qcomp.org

# Competing with Probabilities: Challenges and Outcomes of QComp

Arnd Hartmanns

University of Twente

*based on joint work with Carlos E. Budde, E. Moritz Hahn, Christian Hensel, Sebastian Junges, Michaela Klauck, Joachim Klein, Jan Křetínský, David Parker, Tim Quatmann, Enno Ruijters, Marcel Steinmetz, Andrea Turrini, and Zhen Zhang*

# QComp: A Quantitative Competition

"Friendly competition": no ranking

Semantic formalisms:
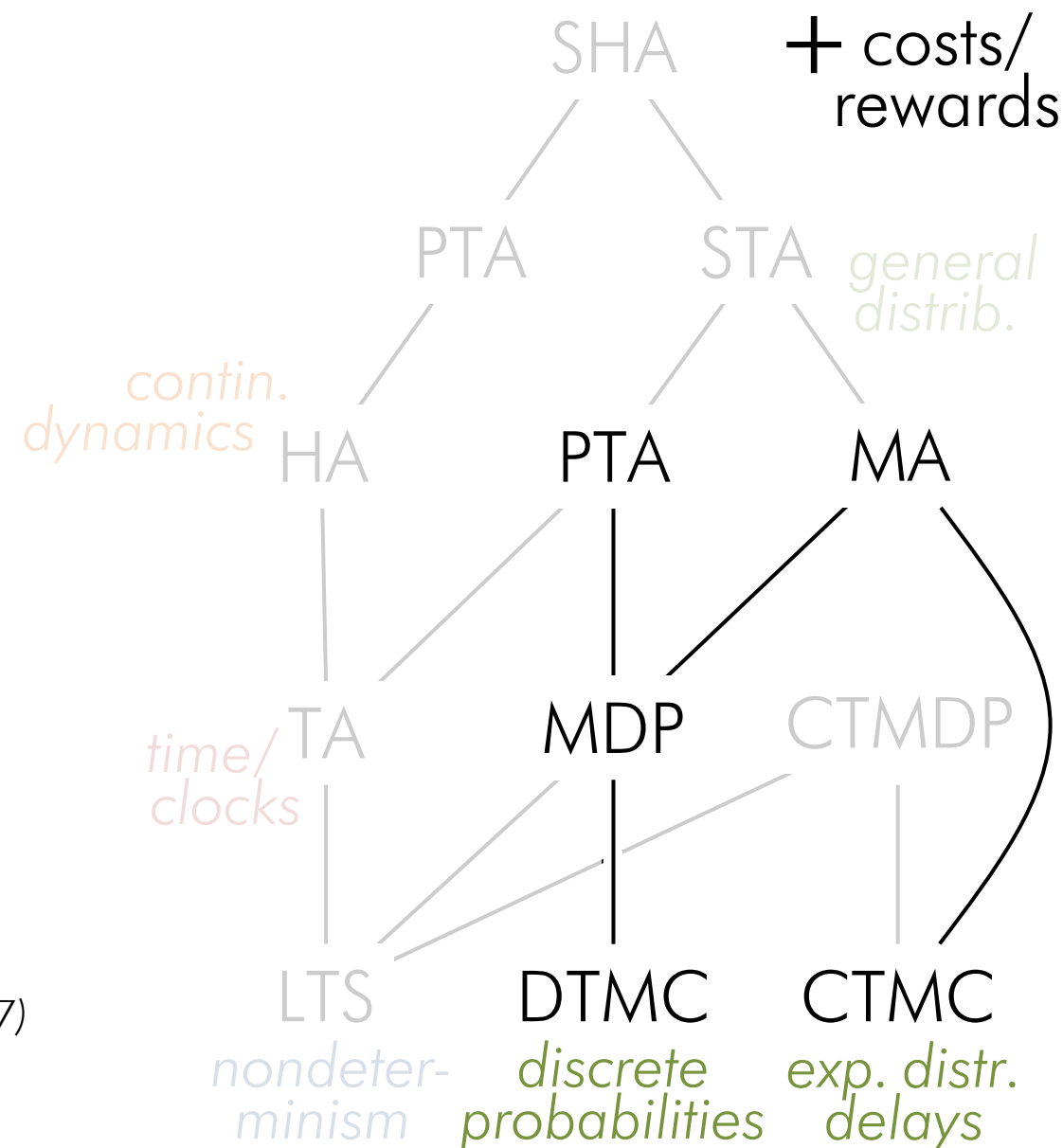DTMC, CTMC, MDP, MA, PTA

