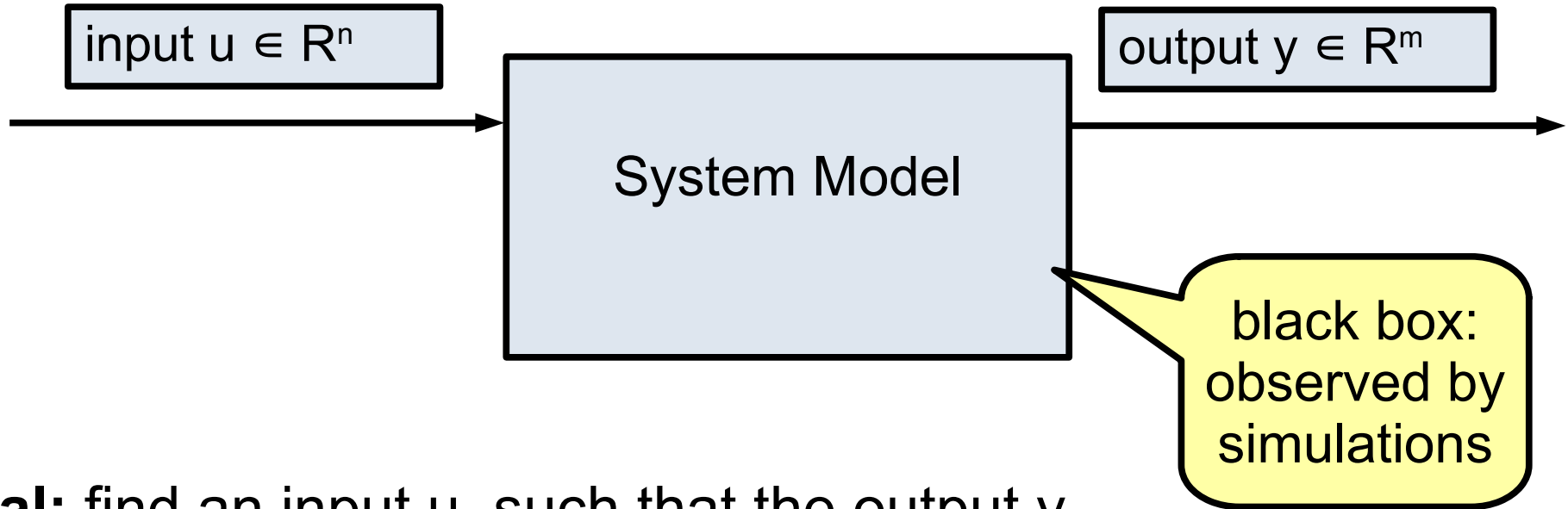


ARCH-COMP 2021: Falsification



[Gidon Ernst <gidon.ernst@lmu.de>](mailto:gidon.ernst@lmu.de), Paolo Arcaini, Ismail Bennani, Aniruddh Chandratre, Alexandre Donze, Georgios Fainekos, Goran Frehse, Khouloud Gaaloul, Jun Inoue, Tanmay Khandait, Logan Mathesen, Claudio Menghi, Giulia Pedrielli, Marc Pouzet, Masaki Waga, Shakiba Yaghoubi, Yoriyuki Yamagata, and Zhenya Zhang

