## 8th Competition on Software Verification

#### Dirk Beyer (Competition Chair)

Supported by:







# The content of this presentation is available in the SV-COMP 2019 report: https://doi.org/10.1007/978-3-030-17502-3\_9

- 1. Community suffers from unreproducible results  $\rightarrow$  Establish set of benchmarks
- 2. Publicity for tools that are available  $\rightarrow$  Provide state-of-the-art overview
- 3. Support the development of verification tools  $\rightarrow$  Give credits and visibility to developers
- 4. Establish standards
  - $\rightarrow$  Specification language, Witnesses, Benchmark definitions, Validators

#### Session 1:

- Competition Report, by organizer
- System Presentations, 4 min by each team
- Short discussion

#### Session 2:

 Open Jury Meeting, Community Discussion, moderated by organizer Three Steps - Three Deadlines:

- Benchmark submission deadline
- System submission
- Notification of results (approved by teams)

Input:

- $\blacktriangleright$  C program  $\rightarrow$  GNU/ANSI C standard
- Property
  - $\rightarrow$  Reachability of error label, of overflows
  - $\rightarrow$  Memory safety (inv-deref, inv-free, memleak)
  - $\rightarrow$  Termination

Output:

- ► TRUE + Witness
- ► FALSE + Witness
- UNKNOWN

(property holds)

(property does not hold)

(failed to compute result)

Machines (1000 \$ consumer machines):

- ► CPU: 3.4 GHz 64-bit Quad-Core CPU
- RAM: 33 GB
- ▶ OS: GNU/Linux (Ubuntu 18.04)

Resource limits:

- ▶ 15 GB memory
- ▶ 15 min CPU time (consumed 461 days)

Volume: 178 674 ver. runs, 517 175 val. runs Total: 5 880 071 runs using 15 years and 182 days of CPU time Common principles: Ranking measure should be

- easy to understand
- reproducible
- computable in isolation for one tool

SV-COMP:

- Ranking measure is the quality of verification work
- Expressed by a community-agreed score
- Tie-breaker is CPU time

## Scoring Schema (2019)

(from 2020 onwards: only confirmed results count)

Reported result	Points	Description
UNKNOWN	0	Failure, out of ressources
FALSE correct	+1	Error found and confirmed
FALSE incorrect	-16	False alarm (imprecise analysis)
TRUE correct	+2	Proof found and confirmed
TRUE unconfirmed	+1	Proof found but unconfirmed
TRUE incorrect	-32	Missed bug (unsound analysis)

Jury:

- Team: one member of each participating candidate
- Term: one year (until next participants are determined)

Systems:

- All systems are available in open GitLab repo
- Configurations and Setup in GitHub repository
   Integrity and reproducibility guaranteed

Qualification:

- 31 Qualified (out of 31 Submitted)
   1 verifier disqualified from several categories (rule viol.)
- One person can participate with different tools
- One tool can participate with several configurations (frameworks, no tool-name inflation)

Benchmark quality:

Community effort, documented on GitHub

Role of organizer:

Just service: Advice, Technical Help, Executing Runs

Everybody can submit benchmarks (conditions apply)

- Eight categories when closed (scores normalized):
  - Reachability: 3831 tasks
  - Memory Safety: 434 tasks
  - Concurrency: 1082 tasks
  - NoOverflows: 359 tasks
  - Termination: 2007 tasks
  - Software Systems: 2809 tasks
  - Overall: 10522 tasks
  - Java: 368 tasks

SV-Benchmarks: https://github.com/sosy-lab/sv-benchmarks

- SV-COMP Setup: https://github.com/sosy-lab/sv-comp
- Resource Measurement and Process Control: https://github.com/sosy-lab/benchexec
- Archives:

https://gitlab.com/sosy-lab/sv-comp/archives-2019

#### • Witnesses:

https://sv-comp.sosy-lab.org/2019/results/
results-verified

	Selec	t Colu	mns	Filt	er Rows	C	uantile P	lot	Scat	ter Plot	Shrink Header							Gener	ated with	BenchExec
												2LS 0.5.0								
	timelimit: 900 s, memlimit: 15000 MB, CPU core limit: 8																			
																apollon*				
																Linux 4.4.0-57-generic				
											CPU: Intel X	eon E3-	1230 v	5@3.4	) GHz, ci	ores: 8, frequency: 3.8 GHz, 1	Turbo	Boost	disable	1; RAM: 335
										2	017-01-10 17:21:21 CET [[ 20	17-01-	14 18:0	10:17 CE	T ]] [[ 20	17-01-14 20:02:31 CET ]] [[ 2	2017-0	)1-14 1	8:18:08	CET ]] [[ 20
															sv-con	p17.ReachSafety-ControlFic	w			
	I mental and a second s																			
	verifier status	score	witness	inspect witness	cpu (s)	wall (s)	energy (J)	mem (MB)	blkio-w (MB)	blkio-r (MB)	validator cpachecker violation t<90s status	cpu (s)	wall (s)	energy (J)	mem (MB)	validator uautomizer violation t<90s status	cpu (s)	wall (s)	energy (J)	mem (MB) c
ination.cil.c	false(unreach-call)	1	wit	inspect	1.3	1.3	13	370	.0041	0	false(unreach-call)	8.2	4.4	120	320	false(unreach-call)	17	9.1	320	520
ation.cil.c	false(unreach-call)	1	wit	inspect	.35	.34	3.3	60	.0041	0	false(unreach-call)	8.1	4.3	170	310	false(unreach-call)	13	6.6	240	450
mation.cil.c	false(unreach-call)	1	wit	inspect	.55	.53	4.9	120	.0041	0	false(unreach-call)	8.3	4.4	100	330	false(unreach-call)	12	6.6	210	500
nation.cil.c	false(unreach-call)	1	wit	inspect	.29	.28	2.4	39	.0041	0	false(unreach-call)	7.8	4.2	80	380	false(unreach-call)	13	6.7	180	410
ation.cil.c	true	2	wit	inspect	1.4	1.4	14	410	.0041	12	100 B	-	-	-	-		-	-	-	-
ination.cil.c	true	2	wit	inspect	.53	.53	4.9	110	.0041	0		-	-	-	-		-	-	-	-
tion.cil.c	true	2	wit	inspect	.36	.35	3.3	67	.0041	0		-	-	-	-		-	-	-	-
tion.cil.c	true	2	wit	inspect	.59	.58	5.6	130	.0041	0		-	-	-	-		-	-	-	-
mation.cil.c	true	2	wit	inspect	.19	.19	1.4	26	.0041	0	-	-	-	-	-		-	-	-	-
ation.cil.c	true	2	wit	inspect	.29	.29	2.3	45	.0041	0		-	-	-	-	-	-	-	-	-
	talse(unreach-call)	1	wit	inspect	21	21	200	200	.0041	0	taise(unreach-call)	7.4	3.9	170	310	taise(unreach-call)	12	6.8	210	460
	false(unreach-call)	1	wit	inspect	23	23	200	200	.0041	0	false(unreach-call)	7.4	3.9	130	310	false(unreach-call)	14	7.1	200	480



'Value' of result is defined by Scoring Schema

## Impact / Achievements

- ► Large benchmark set of verification tasks → established and used in many papers for experimental evaluation
- ► Good overview over state-of-the art → covers model checking and program analysis
- Participants have an archived track record of their achievements
- Infrastructure and technology for controlling the benchmark runs (cf. StarExec)

[Competition Report and System Descriptions are archived in Proceedings TACAS '19] https://doi.org/10.1007/978-3-030-17502-3\_9

## Alternative Rankings — Definitions

#### Correct Verifiers — Low Failure Rate:

number of incorrect results total score

with unit E/sp.

Green Verifiers — Low Energy Consumption:

total CPU energy total score

with the unit J/sp.

 New Verifiers — High Quality: quality in score points as rank measure.

## Alternative Rankings — Results

Table 9: Alternative rankings; quality is given in score points (sp), CPU time in hours (h), energy in kilojoule (kJ), wrong results in errors (E), rank measures in errors per score point (E/sp), joule per score point (J/sp), and score points (sp)

$\mathbf{Rank}$	Verifier	Quality	$\mathbf{CPU}$	CPU	Solved	Wrong	Rank
			Time	Energy	Tasks	Results	Measure
		(sp)	(h)	(kJ)		(E)	
Correct	$t \ Verifiers$						(E/sp)
1	CPA-Seq	9329	120	4300	2811	0	.0000
2	Symbiotic	6129	9.7	390	2519	0	.0000
3	PeSCo	8466	120	3900	2431	9	.0011
worst							.3836
Green	Verifiers						$(\mathrm{J/sp})$
1	Symbiotic	6129	9.7	390	299	0	64
2	CBMC	4341	11	380	296	14	88
3	DIVINE-EXPLICIT	1547	4.4	180	256	10	120
worst							4 200
New V	<i>Verifiers</i>						(sp)
1	PeSCo	8466	120	3900	1026	9	8466
2	CBMC-Path	1587	8.9	380	1006	69	1587