Modelling languages:
GreatSPN        *stochastic Petri nets*
PPDDL           *planning domains*
PRISM           *general, low-level*
…               *and several others*
+ JANI          *model exchange format*
                (*jani-spec.org*, *TACAS'17*)

SHA        + costs/
                rewards

PTA          STA    *general distrib.*

*contin. dynamics*    HA          PTA          MA

*time/ clocks*    TA          MDP          CTMDP

LTS          DTMC          CTMC

*nondeter-minism*    *discrete probabilities*    *exp. distr. delays*

# QComp: A Quantitative Competition

Properties to check:

reachability probability $\quad \mathbb{P}(\diamond\, G)$

expected reward $\quad\quad\quad\; \mathbb{E}(\text{cost} \to G)$

steady-state probability $\quad \mathbb{S}(G)$

} unbounded, time-, reward-bounded

Benchmarks from the Quantitative Verification Benchmark Set
*all QVBS entries must have a JANI version*

**QVBS**

*TACAS 2019*

qcomp.org/benchmarks

The Quantitative Verification Benchmark Set*

Arnd Hartmanns[1], Michaela Klauck[2], David Parker[3],
Tim Quatmann[4], and Enno Ruijters[1]

[1] University of Twente, Enschede, The Netherlands
[2] Saarland University, Saarbrücken, Germany
[3] University of Birmingham, Birmingham, United Kingdom
[4] RWTH Aachen, Aachen, Germany

QComp - Quantitative Verification ✕ QComp - Quantitative Verification ✕ +

QComp - Quantitative Verification ✕

with a [Pb ⌄] property and [zero] - [infinity] states

earch

Show all models      of type [(all) ⌄] / [GreatSPN ⌄]

| Models | | Type | Original | Params | States | Properties | Notes |
|---|---|---|---|---|---|---|---|
| **Model ▾** | **Name** | | | | | 6 (2×Pb, 2×... | (small symbolic r... |
| ☑ flexible-... | Flexible Manufacturing... | MA | GreatSPN | 2 (1/1) | 2.44 k - 2.70 M | 3 (1×P, 1×Pb,... | (small symbolic r... |
| ☑ philosop... | Dining Philosophers | CTMC | GreatSPN | 2 (1/1) | 34 - 1.77 T | 4 (2×P, 1×Pb... | (standard GSPN ... |
| | readers-... | Readers and Writers S... | MA | GreatSPN | 1 (1/0) | 1.61 k - 101 M | | |

↳ compare results

(CLOSE)

# The Competitors: Algorithms

# Probabilistic Model Checking $^{\textbf{PMC}}$

= numeric algorithm on full state space

**—** limited by state space explosion

**+** $\epsilon$-correct results: $|v - \bar{v}|/v \leq \epsilon$

**(unknown)** true value ↑   ↑ computed result

# Statistical Model Checking $^{\textbf{SMC}}$
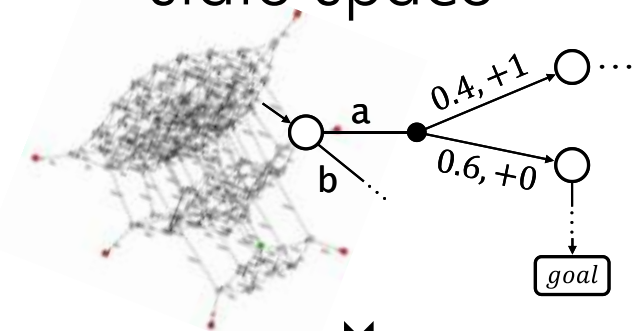
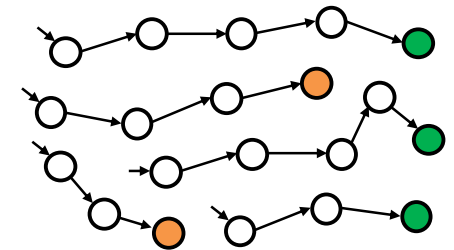**PMC:**

formal model

```
process P() {
  alt {
  :: stop {= fail = true =}
  :: send palt {
      :95: {= done = true =}
      : 5: reset; P()
} } }
```