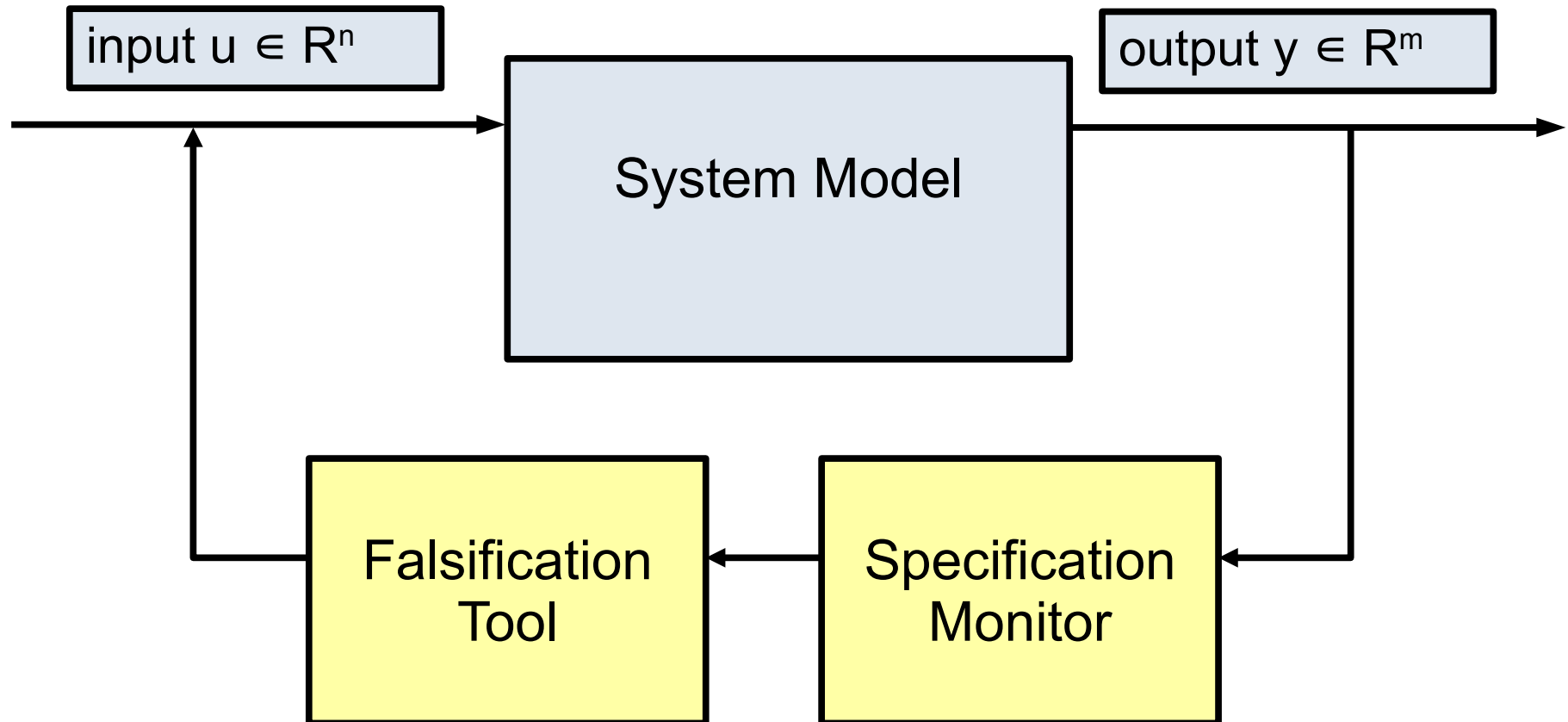
Falsification



Goal: find an input u , such that the output y violates a given specification in temporal logic

how to crash a car into a wall

Falsification



Δ 2021

- 1) two new tools, some new people (welcome!)
- 2) new conjunctive requirements (thanks, Logan!)
- 3) effort towards validating results (thanks everyone!)

Participants

- **FalCAuN** (M. Waga) discrete abstraction + model cecking
- **ForSee** (Z. Zhang, P. Arcaini, I. Hasuo) MCTS, robustness, optimization

- **ARIsTEO** (K. Gaaloul, C. Menghi) surrogate model, refinement
- **Breach** (A. Donze) robustness + stoch. optimization
- **falsify** (J. Inoue, Y. Yamagata) reinforcement learning
- **FalStar** (G. Ernst, S. Sedwards) Las-Vegas tree search
- **S-TaLiRo** (S. Yaghoubi, L. Mathesen, A. Chandratre, T. Khandait, G. Pedrielli, G. Fainekos) robustness + stoch. optimization

