

Benchmarking of Java Verification Tools at the Software Verification Competition (SV-COMP)



The University of Manchester

Lucas Cordeiro
Daniel Kroening
Peter Schrammel



What is SV-COMP?

<https://sv-comp.sosy-lab.org>

Annual comparative evaluation of fully automatic software verifiers

- Reflect state of the art w.r.t. effectiveness and efficiency
- Promote reproducibility and validity of experimental results
- Increase the visibility and credits for tool developers
- Establish set of benchmarks for software verification community

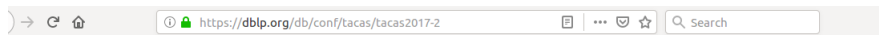
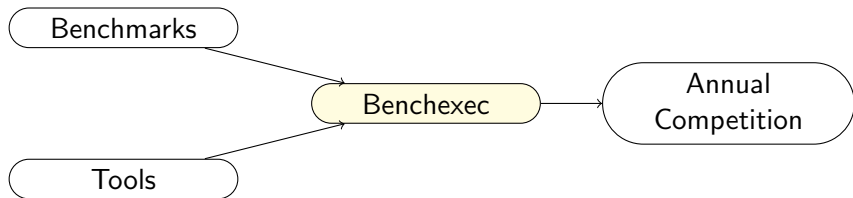
Started in 2012 for C programs, 8th edition in 2019:

- More than 10000 benchmarks
- More than 30 participants
- **NEW:** Java track

Goals of this talk

- Understand how SV-COMP works
- Know about the Java track at SV-COMP 2019
- Be able to use it for your own research
- Be able to contribute tools and benchmarks

How does SV-COMP work?



SV COMP

-     Dirk Beyer:
Software Verification with Validation of Results - (Report on SV-COMP 2017). 331-349
-     Jera Hensel, Frank Emrich, Florian Frohn, Thomas Ströder, Jürgen Giesl:
AProVE: Proving and Disproving Termination of Memory-Manipulating C Programs - (Competition Contribution). 350-354
-     Pavel Andrianov, Karlheinz Friedberger, Mikhail U. Mandrykin, Vadim S. Mutilin, Anton Volkov:
CPA-BAM-BnB: Block-Abstraction Memoization and Region-Based Memory Models for Predicate Abstractions - (Competition Contribution). 355-359
-     Williane Rocha, Herbert Rocha, Hussama Ismail, Lucas C. Cordeiro, Bernd Fischer:
DepthK: A k-Induction Verifier Based on Invariant Inference for C Programs - (Competition Contribution). 360-364
-     Lukás Holík, Martin Hruska, Ondrej Lengál, Adam Rogalewicz, Jirí Simáček, Tomáš Vojnar:
Forester: From Heap Shapes to Automata Predicates - (Competition Contribution). 365-369

Benchmarks

<https://github.com/sosy-lab/sv-benchmarks>

For each **verification task** (aka benchmark)

- Source files (open source license)
- Descriptor (.yml file)
 - File name is name of the benchmark
 - Reference to one or more **properties** (.prp files)
 - unreachable-call:
`CHECK(init(main()), LTL(G !call(__VERIFIER_error())))`
 - termination:
`CHECK(init(main()), LTL(F end))`
 - no-overflow, valid-memsafety, ...
 - Expected **answer**: true, false(property), unknown

Benchmarks

<https://github.com/sosy-lab/sv-benchmarks>

For each **verification task** (aka benchmark)

- Source files (open source license)
- Descriptor (.yml file)
 - File name is name of the benchmark
 - Reference to one or more **properties** (.prp files)
 - unreachable-call:
`CHECK(init(main()), LTL(G !call(__VERIFIER_error())))`
 - termination:
`CHECK(init(main()), LTL(F end))`
 - no-overflow, valid-memsafety, ...
 - Expected **answer**: **true**, **false(property)**, **unknown**

Categories defined as subsets (.set files)

- ReachSafety, ConcurrencySafety, MemorySafety, NoOverflows, Termination, ...
- There are sub-categories (loops, arrays, heap, ...).

