Abstract Semantic Differencing via Speculative Correlation

NIMROD PARTUSH ERAN YAHAV TECHNION (HAIFA, ISRAEL) [OOPSLA'14]



European Research Council

Established by the European Commission

Supporting top researchers from anywhere in the world

Problem: Semantic Differencing

For two procedures P,P':

Prove that P and P' are equivalent
 Produce the same output for the same input

Otherwise, produce a useful description of the difference

Motivation

"Successful software always gets changed."
 – Frederick P. Brooks

 Program understanding & Debugging





Test generation & pruning



"Big Code"



Example: Sequence Printing



Prints a sequence of numbers from *first* to *last* in increments of *STEP*

Simplified code from Coreutils seq.c version 6.9

For instance for the input *first* = 1, *last* = 13, *STEP* = 2 the procedure will print the sequence 1, 3, 5, 7, 9, 11, 13

Example: Sequence Printing



- Does version 6.10 print the same sequence?
- Syntactic difference is vast
- But the output is the same: x is equivalent at the point of printing (9,10')

Test Equivalence

Run the procedures with the same inputs and check outputs for equivalence





Abstract Semantic Differencing

Use abstract interpretation to prove equivalence between two program versions

Or characterize their difference

Find (an abstraction of) differing programs states that come from the same input

Sound

Equivalence is guaranteed

Precise

Report few false differences

Equivalence Under Abstraction

Analyzing the programs separately may result in false equivalence

For instance, interpreting with a numerical relational abstraction:



Equivalence under abstraction does not entail equivalence between the concrete values it represents.

Our Approach

Analyze P and P' together

- Define a correlating semantics that interprets the programs together
- Interpretation is done in some interleaving of their steps
- Abstract the correlating semantics to handle infinite-state programs
- Search for the interleaving that allows the abstraction to best capture equivalence

Correlating Semantics

Maintain direct correlation between values in the programs P and P'

Use a relational abstraction that captures equivalences



Order Matters!

Analysis order determines the abstract correlating semantics' ability to track equivalence

- For example, in sequential order:
 - By the time one program's analysis is finished, the values have been abstracted and equivalence is lost

while
$$(i < n)$$
{
 $i \neq 1;$
 $print(i);$
}
Loop invariant:
 $i \neq 2$ }
 $i \neq 2$ }

In which order should the programs be analyzed?

Sequential?



In which order should the programs be analyzed?

Lock-Step?



```
static void
    print numbers v6 10(long first, long last, ...)
 3
      bool out of range = (last < first);</pre>
 4
      if (!out of range) {
 5
 6
        long x = first;
                                                         1
 7
        long i;
        for (i = 1; /* empty */ ; i++) {
8
          long x0 = x;
9
          printf (fmt, x);
10
11
          x = first + i * STEP;
12
          out of range = (last < x);</pre>
          if (out of range) break;
13
14
          fputs (separator, stdout);
15
16
        fputs (terminator, stdout);
17
      }
18
        print(1)
     static void
  1
     print numbers v6 9(long first, long las
  2
  3
  4
       long i;
  5
       for (i = 0; /* empty */; i++) \{
  6
         long x = first + i * STEP;
  7
         if (last < x) break;</pre>
         if (i) fputs (separator, stdout);
  8
 9
         printf (fmt, x);
 10
 11
       if (i)
 12
         fputs (terminator, stdout);
 13 }
```

In which order should the programs be analyzed?

- All possible interleavings?
 - Will result in an exposential blow-up



static void

The challenge is to find an interleaving that allows maintaining equivalence under abstraction

While avoiding exponential blow-up



Speculative Exploration

The search for an interleaving is part of the fixed-point abstract in Speculatively explore interpreting k steps, distributed over both programs
 The search driv
 0 over P & k over P'

- 1 over P & k 1 over P'
- Etc. k is the speculation window
- a parameter to the algorithm

while (worklist $\neq \emptyset$)

The algorithm is c

▶ results ← Speculate (P, P', worklist, state_⋈, k)

(worklist,state_⋈) ←
 Find Minimal Diff Result(P,P',results)

At the end of each speculative step, compare the k+1 results and pick the one with minimal difference

A greedy approach

Speculative Exploration: Example

> Speculate (k = 7)

```
static void
    print numbers v6 9(long first, long last, ...)
      long i;
      for (i = 0; /* empty */; i++) {
        long x = first + i * STEP;
        if (last < x) break;</pre>
        if (i) fputs (separator, stdout);
 9
        printf (fmt, x);
10
      }
11
      if (i)
12
        fputs (terminator, stdout);
13
   }
```



0 steps over v6.9, 7 steps over v6.10	•••	3 steps over v6.9, 4 steps over v6.10	 7 steps over v6.9, 0 steps over v6.10
$(4,11') \mapsto$ $\{i' = 1, i = ?,$ x' = first' + STEP', x = ?, $\}$	••••	$(7,7') \mapsto$ $\{i' = 1, i = 0, x' = x = first, \dots\}$	 $(6,4') \mapsto$ $\{i' = ?, i = 1, x' = ?, x = first,\}$ 17

Comparing Abstractions to find Minimal Difference

We explored two strategies

- 1. Scale-oriented: count the number of equivalent variables (denoted $\equiv_{\{v\}}$) per state
- 2. Precision-oriented: Use the abstract domain (APRON Polyhedra) geometrical representation to compute difference inclusion