- **zlscheck** (I. Bennani, M. Pouzet, G. Frehse) whitebox gradient descent

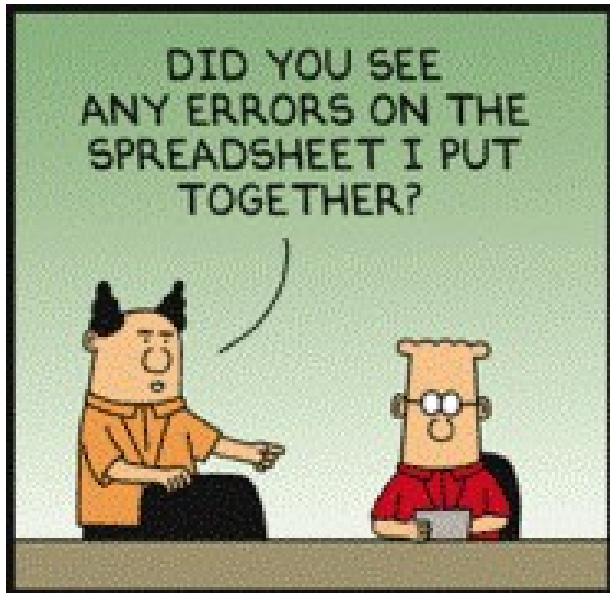
Requirements

- Example (automatic transmission):
 - within 20s the speed remains below 120
 - formally: $\square_{[0,20]} (v < 120)$
- Example (magnetic levitation controller):
 - the position stabilizes in a certain way
 - $\diamond_{[0.0, 1.0]} (Pos > 3.2) \wedge$
 $\diamond_{[1.0, 1.5]} (\square_{[0, 0.5]} 1.75 < Pos < 2.25) \wedge$
 $\square_{[2.0, 3.0]} (1.825 < Pos < 2.175)$