Benchmark execution software

(Beyer et al SPIN'15)

- Implemented in Python 3
- Resource-limited execution (8 cores, 15GB, 900s CPU time)
- Interfaces to **competition candidates** (aka tools) via **tool-info** modules
 - Name, version
 - Build command line
 - Parse answer (**true**, **false**, **unknown**) from tool output

Benchmark execution software

(Beyer et al SPIN'15)

- Implemented in Python 3
- Resource-limited execution (8 cores, 15GB, 900s CPU time)
- Interfaces to **competition candidates** (aka tools) via **tool-info** modules
 - Name, version
 - Build command line
 - Parse answer (**true**, **false**, **unknown**) from tool output
- Table-generator to generate HTML table of results

Benchmark execution software

(Beyer et al SPIN'15)

- Implemented in Python 3
- Resource-limited execution (8 cores, 15GB, 900s CPU time)
- Interfaces to **competition candidates** (aka tools) via **tool-info** modules
 - Name, version
 - Build command line
 - Parse answer (**true**, **false**, **unknown**) from tool output
- Table-generator to generate HTML table of results
- SV-COMP 2019 runs on Ubuntu 18.04

Benchmark execution software

(Beyer et al SPIN'15)

- Implemented in Python 3
- Resource-limited execution (8 cores, 15GB, 900s CPU time)
- Interfaces to **competition candidates** (aka tools) via **tool-info** modules
 - Name, version
 - Build command line
 - Parse answer (**true**, **false**, **unknown**) from tool output
- Table-generator to generate HTML table of results
- SV-COMP 2019 runs on Ubuntu 18.04

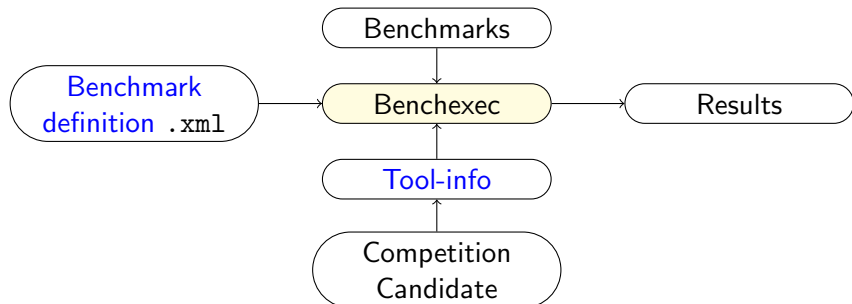
E.g. run CBMC on sub-category ReachSafety-BitVectors:
`bin/benchexec cbmc.xml -t ReachSafety-BitVectors`

Benchmark definition

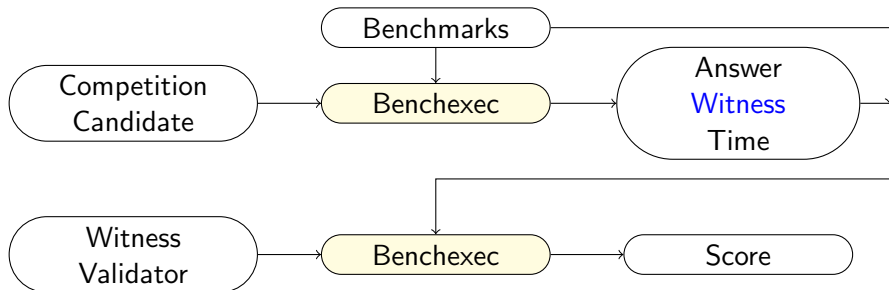
<https://github.com/sosy-lab/sv-comp>

Benchmark definition (tool.xml):