Speculative Exploration: Visualization (k = 7)







Producing Description of Difference

- The speculative algorithm is geared towards finding minimal difference (and not strictly equivalence)
 - A usable description of difference is produced even if full equivalence does not hold



Evaluation

	Function	#LOC	#Patch	#Loops	Time
Corestilis	print_numbers	23	7-,13+	1	00:11 (k=2)
	cache_fstatat	17	2-,4+	0	00:03 (k=1)
	set_owner	51	2-,4+	0	00:02 (k=2)
	fmt	42	5-,5+	1	00:22 (k=2)
	md5sum	40	0-,3+	3	13:31 (k=2)
	char_to_clump	111	2-,12+	3	19:09 (k=2)
	savewd	86	0-,1+	0	00:46 (k=2)
	addr	77	1-,2+	0	00:17 (k=1)
	SetTextInternal	47	0-,3+	1	11:28 (k=3)
Ř-	get_sha1_basic v1	145	3-,10+	2	118:01 (k=2)
	get_sha1_basic v2	149	2-,20+	2	TO (2H)
	get_path_prefix	22	2-,3+	1	29:12 (k=3)
	boot_attr v1	77	7-,4+	0	08:08 (k=4)
	boot_attr v2	74	5-,7+	0	06:04 (k=4)
	read_attr	32	1-,4+	1	05:42 (k=2)
	ll_binary_merge	37	8-,24+	1	00:53 (k=1)
	write zip entry	340	1-,4+	3	07:32 (k=2)
Ned -	DDEC	10	3-,3+	1	00:13 (k=1)
	DSE	7	2-,3+	1	00:09 (k=1)
	RegVer	10	4-,4+	1	00:07 (k=1)
	SymDiff	32	5-,4+	0	00:04 (k=1)

Conclusion

Contact us! {nimi,yahave}@cs.technion.ac.il

 \blacktriangleright The interleaving matching problem is closely coupled with semantic diff & equivalence Previous approach mainly use syntactic cues New approach: search for an interleaving With an equivalence criteria to drive the search A description of difference is produced, instead of a binary yes/no for equivalence Useful for program understanding, debugging etc. Available on github

Proposed Questions

- 1. <u>Shouldn't there (still) be some exponential blow-up here??</u>
- 2. Where does the variable matching come from?
- 3. Can you talk about how related work find their interleaving?
- 4. How do you handle function calls?
- 5. <u>How do you handle heap\array data?</u>
- 6. Why were the k's in your evaluation so small?
- 7. Can you elaborate on the precise method for comparing abstractions?
- 8. You seem to have used a disjunctive domain, how did you \square and \square it?
- 9. Won't you miss interleavings?
- 10. How does the approach scale (inter-procedurally)?
- 11. Is the analysis proportional to the size of change\program\both?

Un-interpreted Functions

Function calls are handled modularly

$$x = foo(y,z)$$

- If foo was proven to be equivalent, and so are y & z, the result will be equivalent
 - Otherwise, anything (T) is ossible for x

$$\equiv_{foo} \equiv_{\{v_1, v_2\}}$$
$$\equiv_{foo}(v_1, v_2)$$

Array and heap access are modeled similarly

$$\equiv_{array} \equiv_{idx}$$

= array[idx]

Related work: program matching

Symbolic Execution based approaches [DSE, UC-KLEE]:

- Bound loop iteration
- try to explore all bound paths
- Check equivalence on each path
- Recursion-Rule based approaches [SymDiff]:
 - Transform loops to function calls
 - Use function calls as matching locations
 - i.e. try and prove the inputs to each function call are equivalent
- Data\Trace based approaches [DDEC]:
 - Use run traces variable values to infer a bi-simulation relation of the two programs

Equivalence-based Partitioning

We use a disjunctive domain to allow separating equivalent paths from differing paths

{ $print_extra_number, x' = x, i' = i + 1, x \le last$ }

 $\{\neg print_extra_number, x' = x + STEP, i' = i + 1, x \le last\}$

How do you Join such a domain?

which sub-states are joined with which?

We Join and Widen abstract states based on the equivalences they preserve

- the set of variables that hold equivalence
- disjunction size bound at 2^{|VAR|}
- Iose some information, but maintain what's important (equivalence)





Partial Order Reduction



- Since each step can be performed on either P or P'
- For the most part, a single representative is sufficient
 - Meaning foreach $0 < i \le k$ we explore first interpreting all *i* steps over P and then k i over P'
 - No alternation in-between
 - Resulting in k + 1 results
 - Since the domain is commutative and join\widen is only performed after each speculative step
- This does mean we miss some interleavings
 - But did not pose an issue in our benchmarks
 - A small sacrifice to make for scaling

Full Program Equivalence Checking

- For full program equivalence, we currently use the modular approach
 - And in the case the callee differs among versions, we assume T
 - ▶ Future: use the description of difference instead, somehow
- This is the state of the art for equivalence checking
 - The only full-program approaches (known to this presenter :) are symbolic execution based ones, that try to pin-point differing paths
 - Usually get poor coverage
 - Generally unable to prove equivalence

Interesting Diff