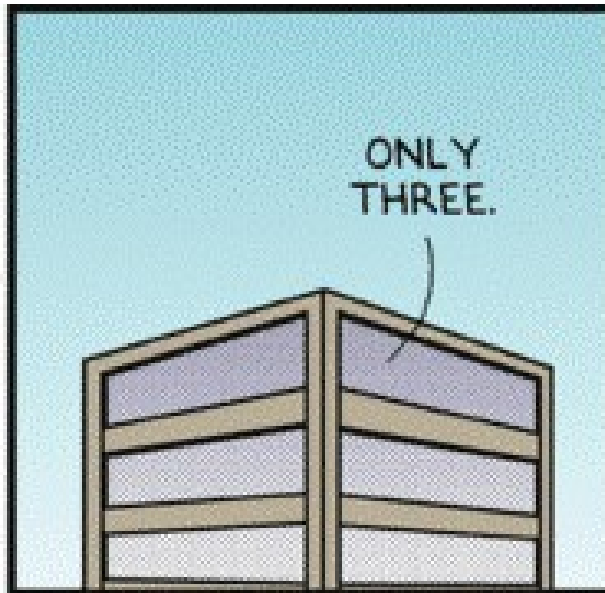
check multiple
requirements
at once

but possibly
conflicting
search guidance

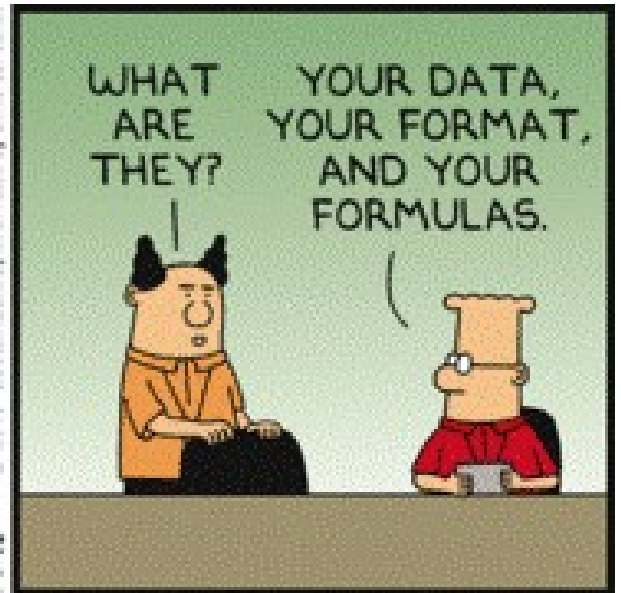
Validation



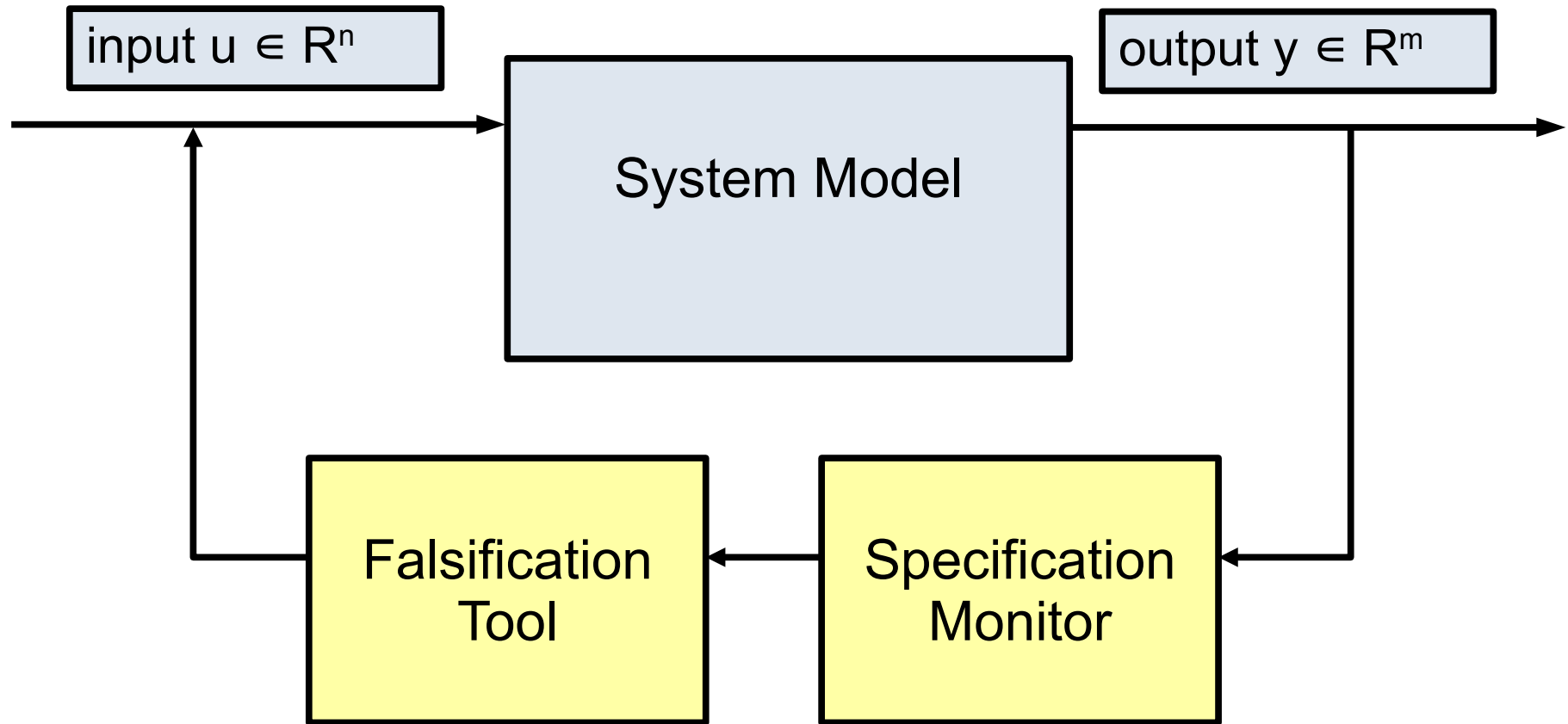
Dilbert.com @ScottAdamsSays



1-4-16 © 2016 Scott Adams, Inc. /Dist. by Universal Uclick



Error Sources



Error Sources

search space
definition

computational
differences

input $u \in \mathbb{R}^n$

output $y \in \mathbb{R}^m$

System Model

configuration
mismatch

formalization
mistakes

Falsification
Tool

Specification
Monitor

abstraction
techniques

computational
differences

2021/FALS/Validation.md