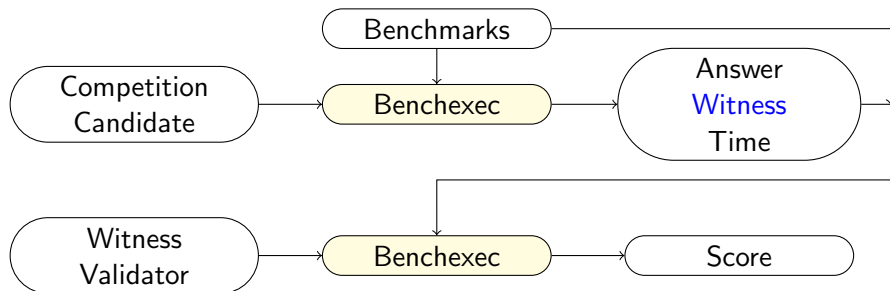
- Reference to [tool-info](#) module
- Resource limits
- Global options for tool
- Definition of (sub-)categories:
 - Reference to [category](#) .set files
 - Reference to [property](#) .prp file



Annual Competition



Annual Competition



Witness Validation

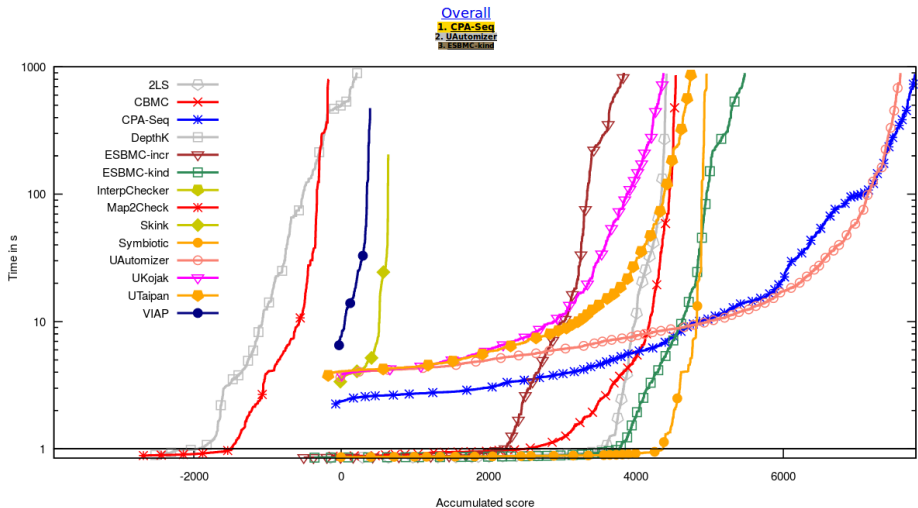
(Beyer et al FSE'15, FSE'16)

- Correctness and refutation witnesses (`.graphml`)
- Validated by witness validation tools

Scores

Points	Reported result	Description
0	UNKNOWN	Failure to compute verification result, out of resources, program crash.
+1	FALSE correct	The error in the program was found and a violation witness was confirmed.
-16	FALSE incorrect	An error is reported for a program that fulfills the specification (false alarm, incomplete analysis).
+2	TRUE correct	The program was analyzed to be free of errors and a correctness witness was confirmed.
+1	TRUE correct, witness unconfirmed	The program was analyzed to be free of errors but the correctness witness was not confirmed.
-32	TRUE incorrect	The program had an error but the competition candidate did not find it (missed bug, unsound analysis).