⇓

state space



⇓

precise results

$\mathbb{P}_{\min}(\diamond\, a) = 0.2035$

$\mathbb{P}_{\max}(\diamond\, a \wedge b) = 0.89$

$\mathbb{E}_{\min}(\#s \mid b) = 12.5$

# Probabilistic Model Checking [PMC]

= numeric algorithm on full state space

− limited by state space explosion

+ $\epsilon$-correct results: $|v - \bar{v}|/v \leq \epsilon$

(**unknown**) true value ↑ ↑ computed result

# Statistical Model Checking [SMC]

= formal Monte Carlo simulation

+ constant memory usage

− rare events, nondeterminism

PAC guarantee: $\mathbb{P}(|v - \hat{v}| > \epsilon) < 1 - \delta$

↑ estimate ↑ confidence, e.g. 95 %

**SMC:**

formal model

```
process P() {
 alt {
 :: stop {= fail = true =}
 :: send palt {
    :95: {= done = true =}
    : 5: reset; P()
} } }
```

⇩

sample runs



⇩

estimated results

$\mathbb{P}_{\min}(\diamond a) \approx 0.2$

$\mathbb{P}_{\max}(\diamond a \wedge b) \approx 0.9$

$\mathbb{E}_{\min}(\#s \mid b) \approx 12$

# Probabilistic Model Checking $^{PMC}$

= numeric algorithm on full state space

**—** limited by state space explosion

**+** $\epsilon$-correct results: $|v - \bar{v}|/v \le \epsilon$

**(unknown)** true value ↑    ↑ computed result

# Statistical Model Checking $^{SMC}$

= formal Monte Carlo simulation

**+** constant memory usage

**—** rare events, nondeterminism

PAC guarantee: $\mathbb{P}(|v - \hat{v}| > \epsilon) < 1 - \delta$

estimate ↑    ↑ confidence, e.g. 95%

**+ Hybrid Approaches**

*reinforcement*
*learning*

*deep*
*learning*

*truncation*

*partial exploration,*
*guided by simulation*

*probabilistic*
*planning*

# Challenges to Correctness

# Challenges to Correctness

## 1. Bugs in algorithms:

*the algorithm itself is incorrect*

*e.g. sound value iteration: small bug in helper method pseudocode in original paper, wrong in 1 of 79 test models*

## 2. Bugs in implementations:

*the algorithm is correct, but the implementation is not*

Acceptable?

Solutions:  verify the algorithm with a theorem prover

correct-by-construction implementations

program verification

*Isabelle*

**not specific to the quantitative setting**

# Challenges to Correctness

3. Unsound algorithms:

*often but not always deliver $\epsilon$-correct result*

$\rightarrow$ value iteration and derived algorithms with one-sided approximation of the fixpoint only



Solutions: interval iteration, optimistic value iteration, BRTDP, …

4. Floating-point implementations:

*results unpredictably affected by rounding, cancellation, …*

Solutions: exact rational arithmetic, safe rounding

does not scale        new

**specific to probabilistic model checking**

# Challenges to Correctness

5. The statistical error in SMC:

*up 5% of the results may be totally wrong, and that's okay*

Recall PAC guarantee: $\mathbb{P}(|v - \hat{v}| > \epsilon) < 1 - \boxed{\delta}$

true value     estimate     error     $\boxed{\text{confidence}}$ e.g. 95%

$\rightarrow$ unavoidable in a statistical approach,
  quantifiable (user-selectable confidence level)

*How can we deal with these challenges
in a tool competition?*

# Correct Quantitative Competitions

Option N:   *Disqualify any tool that produces just a single*
            *($\epsilon$-)incorrect result and publicly shame its authors*

**maybe that's a good idea?**

Consequences:   All SMC tools disqualified
                No unsound algorithms allowed
                Floating-point implementations out

                $\rightarrow$ only STORM and PRISM remain,
                using their limited *exact* engines

**not representative of today's quantitative verification tools**

Option SMC:

*Use statistical test on statistical tools*
*to assure confidence $\delta$ is adhered to*

**evaluation time explosion**

# The QComp 2020 Approach

# QComp 2020: Tracks

Option QC20: *Use different tracks for different guarantees*

**correct**: must match true rational value where known          $\epsilon = 0$

**floating-point correct**: must use algorithm that gives          $\epsilon = 10^{-14}$
   exact result, but may use floating-point arithmetic

$\epsilon$**-correct**: unconditionally require $|v - \bar{v}|/v \leq \epsilon$          $\epsilon = 10^{-6}$

**probably $\epsilon$-correct**: require $\mathbb{P}(|v - \hat{v}| > \epsilon) < 1 - \delta$          $\epsilon = 5 \cdot 10^{-2}$
   from algorithm, but we do not check this statistically

**often $\epsilon$-correct**: should ensure $|v - \bar{v}|/v \leq \epsilon$, but          $\epsilon = 10^{-3}$
   may sometimes be wrong (also with **10′** bound)