formalization

property: **AT1** ↘
formula: $\square_{[0, 20]} (\text{speed} < 120)$

search space

input is within bounds

using stop time as provided: 20.0

reported verdict

falsified is correct

expected robustness -0.014

computed robustness -0.013

simulation &
computation

robustness error 0.001

[...]

Validity via Robustness Score

by how much did we miss the wall?

reported
(by participants)

computed
(by validator)

confirmed
falsification?

$r < 0$

$r < 0$

$r < 0$

$r < \text{threshold}$



$r \geq 0$

—





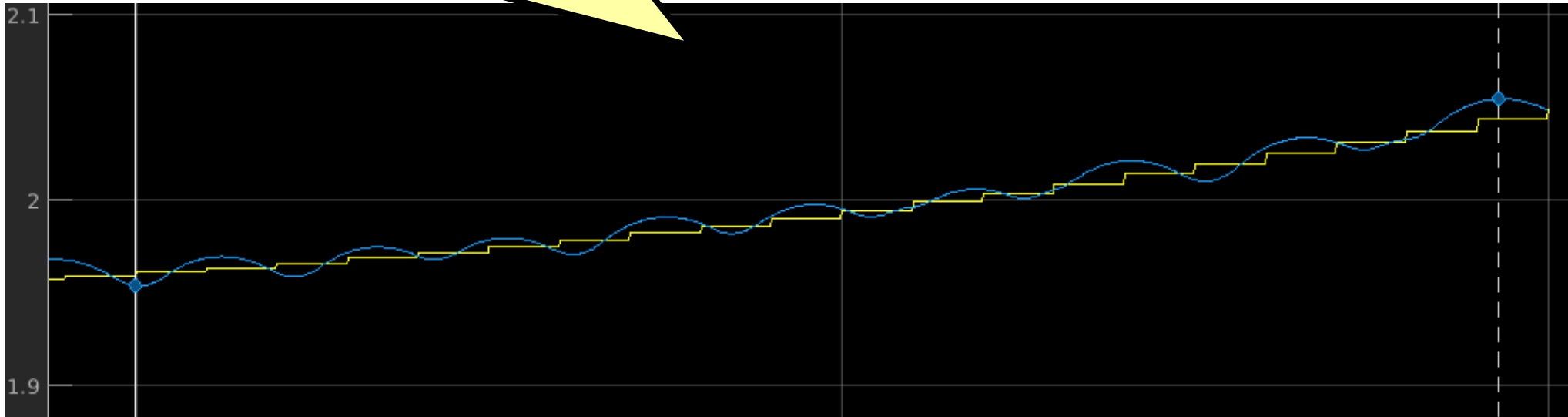
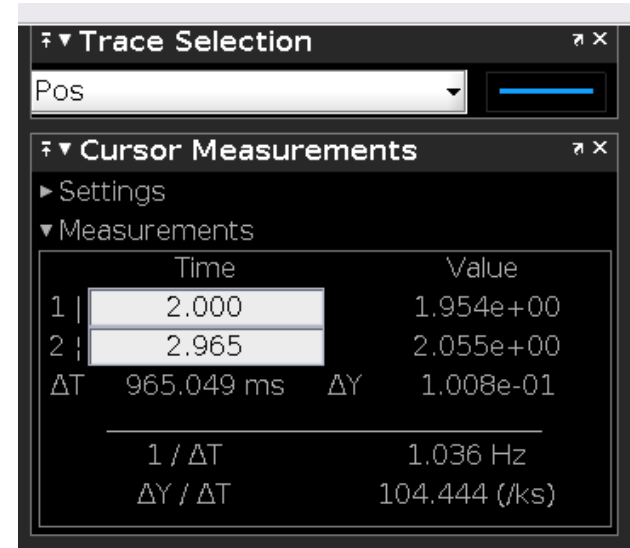
Erwin Wurm: Truck, 2015