Annual Competition



<https://sv-comp.sosy-lab.org/2018/results/results-verified/>

Annual Competition

Tool	JPF 32				JayHorn 5.1				JBMC 5.8			
Limits	memlimit: 180 s, memlimit: 15000 MB, CPU core limit: 8											
Host	localhost											
OS	Linux 4.4.0-116-generic x86_64											
System	CPU: Intel Core i7-6700HQ CPU @ 2.60GHz with 8 cores, frequency: 3500 MHz, Turbo Boost enabled; RAM: 16014056 kB											
Date of execution	2018-03-10 14:09:57 GMT				2018-03-10 11:39:56 GMT				2018-03-10 16:54:00 GMT			
Run set	jpf.sv-comp18				jayhorn.sv-comp18				jbmc.sv-comp18			
./jpf-sv-benchmarks/jpf/	status	cpulime	walltime	memUsage	status	cpulime	walltime	memUsage	status	cpulime	walltime	memUsage
lmc-regression/ArithmeticException1_false-assert.jar	false(reach)	1.2524	0.4924	8929488	true	3.7484	1.3524	28377688	false(reach)	0.2874	0.2244	17719298
lmc-regression/ArithmeticException5_true-assert.jar	false(reach)	1.2888	0.4798	85815296	true	3.5644	1.1694	184410112	true	0.2124	0.2264	12799888
lmc-regression/ArithmeticException6_false-assert.jar	true	1.1724	0.4524	82345984	true	3.4484	1.1644	174071088	false(reach)	0.2514	0.2644	17477532
lmc-regression/ArrayIndexOutOfBoundsException1_false-assert.jar	true	1.1644	0.4754	83744528	false(reach)	6.5164	1.9344	228288448	false(reach)	0.2564	0.2688	16728864
lmc-regression/ArrayIndexOutOfBoundsException2_false-assert.jar	true	1.1484	0.4434	81854464	false(reach)	6.7484	2.8384	238465536	false(reach)	0.2584	0.2574	18628688
lmc-regression/ArrayIndexOutOfBoundsException3_false-assert.jar	true	1.2564	0.5014	82014208	false(reach)	7.2044	2.8234	233783296	false(reach)	0.2234	0.2364	18378952
lmc-regression/BufferedReaderReadLine_false-assert.jar	true	1.2064	0.4644	85873928	false(reach)	11.4844	2.9924	264462568	false(reach)	0.3734	0.3814	29315872
lmc-regression/CharSequenceBug_false-assert.jar	true	1.1364	0.4444	79372288	false(reach)	7.3364	2.1914	233285828	false(reach)	1.4874	1.4224	26898248
lmc-regression/CharSequenceToString_true-assert.jar	true	1.1364	0.5014	80220168	false(reach)	7.4124	2.1864	235995024	true	1.8474	1.8594	36552784
lmc-regression/ClassCastException1_false-assert.jar	false(reach)	1.1724	0.4624	86487848	false(reach)	6.1244	1.8984	238752256	false(reach)	0.3054	0.3194	18665472
lmc-regression/ClassCastException2_true-assert.jar	true	1.2164	0.5964	77561856	true	6.5124	1.9074	233218048	true	0.2334	0.2534	13275136
lmc-regression/ClassCastException3_false-assert.jar	false(reach)	1.1484	0.4414	84291584	false(reach)	6.8524	2.8884	233496576	false(reach)	0.2354	0.2424	17874944
lmc-regression/Class_method1_true-assert.jar	true	1.1684	0.4574	80973824	true	5.1484	1.6994	236470272	true	0.2364	0.2564	13535664
lmc-regression/Inheritance1_true-assert.jar	true	1.1884	0.4824	83697864	true	7.2484	2.8464	238694912	true	0.2684	0.2764	13492224
lmc-regression/NegativeArraySizeException1_false-assert.jar	false(reach)	1.2364	0.4924	87408448	true	3.7284	1.2384	185192088	false(reach)	0.2134	0.2254	16687104
lmc-regression/NegativeArraySizeException2_false-assert.jar	false(reach)	1.1764	0.4724	83542016	true	3.9884	1.2454	187536688	false(reach)	0.2414	0.2474	18546688
lmc-regression/NullPointerException1_false-assert.jar	false(reach)	1.1524	0.4484	84785288	unknown	1.8644	0.6884	136836352	false(reach)	0.2434	0.2454	16643568
lmc-regression/NullPointerException2_false-assert.jar	false(reach)	1.2124	0.4664	85893128	unknown	1.7484	0.6764	134981768	false(reach)	0.2444	0.2514	16654336
lmc-regression/NullPointerException3_false-assert.jar	false(reach)	1.3684	0.5324	88222752	unknown	1.6684	0.6184	133578752	false(reach)	0.2434	0.2554	16982016
lmc-regression/NullPointerException4_false-assert.jar	false(reach)	1.2724	0.5354	77479936	unknown	1.6364	0.6264	128978752	false(reach)	0.2484	0.2484	18771968
lmc-regression/RegexMatches01_true-assert.jar	true	1.3084	0.5374	80244736	false(reach)	8.3444	2.3934	246411264	true	1.2884	1.3014	39448968
lmc-regression/RegexMatches02_false-assert.jar	true	1.1924	0.4814	85598288	false(reach)	9.8244	2.6214	254148416	false(reach)	1.8334	1.8444	64978844
lmc-regression/RegexSubstitution01_true-assert.jar	true	1.2524	0.5034	83034112	false(reach)	14.2684	3.5354	394878976	true	72.7744	72.8584	757209440
lmc-regression/RegexSubstitution02_false-assert.jar	true	1.1164	0.4424	81616896	false(reach)	17.1364	4.2844	398799368	true	48.9734	49.0114	526553888
lmc-regression/RegexSubstitution03_true-assert.jar	true	1.2324	0.4654	86216784	false(reach)	12.6884	3.1584	338695952	true	1.7544	1.7664	46884992
lmc-regression/StaticCharMethods01_true-assert.jar	true	1.1764	0.5024	82983048	false(reach)	6.7444	2.8414	228970496	true	0.2614	0.2764	14378968
lmc-regression/StaticCharMethods02_false-assert.jar	true	1.1964	0.4844	83034112	false(reach)	8.5244	2.4764	237215744	false(reach)	0.8294	0.8544	50972748
lmc-regression/StaticCharMethods03_false-assert.jar	true	1.1924	0.4864	82286788	false(reach)	10.1884	2.8884	243417888	false(reach)	0.8474	0.8584	50866784
lmc-regression/StaticCharMethods04_false-assert.jar	true	1.1964	0.4874	84496384	false(reach)	8.6644	2.5294	231288648	false(reach)	0.8214	0.8294	51421992

