

Algorithm 1 DSI usage in our preliminary study

Inputs: P : a program, T : set of unit tests for P ,
 $T_{map} : \{t \rightarrow S\}$
 Each S_i in T_{map} is the set of specs mined from test t_i

Outputs: $S^f = \{s_1^f, s_2^f, \dots\}$, a set of spurious specs,
 $S^t = \{s_1^t, s_2^t, \dots\}$, a set of likely true specs

```

1:
2: procedure runDSI( $P, T, T_{map}$ )  $\triangleright$  Main procedure
3:   results  $\leftarrow \{\}$   $\triangleright$  Set of Tuples ( $S_{test}^f, S_{test}^t$ )
4:   for all  $\ell$  in {all, class, method} do
5:     for all test in tests( $\ell, T$ ) do  $\triangleright$  tests( $\ell, T$ ): level  $\ell$  tests
6:        $S_{test}^f, S_{test}^t \leftarrow \text{validate}(\text{test}, T_{map}[\text{test}], P)$ 
7:       results  $\leftarrow$  results  $\cup (S_{test}^f, S_{test}^t)$ 
8:   return  $S^f, S^t$ 
9:
10: procedure validate(test,  $T_{map}[\text{test}], P$ )  $\triangleright$  DSI procedure
11:    $S^f \leftarrow \{\}$ ,  $S^t \leftarrow \{\}$   $S^u \leftarrow \{\}$   $\triangleright$   $S^f$ : spurious specs,  $S^t$ :
    likely specs,  $S^u$ : specs that DSI could not invalidate
12:   for all  $s$  in  $T_{map}[\text{test}]$  do
13:     out1  $\leftarrow$  run( $P, \text{test}$ )  $\triangleright$  out*  $\in \{\text{pass}, \text{fail}, \text{crash}\}$ 
14:     out2, trace1  $\leftarrow$  InstrRun( $P, \text{test}$ )  $\triangleright$  1st instrumented run
15:     out3, trace2  $\leftarrow$  InstrRun( $P, \text{test}$ )  $\triangleright$  2nd instrumented run
16:     if sanityCheck( $s, \text{out1}, \text{out2}, \text{out3}, \text{trace1}, \text{trace2}$ ) then
17:       exps  $\leftarrow$  experiments( $s, \text{trace}$ )  $\triangleright$  Algorithm 1 in [1]
18:     else
19:        $S^u \leftarrow S^u \cup \{s\}$ ; continue
20:   runExperiments( $s, P, \text{test}, S^u, S^f, S^t, \text{exps}$ )
21:   return  $S^f, S^t$ 
22:
23: procedure runExperiments( $s, P, I, S^u, S^f, S^t, \text{exps}$ )
24:   for all  $e : (\text{idx}, \text{len})$  in exps do
    Violate  $s$ : delay call at trace index 'e.idx' by 'e.len' calls
25:   out4, stage  $\leftarrow$  runExp( $P, I, e$ )
26:   if stage == 3  $\wedge$  out4 == pass then
27:      $S^f \leftarrow S^f \cup \{s\}$ ; continue
28:   if stage == 0 then  $S^u \leftarrow S^u \cup \{s\}$ ; continue
29:   else  $S^t \leftarrow S^t \cup \{s\}$ 
30:
31: procedure sanityCheck( $s, o1, o2, o3, t1, t2$ )
32:   if o1  $\in \{\text{fail}, \text{crash}\}$  then return false  $\triangleright$  Base run crash
33:   if o2  $\in \{\text{fail}, \text{crash}\}$  then return false  $\triangleright$  Instrument fail
34:   if o2  $\neq$  o3 then return false  $\triangleright$  Nondeterministic runs
35:   if  $s \not\models t1 \vee s \not\models t2$  then return false  $\triangleright$  Traces not in  $s$ 
36:   return true

```

Algorithm 2 DSI usage in TEMARI

Inputs: P : a program, T : the whole test suite for P ,
 S : a two-letter spec

Outputs: “buggy” if S is likely buggy, and
 “unknown” otherwise.

```

1:
2: procedure runDSI( $P, T, S$ )  $\triangleright$  Run DSI with unit tests
3:   out1  $\leftarrow$  run( $P, T$ )  $\triangleright$  out*  $\in \{\text{pass}, \text{fail}, \text{crash}\}$ 
4:   out2, trace1  $\leftarrow$  InstrRun( $P, T$ )  $\triangleright$  First instrumented run
5:   out3, trace2  $\leftarrow$  InstrRun( $P, T$ )  $\triangleright$  Second instrumented run
6:   if sanityCheck( $S, \text{out1}, \text{out2}, \text{out3}, \text{trace1}, \text{trace2}$ ) then
7:     exps  $\leftarrow$  experiments( $S, \text{trace}$ )  $\triangleright$  Algorithm 1 in [1]
8:   else
9:     return unknown
10:  return runDSIExperiments( $S, P, T, \text{exps}$ )
11:
12: procedure runDSIExperiments( $S, P, T, \text{exps}$ )
13:    $V \leftarrow \{\}$   $\triangleright$  Collect the set of DSI verdicts
14:   for all  $e : (\text{idx}, \text{len})$  in exps do
     $\triangleright$  Violate  $S$ ; delay call at trace index 'e.idx' by 'e.len' calls
15:   out4, stage  $\leftarrow$  runExp( $P, T, e$ )  $\triangleright$  mutate  $P$ , does  $T$  pass?
16:   if out4 == OK then  $V \leftarrow V \cup \{\text{buggy}\}$ 
17:   if stage  $\in \{0, 1, 2, 3\}$  then  $V \leftarrow V \cup \{\text{unknown}\}$ 
18:   if  $V == \{\text{buggy}\}$  then
19:     return buggy
20:   else
21:     return unknown

```

Algorithm 1 shows how we run DSI in our preliminary study; it extends DSI to also use assertions and to use multiple test runs. Algorithm 1 takes a program P , a set of tests T , and a map from each test to the specs that were mined from that test (T_{map}). The outputs are sets of likely buggy (S^f) and likely true (S^t) mined specs. Note that in Algorithm 1 we run tests across all granularity levels: test methods, test classes, and the whole test suite. The reason is that it is not clear a priori what granularity level one should use for mining and validation. But, we only use tests that mine a spec to validate that spec.