# QComp 2020: Tools

| | | | |
|---|---|---|---|
| ePMC | PMC | *modular tool, focus on LTL* | ISCAS |
| MCSTA | PMC | *disk-based, focus on correctness* | Twente |
| PRISM | PMC | *the original probabilistic model checker* | Birmingh. |
| STORM | PMC | *has all the algorithms and languages* | Aachen |
| DFTRES | SMC | *dynamic fault tree rare event simulator* | Twente |
| MODES | SMC | *rare events and nondeterminism* | Twente |
| MFPL | hybrid | *probabilistic planning using LRTDP* | Saarland |
| PET | hybrid | *the partial exploration tool* | Munich |
| STAMINA | hybrid | *truncation for infinite-state CTMC* | Utah |

# QComp 2020: Tools



| Tool | Galileo | GreatSPN | Jani | Modest | PGCL | PPDDL | Prism | DTMC P | DTMC Pr | DTMC E | CTMC P | CTMC Pt | CTMC E | CTMC S | MDP P | MDP Pr | MDP E | MA P | MA Pt | MA E | MA S | PTA P | PTA Pt | PTA E |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DFTRES | ✓ | | ✓ | | | | | + | | | + | ✓ | + | ✓ | | | | + | ✓ | + | ✓ | | | |
| ePMC | | | ✓ | | | | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | | ✓ | | ✓ | | | | | | | |
| mcsta | | | ✓ | ✓ | | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | + | ✓ | ✓ | ✓ | ✓ | ✓ | + | ✓ | ✓ | ✓ | ✓ |
| modes | | | ✓ | ✓ | | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | + | ✓ | ✓ | ✓ | ✓ | ✓ | + | ✓ | ✓ | ✓ | ✓ |
| MFPL | | | ✓ | ✓ | | | | | | | | | | | + | | + | | | | | | | |
| Prism | | | | | | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | | | ✓ | ✓ | ✓ |
| PET | | | | | | | | ✓ | ✓ | | | ✓ | | | ✓ | | | | | | | | | |
| Stamina | | | | | | | | + | | | | + | | | | | | | | | | | | |
| Storm | ✓ | ✓ | ✓ | ✓ | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ |

*Tool capabilities; + marks additions since 2019*

# QComp 2020: Tools

Tool participation in the different tracks:

| track | DFTRES | ePMC | MCSTA | MODES | MFPL | PRISM | PET | STAMINA | STORM |
|---|---|---|---|---|---|---|---|---|---|
| correct | — | — | — | — | — | — | — | — | ✓ |
| floating-p. | — | — | ✓ | — | — | — | — | — | ✓ |
| $\varepsilon$-correct | — | — | ✓ | — | — | ✓ | ✓ | — | ✓ |
| probably $\varepsilon$ | ✓ | — | ✓ | ✓ | — | ✓ | ✓ | ✓ | ✓ |
| often $\varepsilon$ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| often $\varepsilon$ (10') | ✓ | — | ✓ | ✓ | ✓ | — | ✓ | ✓ | ✓ |

→ specialised tools and generalists:
   focus on specific algorithm vs. toolset

# QComp 2020: Tuning the Tools

Some tools provide many options and algorithms.

*Which to use to win the competition?*

Default configuration: evaluate tool like a non-expert user

Specific tuning per instance: showcase the tool's abilities

QComp 2020:

default =  configuration *per track, modelling formalism,*
*and property type* recommended by authors <u>today</u>

(tool defaults
may be
historical)

specific = aggressively tuned *per instance;*
not used by all tools

# QComp 2020: Tuning the Tools



New in STORM:

  *automatic selection of analysis configuration
  based on syntactic aspects of the benchmark*

  …using a decision tree learned from the QComp benchmarks

  → default/specific distinction now pointless

  Q: do we compare tools or algorithms?

Pragmatic solution for QComp 2020:

  STORM + STORM-STATIC
    ↑               ↑
  automatic      as in QComp 2019

# QComp 2020: The Results

# QComp 2020: Results

100 _benchmark instances_, from the QVBS

⟨model, parameters, property⟩

**restricted to intersection**

Quantile plots for overall comparison:



floating-point
correct track

→ observe STORM vs. STORM-STATIC

# QComp 2020: Results

$\epsilon$-correct track:

PMC tools + PET

excl. STORM



MCSTA (wins 34/88)