Competition Timeline

- September: Contribution of benchmarks
- October: Tool registration and qualification
- November: Tool submission
- December: Announcement of winners
- January: Tool paper submission
- April: SV-COMP session at ETAPS

Benchmarks and Properties

368 benchmarks (40LOC on average, 250LOC max)

- jayhorn-recursive, jbmc-regression, jpf-regression, MinePump

Java 1.8

Calls to Java standard library (`java.*`, `javax.*`) allowed;
sources of other dependencies must be part of the benchmark.

1 category for violation of asserts ("ReachSafety")

Property: `CHECK(init(Main.main), LTL(G assert))`

Benchmarks and Properties

`java/jbmc-regression/StringStartEnd02/Main.java`

`java/jbmc-regression/StringStartEnd02.yml`

```
format_version: "0.1"
input_files:
  - ../common/
  - StringStartEnd02/
properties:
  - property_file: ../properties/assert.prp
    expected_verdict: false
```

```
import org.sosy_lab.sv_benchmarks.Verifier;

public class Main {
    public static void main(String[] args) {
        String[] strings = new String[4];
        strings[0] = Verifier.nondetString();
        strings[1] = Verifier.nondetString();
        strings[2] = Verifier.nondetString();
        strings[3] = Verifier.nondetString();

        int i = 0;
        for (String string : strings) {
            if (string.startsWith("te"))
                ++i;
        }
        assert i == 1;
    }
}
```

If a tool requires class files as input it is responsible for compiling the benchmark.