18/34

Result		Т	RUE		False							
	Total	Confi	rmed	Unconf.	Total	Confi	rmed	Unconf.				
CPA-Seq	4417	3968	90%	449	2859	2686	94%	173				
PeSCo	4176	3814	91%	362	2823	2652	94%	171				
UAutomizer	4244	4199	99%	45	1523	1255	82%	268				
Symbiotic	2430	2381	98%	49	1451	1214	84%	237				
CBMC	1813	1702	94%	111	1975	1248	63%	727				
UTAIPAN	3015	2936	97%	79	915	653	71%	262				
2LS	2072	2045	99%	27	1419	945	67%	474				
ESBMC-KIND	3679	3556	97%	123	2141	1753	82%	388				
UKojak	2070	2038	98%	32	553	548	99%	5				
CBMC-Path	1206	1162	96%	44	897	670	75%	727				
DIVINE-explicit	693	673	97%	20	768	353	46%	415				
DIVINE-SMT	645	626	97%	19	943	601	64%	342				
DepthK	612	602	98%	10	1938	1370	71%	568				

Table 10: Confirmation rate of verification witnesses in SV-COMP 2019

#### Number of Participants



Fig. 7: Number of participating teams for each year

- Task definition
- License
- More programs
- LTL properties
- Eliminate pre-processing
- Undefined behavior of C programs
- Witnesses in all categories
- Tests as Witnesses

```
format version: '1.0'
1
2
   # old file name: floppy_true-unreach-call_true-valid-memsafety.i.cil.c
3
   input files: 'floppy.i.cil-3.c'
4
5
   properties:
6
    - property_file: ../properties/unreach-call.prp
7
       expected_verdict: true
8
     - property_file: ../properties/valid-memsafety.prp
9
       expected verdict: true
10
```

Fig. 3: Example task definition for program floppy.i.cil-3.c

#### Practical Impact: Get Tests from Verification Tools



- ► TACAS (PC Chairs + TACAS SC, thanks!)
- ► Jury (32 people)
- Participants (177 people)
- Sponsors: Amazon Web Services and LMU Munich







### **Benchmark Definition**

```
<?xml version = "1.0"?>
<!DOCTYPE benchmark PUBLIC "+//IDN sosy-lab.org//DTD BenchExec benchmark 1.9//EN"
"http://www.sosy-lab.org/benchexec/benchmark -1.9.dtd">
<benchmark tool="cpachecker" timelimit="15 min" hardtimelimit="16 min"</pre>
 memlimit="15 GB" cpuCores="8">
<require cpuModel="Intel Xeon E3-1230 v5 @ 3.40 GHz" cpuCores="8"/>
<resultfiles >**.graphml</resultfiles >
<option name="-sycomp19"/>
<option name="-heap">10000M</option>
<option name="-benchmark"/>
<option name="-timelimit">900 s</option>
<rundefinition name="sv-comp19_prop-reachsafety">
<tasks name="ReachSafety-Arrays">
<includesfile >../sv-benchmarks/c/ReachSafety-Arrays.set</includesfile>
cpropertyfile >../sv-benchmarks/c/properties/unreach-call.prp</propertyfile >
</tasks>
<tasks name="ReachSafety-BitVectors">
<includesfile >../sv-benchmarks/c/ReachSafety-BitVectors.set</includesfile >
< propertyfile > ... / sv-benchmarks / c / properties / unreach-call . prp <math>< / propertyfile >
</tasks>
```

## LICENSE FOR RESEARCH AND EVALUATION

[...]

The licensor grants to the user of this software the following rights, irrevocably:

- 1. Redistribution of exact copies of this archive.
- 2. Execution of this software for research and evaluation purposes.
- 3. Publication of the output and measurement results obtained from executing this software.

The licensor confirms that the licenses of all components contained in this archive are compatible with the above requirements.

#### GUARANTEE OF RIGHTS FOR RESEARCH AND EVALUATION

The licensor guarantees that the licenses of all components contained in this archive grant the following perpetual, worldwide, no-charge, royalty-free, irrevocable rights to everybody:

- 1. Redistribution of exact copies of this archive.
- 2. Execution of this software for research and evaluation purposes.
- 3. Publication of the output and measurement results obtained from executing this software.

For all parts of the software for which the licensor holds the copyright, the licensor grants to the user of this software the above rights. This takes precedence over any contradicting clauses in accompanying licenses.

## Search-Space Reduction for Stepwise Testification



## **Produce Witnesses**



29 / 34

## Search-Space Reduction for Stepwise Testification





## Search-Space Reduction for Stepwise Testification



#### Produce Unit TestsFrom Witnesses



## Search-Space Reduction for Stepwise Testification



<ロト < 部 ト < 言 ト < 言 ト 言 の < で 34/34