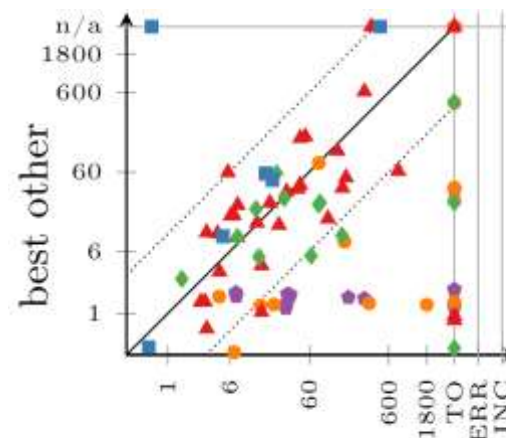PRISM (wins 15/52)

PET (wins 4/24)

ST.-static (wins 34/96)
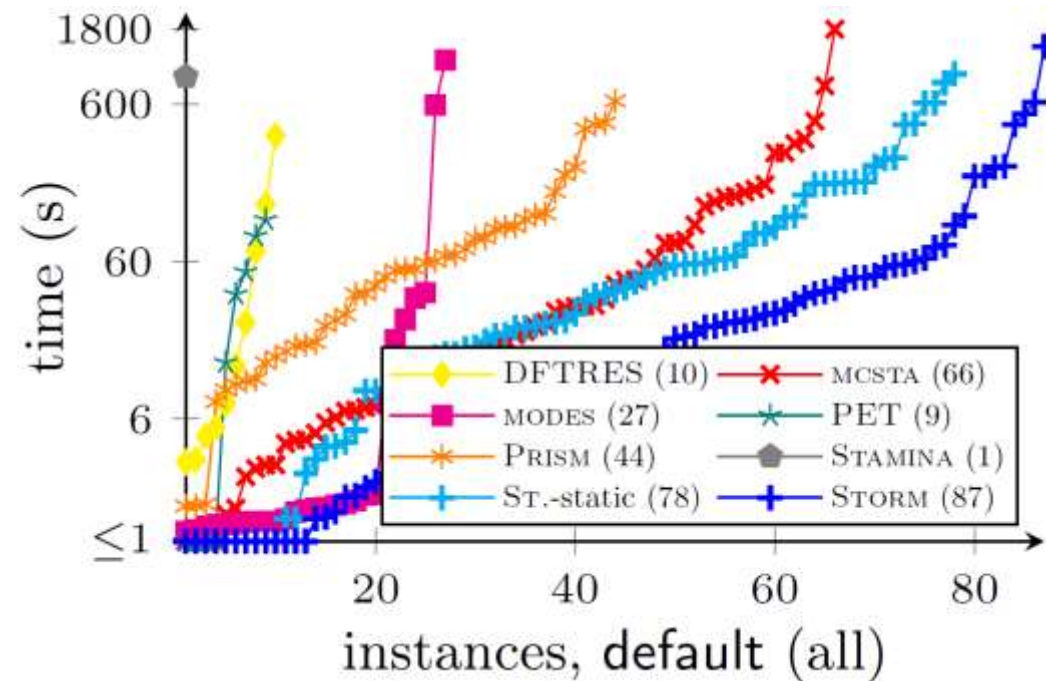
STORM (wins 54*/96)
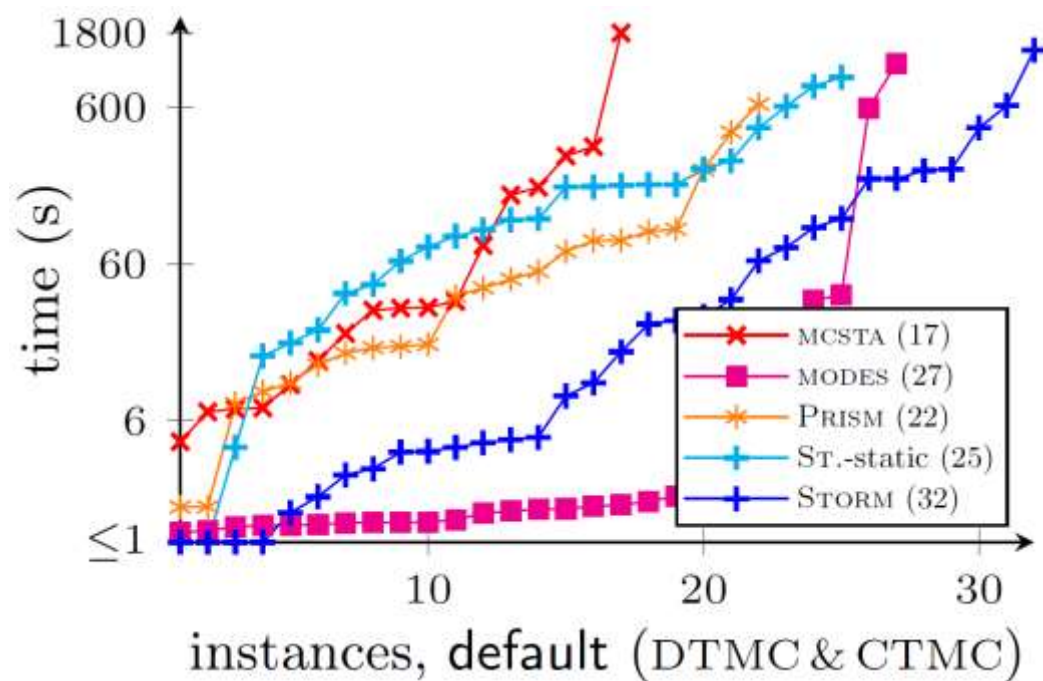
→ scatter plots show more details

