsqldb	43	3048	4134	0.012899	0.073863
jython	338	56336	167040	40.239	55.3311
luindex	167	9931	14671	0.068013	0.636346
lusearch	200	12837	21010	0.163561	1.12788
pmd	357	31648	58025	2.21662	8.92306
xalan	41	2342	2979	0.006626	0.045144

Time-usage comparison (ms)

Benchmark	Nodes		TW		Our Algo		TAL		CFL	
	Lib.	Cl.	Lib.	Cl.	Lib.	Cl.	Lib.	Cl.	Li.	Cl.
helloworld	16003	296	5	3	229	5	1044	31	855	578
check	16604	3347	5	4	228	54	1062	72	821	620
compiler	16190	536	5	3	248	11	995	57	876	572
sample	3941	28	4	1	86	1	258	14	368	113
crypto	20094	3216	5	5	273	66	1451	196	961	776
derby	23407	1106	6	3	389	22	1301	83	1003	1100
mpegaudio	28917	27576	5	24	204	177	5358	253	1864	1586
xml	71474	2312	5	3	489	115	5492	100	1891	2570
mushroom	3858	7	4	1	86	1	230	14	349	124
btree	6710	1103	4	4	144	34	583	111	571	197
startup	19312	621	5	3	279	17	1651	110	1087	946
sunflow	15615	85	5	2	217	1	1073	31	811	549
compress	16157	1483	5	3	240	23	1119	112	783	999
parser	7856	112	4	1	172	3	443	21	572	241
scimark	16270	2027	5	5	220	34	1004	70	805	595
serial	69999	468	8	3	440	9	MEMOUT	MEMOUT	117147	165958

Space-usage comparison (MB)

Benchmark	Nodes		TW		Our Algo		TAL		CFL	
	Lib.	Cl.	Lib.	Cl.	Lib.	Cl.	Lib.	Cl.	Li.	Cl.
helloworld	16003	296	5	3	31	27	321	44	104	126
check	16604	3347	5	4	34	31	336	89	132	184
compiler	16190	536	5	3	31	28	329	44	108	137
sample	3941	28	4	1	19	16	232	59	59	64
crypto	20094	3216	5	5	45	45	261	61	127	188
derby	23407	1106	6	3	46	41	600	88	204	265
mpegaudio	28917	27576	5	24	96	96	516	219	262	397
xml	71474	2312	5	3	108	108	463	153	373	480
mushroom	3858	7	4	1	19	16	230	59	58	58
btrees	6710	1103	4	4	22	19	308	65	72	89
startup	19312	621	5	3	66	66	345	92	178	230
sunflow	15615	85	5	2	30	27	315	43	102	124
compress	16157	1483	5	3	32	29	338	50	105	131
parser	7856	112	4	1	22	19	320	64	73	83
scimark	16270	2027	5	5	32	29	134	49	106	140
serial	69999	468	8	3	130	130	MEMOUT	MEMOUT	3964	4314

Space-usage comparison (MB)

Benchmark	Nodes		TW		Our Algo		TAL		CFL	
	Lib.	Cl.	Lib.	Cl.	Lib.	Cl.	Lib.	Cl.	Li.	Cl.
helloworld	16003	296	5	3	31	27	321	44	104	126
check	16604	3347	5	4	34	31	336	89	132	184
compiler	16190	536	5	3	31	28	329	44	108	137
sample	3941	28	4	1	19	16	232	59	59	64
crypto	20094	3216	5	5	45	45	261	61	127	188
derby	23407	1106	6	3	46	41	600	88	204	265
mpegaudio	28917	27576	5	24	96	96	516	219	262	397
xml	71474	2312	5	3	108	108	463	153	373	480
mushroom	3858	7	4	1	19	16	230	59	58	58
btrees	6710	1103	4	4	22	19	308	65	72	89
startup	19312	621	5	3	66	66	345	92	178	230
sunflow	15615	85	5	2	30	27	315	43	102	124
compress	16157	1483	5	3	32	29	338	50	105	131
parser	7856	112	4	1	22	19	320	64	73	83
scimark	16270	2027	5	5	32	29	134	49	106	140
serial	69999	468	8	3	130	130	MEMOUT	MEMOUT	3964	4314

> 12GB > 12GB