The entry point to Algorithm 1 is `runDSI` (lines 2–8). For all tests at the test method, test class, and test suite granularity levels, `runDSI` calls `validate` (lines 10–21) to check whether each spec that is mined from that test is likely buggy or likely true. For each mined spec in its input, the `validate` procedure (lines 10–21) works in three phases: *preamble* (lines 13–16), *experiment selection* (line 17), and *experiment execution* (line 20).

Preamble. Sanity checks in DSI ensure that the trace obtained from running a test is consistent with the spec, s , being validated. If a sanity check fails, DSI terminates and the outcome on s is unknown—`validate` internally tracks unknown specs in a set, S^u , that can be exposed for debugging. The sanity checks run each test thrice (lines 13–15) and compare the outcomes (lines 16, 31–36). The first run (line 16) executes P

on test and records whether the program crashes. The second and third runs (line 14 and line 15) execute a version of P that is instrumented to collect a new execution trace for the test; both runs record the test result (`pass`, `fail`, or `crash`) and the trace. Procedure `sanityCheck` (lines 31–36) takes these, and returns `false` if a sanity check failed and `true` otherwise.

The `sanityCheck` procedure returns `false` if (1) the first run does not `pass`, so DSI cannot proceed (line 32); (2) the first run `passes` but the second run does not, so instrumentation induced a failure (line 33); (3) the test outcomes in the second and third runs differ, so the runs are nondeterministic (line 34); and (4) the traces from the second and third runs do not satisfy s , so s may be buggy (line 35). The reasons for sanity check failures in our evaluation are concurrency-related issues, state pollution among tests [2], and very few instrumentation failures.

Experiment Selection. DSI creates *experiments* that mutate P to violate the spec by delaying $a()$ and invoking it after $b()$ (line 17). We use DSI’s experiment selection procedure, which aims to perform delays on P ’s execution that are minimally disruptive, i.e., delays should have minimal side effects (Algorithm 1 in [1]). It takes a trace and a spec, and it produces a set of experiments. Each experiment is a pair $\langle \text{idx}, \text{len} \rangle$, where delaying the invocation of $a()$ at idx by len steps in the trace is minimally disruptive. We highlight three important aspects of experiment selection:

(1) *Return values.* Method $a()$ may return a value which intervening code between $a()$ and $b()$, inclusive, may rely on. To reduce crashes due to missing return values of the delayed method, DSI fills in the nearest value of a variable of the same type. If no such variable exists, then DSI uses a default value. Replacing return values needs improvement in DSI (Section V-B).

(2) *Locks.* Experiments may require delaying $a()$ and invoking $b()$.
Stage 0: After delaying $a()$, before calling $b()$ —cannot run experiment and s cannot be validated; add s to S^u (line 28).

ing it at a later location, after needed locks were released. DSI records such locks and tries to reacquire them before invoking $a()$. Runs that deadlock are terminated after a timeout.

(3) *Multiple call sites.* Traces can have multiple pairs of $a()$ and $b()$ with different call sites; unique call sites are called *usage scenarios*. DSI creates an experiment per usage scenario.

Experiment Execution. DSI uses the outcome of running experiments to classify a spec (lines 20, 23–29). DSI runs experiments like so: at location idx , it captures the calling context of $a()$ in a *thunk* (i.e., a function object) and attempts to force the execution of the thunk at the program point that is offset at len from idx . If the thunk runs successfully and the test t does not `fail` or `crash`, then spec s is likely buggy—it may not be necessary for P ’s correctness. So, DSI puts s in S^f (line 27). Otherwise, DSI uses the stage of t `fail` or `crash` to classify s :

Stage 1: While running $b()$, after delaying $a()$ —method $a()$ likely establishes a precondition for $b()$.

Stage 2: After $b()$, but before $a()$ is called— $b()$ puts the program in a state in which $a()$ cannot be run.

Stage 3: After $b()$ and $a()$ run in that order—violating s breaks P .

Stages 1, 2, or 3 mean that s is a likely true; put it S^t (line 29).

Algorithm 2 shows how we use DSI in TEMARI. The `sanityCheck` procedure is the same as in Algorithm 1, but there is no `validate` procedure involved. Different from Algorithm 1, Algorithm 2 outputs whether the spec is buggy or not. So, if the spec is buggy in all DSI experiments, then `runDSI` will output `buggy`; otherwise it outputs `unknown`. That is, for TEMARI, knowledge of whether a constraint is likely true is not useful for finding bugs in the spec.

REFERENCES

- [1] M. Gabel and Z. Su, “Testing mined specifications,” in *FSE*, 2012, pp. 1–11.
- [2] A. Gyori, A. Shi, F. Hariri, and D. Marinov, “Reliable testing: Detecting state-polluting tests to prevent test dependency,” in *ISSTA*, 2015, pp. 223–233.