# QComp 2020: Results

Probably $\epsilon$-correct track: showcase for statistical model checkers
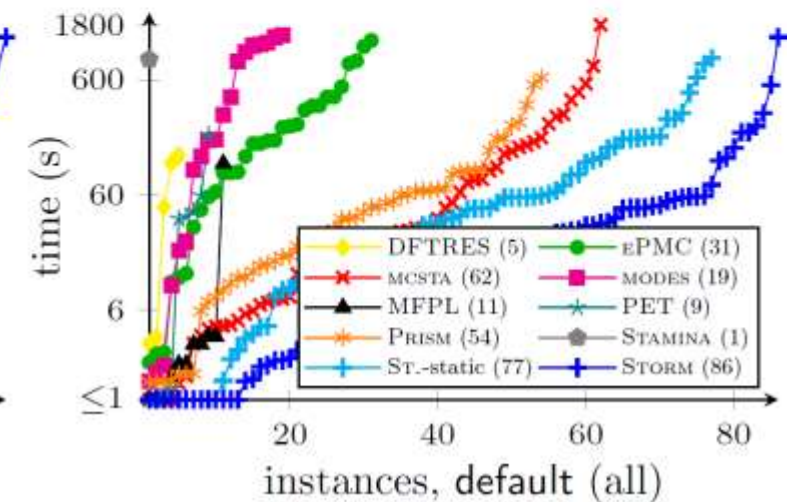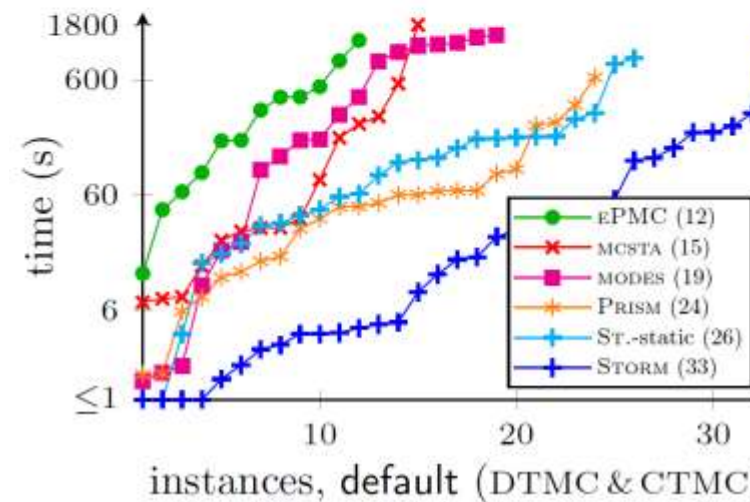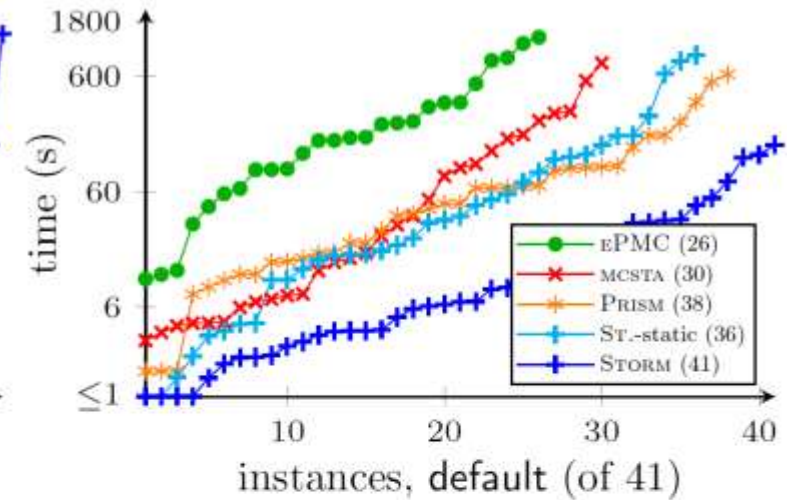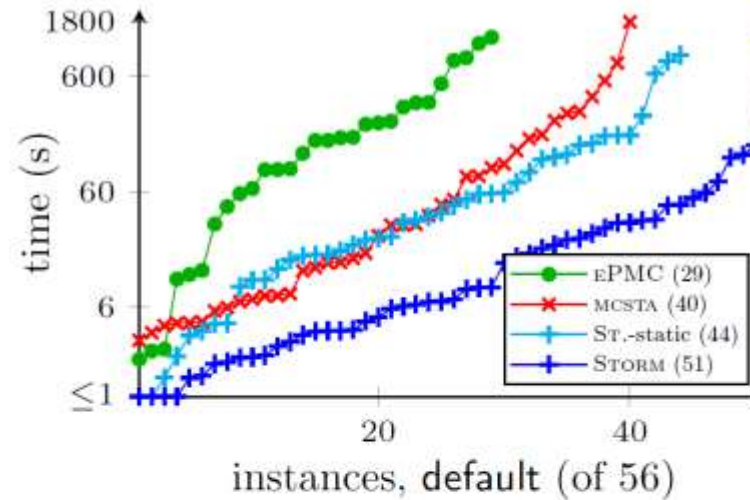


