The Work at a Glance

Problem

Why Interprocedural Analysis?

NULL BUG CHECKER [1] void function1() { int * ptr1; if (<*condition*>) { ptr1 = fn_explct_ret_null(); else { ptr1 = fnA();if (ptr1!=NULL) { Absence of this int b = *ptr1; check
is a bug.

Pattern Checker can detect that ptr1 can be NULL.



Pattern Checker cannot detect that **ptr1** can be NULL, we need interprocedural analysis, e.g. dataflow analysis.

The Scalability Problem

The scalability problem of interprocedural analysis stems from producing distinct solutions for different calling contexts.



No. of calling contexts grows **exponentially** with program size. A moderate-sized program can have **10¹⁴ distinct contexts**. [2]

Non-Parallelizability Most existing techniques frequently involve decision making based on information they discover dynamically.

Implementation Difficulty

To address scalability, practitioners approximate their analyses, however, implementing them is complicated. In Sridharan's and Bodik's [4] work, more than 75% of their entire code was dedicated to tuning the analysis.

Graspan A Big Data System for Big Code Analysis

Kai Wang, Aftab Hussain, Zhiqiang Zuo, Guoqing Xu, Ardalan Amiri Sani Department of Computer Science, University of California, Irvine

We address the scalability problem of inter-procedural static analysis for bug detection in big code. We built Graspan, a graph processing system that helped us uncover 85 new Null pointer bugs in Linux 4.4.0.

sum (a,b) return a + b; (11) sum (a,b) return a + b;

Graphs and Program Analysis? Thomas Reps et al. [3] showed most interprocedural analyses, like dataflow analysis and pointer analysis, can

be transformed to a graph reachability problem.

Queries can be solved using





Find more bugs on **big code** In your machine

Performance-wise Scalable Parallelizable RULES

Easy to Implement

Challenges Our Design Duplicate Edges

Overburdens memory. How to terminate edge addition process?

Repartitioning

How do we partition the graph, and then after computation, how do we repartition oversized partitions?

Future Ambitions

Extend system support for path-sensitive analysis, constraintbased analyses by encoding constraints into edge values.

Graspan

Quick Note:

dynamic transitive computation. Dataflow Analysis: Helps us find how data flows in a program across statements

> **Pointer Analysis**: Helps us find aliases, i.e., variables that point to the same memory location.



gives us dataflow info. and alias info., which can be fed into checkers to find more bugs.

Benefits How it works

Preprocessing

Partition the graphs.

Edge-Pair Centric Computation Analyze pairs of edges, and perform DTC computation.

Post-Processing Repartition oversized partitions according to a threshold.

Our Analysis

We implemented fully context-sensitive dataflow analysis and pointer analysis on these programs:

Program Linux PostgreSQL Apache httpd

Research Q&As

Can Interprocedural Analysis improve checkers? 85 new bugs in Linux 4.4.0.

Is Graspan Efficient and Scalable? Computations took ~2 to less than 12 hrs, largest graph generated had 1.1B edges.

Graspan implementation v/s old implementation? Use an API to provide a grammar.

Graspan v/s existing graph systems? GraphChi crashed in 133 secs with 65M edges added.

> **Example Bug** (missed by original checker)

References

analysis via graph reachability, POPL '95 to analysis for Java, PLDI '06

This work was done under the supervision of Guoqing Xu, University of California, Irvine. Other contributors of this work are Kai Wang, Zhiqiang Zuo, and Ardalan Amiri Sani. The work was submitted to OSDI 2016.

Results

	_			
		-		
•	1	l		

Version 4.4.0-rc5 8.3.9 2.2.18

#LOC 16M 700K 300K



[1] Nicolas Palix, Gaël Thomas, Suman Saha, Christophe Calvès, Julia Lawall, and Gilles Muller, Faults in linux: ten years later, ASPLOS '11 [2] John Whaley and Monica S. Lam, **Cloning-based context-sensitive pointer alias**

analysis using binary decision diagrams, PLDI '04 [3] Thomas Reps, Susan Horwitz, and Mooly Sagiv, **Precise interprocedural dataflow**

[4] Manu Sridharan and Rastislav Bodík, Refinement-based context-sensitive points-

Acknowledgements