A benchmark must be compilable by passing all `.java` files within the directories listed in `input_files` to `javac`.

Rules for Nondeterminism

Only source of nondeterminism:
return values of methods defined in
`org.sosy_lab.sv_benchmarks.`
`Verifier` class.

Must not be used:

- Arguments of `main`
- Library methods that make system calls

```
package org.sosy_lab.sv_benchmarks;

import java.util.Random;

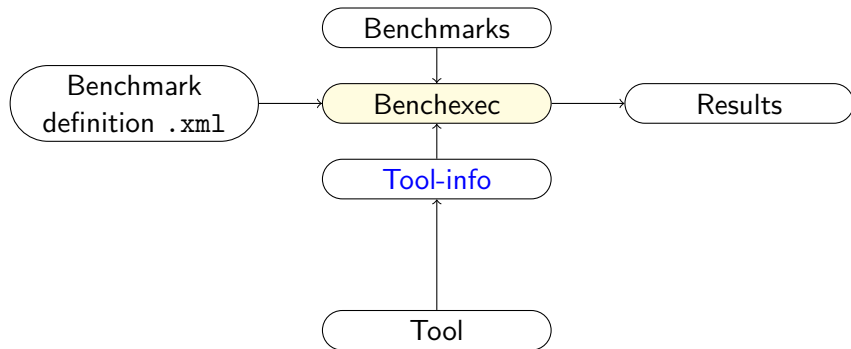
public final class Verifier
{
    public static void assume(boolean condition)
    {
        if(!condition) {
            Runtime.getRuntime().halt(1);
        }
    }

    public static boolean nondetBoolean()
    {
        return new Random().nextBoolean();
    }

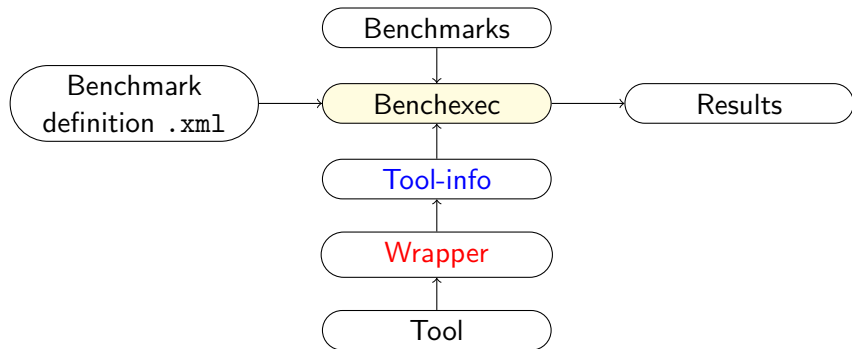
    public static byte nondetByte()
    {
        return (byte)(new Random().nextInt());
    }

    . . .
}
```

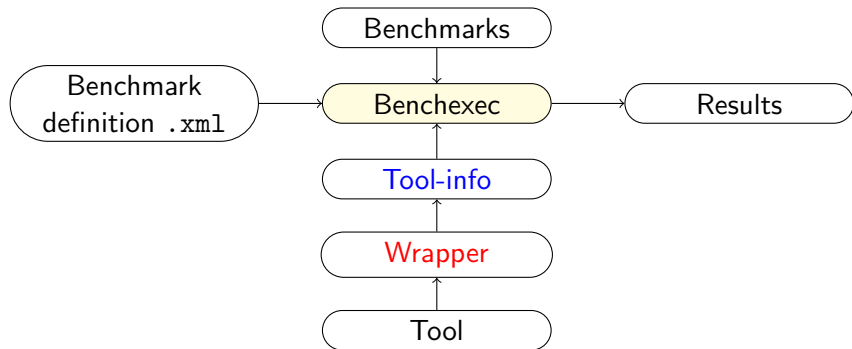
Tool-info and Wrapper Scripts



Tool-info and Wrapper Scripts



Tool-info and Wrapper Scripts



E.g. command line produced for JPF:

```
./jpf-sv-comp  
--graphml-witness witness.graphml  
--propertyfile ../sv-benchmarks/java/properties/assert.prp  
../sv-benchmarks/java/common/org/sosy_lab/sv_benchmarks/Verifier.java  
../sv-benchmarks/java/jbmc-regression/StringStartEnd02/Main.java
```

How can I use it?