MCSTA (wins 24/88)

MODES (wins 24/32)

- ◆ DTMC  ● CTMC
- ▲ MDP   ◆ MA
- ■ PTA



instances, default (DTMC & CTMC)

- MCSTA (17)
- MODES (27)
- PRISM (22)
- ST.-static (25)
- STORM (32)

instances, default (all)

- DFTRES (10)
- MODES (27)
- PRISM (44)
- ST.-static (78)
- MCSTA (66)
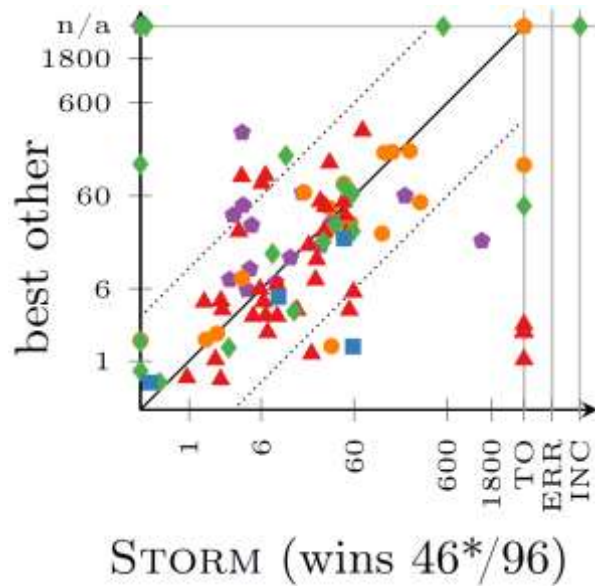- PET (9)
- STAMINA (1)
- STORM (87)

→ quantile plots show whatever you want

# QComp 2020: Results

Often $\epsilon$-correct track:
– compare with 2019
– 10' version useless

Who is the *winner*?

# Summary

Quantitative verification: PMC, SMC, and hybrid approaches

Challenges:    *algorithm bugs*          *unsound algorithms*          **specific to quantitative setting**

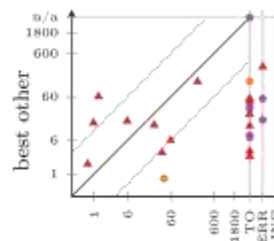*implementation bugs*                    *statistical error*
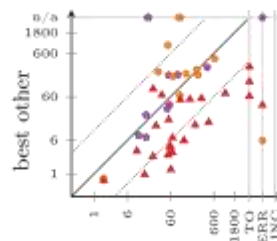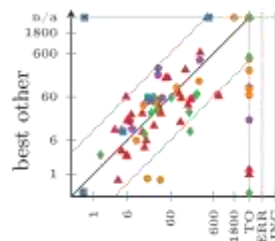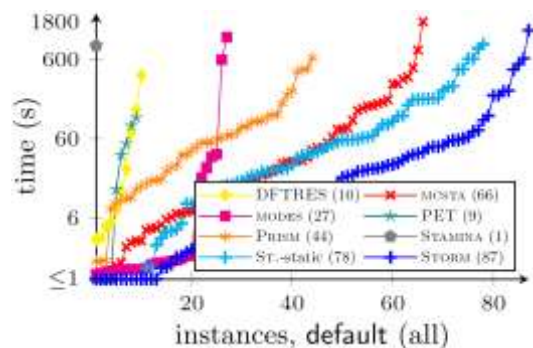
no exact results      $|v - \bar{v}|/v \leq \epsilon$          *floating-point errors*