Re-use existing benchmarking infrastructure

- Stop writing benchmarking scripts
- Use it for running your tests

'Standard' benchmark set

- Use it for running comparisons
- Contribute your benchmarks

Compare with the 'best' configuration of a tool

- Take the most recent competition candidate:
Download from
<https://sv-comp.sosy-lab.org/2018/systems.php>

How can I use it?

Reproduce the competition results:

- Download:

```
git clone https://github.com/sosy-lab/sv-benchmarks
git clone https://github.com/sosy-lab/benchexec
git clone https://github.com/sosy-lab/sv-comp
git clone https://gitlab.com/sosy-lab/sv-comp/archives-2019
```

- Run:

```
cd benchexec
for tool in jpf spf jayhorn jbmc
do
    unzip ../archives-2019/2019/$tool.zip; mv $tool/* .
    bin/benchexec ../sv-comp/benchmark-defs/$tool.xml
done
bin/table-generator results/*.xml.bz2
```

Currently (5 Nov 2018) only runs with benchexec's yam1 branch.

How can I use it?

Run JPF on your own benchmarks:

- Add descriptor yml file for each benchmark
- List descriptor files in `MyBenchmarks.set`
- Add `MyBenchmarks.set` to benchmark definition `jpf.xml`:

```
<tasks name="MyBench">
  <includesfile>
    ../sv-benchmarks/java/MyBenchmarks.set</includesfile>
  <propertyfile>
    ../sv-benchmarks/java/properties/assert.prp</propertyfile>
</tasks>
```

- Run with `-t MyBench`

```
bin/benchexec ../sv-comp/benchmark-defs/jpf.xml -t MyBench
bin/table-generator results/*.xml.bz2
```


Outlook

What is needed for 2020:

- More benchmarks:
Fork <https://github.com/sosy-lab/sv-benchmarks> and
create PR with your benchmarks

Outlook

What is needed for 2020:

- More benchmarks:
Fork <https://github.com/sosy-lab/sv-benchmarks> and
create PR with your benchmarks
- Witness validators

Witness Validation for Java

Refutation witnesses (for ReachSafety property):

- Witness contains counterexample trace annotated with evaluated assignments and conditionals
→ Check whether counterexample trace is feasible and violates the property
- Proposed implementation: generate harness, compile and execute

Correctness witnesses (for ReachSafety property):

- Witness contains dynamic CFG annotated with invariants
→ Check whether invariants hold and imply properties
- ???

Outlook

What is needed for 2020:

- More benchmarks:
Fork <https://github.com/sosy-lab/sv-benchmarks> and
create PR with your benchmarks
- Witness validators

Outlook

What is needed for 2020:

- More benchmarks:
Fork <https://github.com/sosy-lab/sv-benchmarks> and
[create PR with your benchmarks](#)
- Witness validators
- Encourage participation of more tools and tool variants

Subscribe to sv-comp@googlegroups.com

Up-to-date version of paper: <http://arxiv.org/abs/1809.03739>

Outlook

What is needed for 2020:

- More benchmarks:
Fork <https://github.com/sosy-lab/sv-benchmarks> and
[create PR with your benchmarks](#)
- Witness validators
- Encourage participation of more tools and tool variants

Subscribe to sv-comp@googlegroups.com

Up-to-date version of paper: <http://arxiv.org/abs/1809.03739>

www.diffblue.com

- Jobs in program analysis,
verification and machine learning!