$\mathbb{P}(|v - \hat{v}| > \epsilon) < 1 - \delta$

# QComp 2020:    5 tracks    100 benchmarks

9 tools    default + specific



qcomp.org

On Correctness, Precision, and Performance in Quantitative Verification*
QComp 2020 Competition Report

Carlos E. Budde[1], Arnd Hartmanns[1], Michaela Klauck[2],
Jan Křetínský[3], David Parker[4], Tim Quatmann[5],
Andrea Turrini[6,7], and Zhen Zhang[8]

[1] University of Twente, Enschede, The Netherlands
University, Saarland Informatics Campus, Saarbrücken, Germany
University of Munich, Munich, Germany
Birmingham, Birmingham, UK
Germany
of Software,

+ a tuned STORM

# FormaliSE 2022

*10th Int. Conference on Formal Methods in Software Engineering*

Co-located with ICSE 2022
May 22-23, Pittsburgh, USA

Deadlines (tentative):
Jan 20: paper submission
Jan 27: artifacts (voluntary)

Papers: 10 pages, ACM format

*…more info at*
**formalise.org**

# RRRR 2022

*1st Workshop on Reproducibility & Replication of Research Results*

Co-located with ETAPS 2022
April 2, Munich, Germany

Deadlines (tentative):
Feb 1:   short papers (6 pages)
Feb 15: extended abstracts

Informal proceedings,
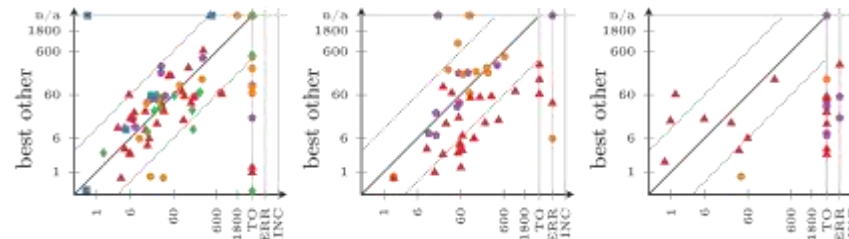extended papers in STTT
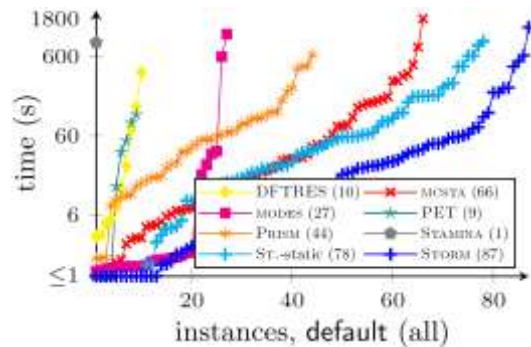
*…see* **qcomp.org/rrrr/2022**

# Summary

Quantitative verification: PMC, SMC, and hybrid approaches

Challenges: *algorithm bugs* *unsound algorithms*

*implementation bugs* *statistical error*

**specific to quantitative setting**

no exact results

$$|v - \bar{v}|/v \leq \epsilon$$
$$\mathbb{P}(|v - \hat{v}| > \epsilon) < 1 - \delta$$

*floating-point errors*

## QComp 2020: 5 tracks   100 benchmarks

9 tools   default + specific



**qcomp.org**

On Correctness, Precision, and
Performance in Quantitative Verification*
QComp 2020 Competition R...

Carlos E. Budde[1], Arnd Hartmanns[...], Mich...
Jan Křetínský[3], David Parker[4], Tim Qu...
Andrea Turrini[6,7], and Zhen Zhang...

[1] University of Twente, Enschede, The N... ...nds
...University, Saarland Informatics Campus, ...ücken, Germany
...University of Munich, Munich, Germany
...ingham, Birmingham, UK
...chen, Germany ...of Software,

*+ a tuned STORM*