Findings

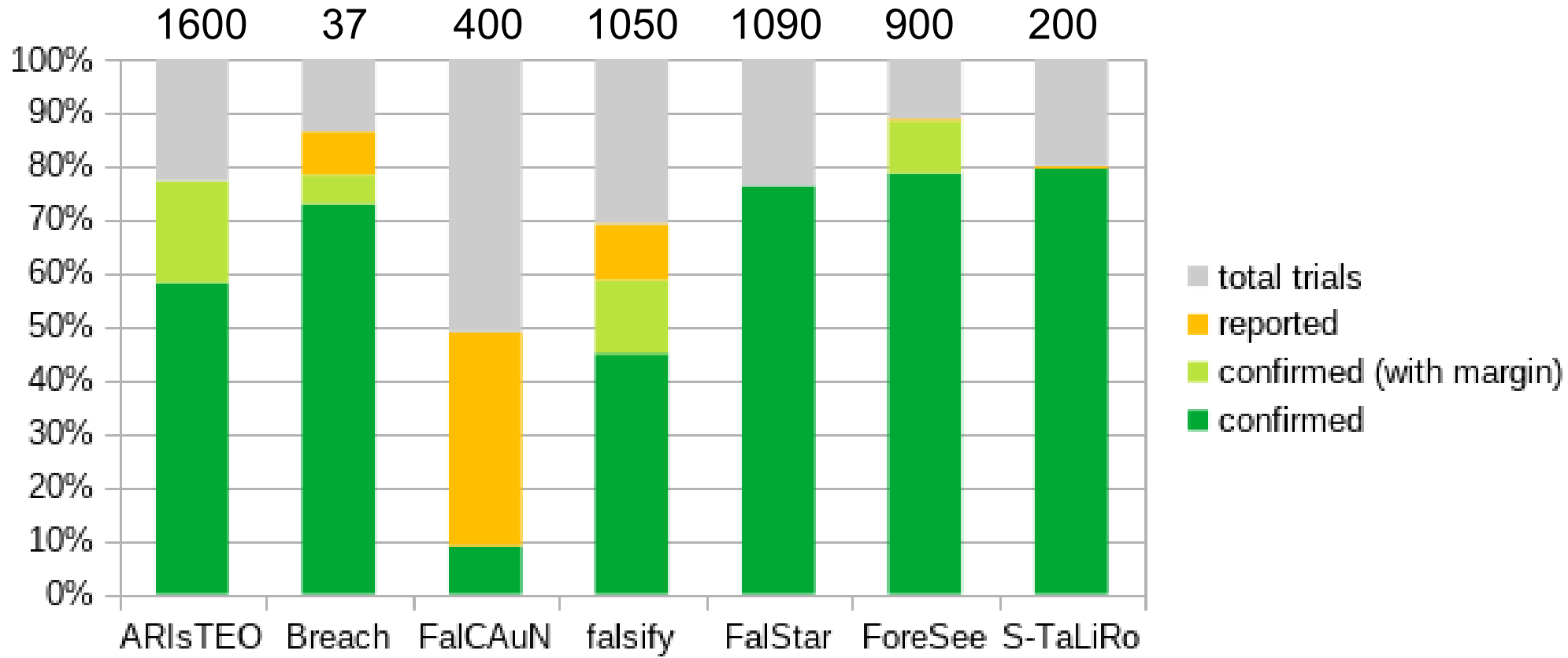
- Discrepancies **discovered** everywhere
 - documentation somewhat lacking/imprecise/inconsistent
 - some mistakes: formalization, input ranges
 - signal sampling and reporting (granularity, interpolation)
- Summary
 - **majority of new results confirmed** | **some old results wrong**
 - **obvious & large discrepancies now fixed**
 - **small discrepancies: computational mismatch or actual errors?**¹³

Troubleshooting

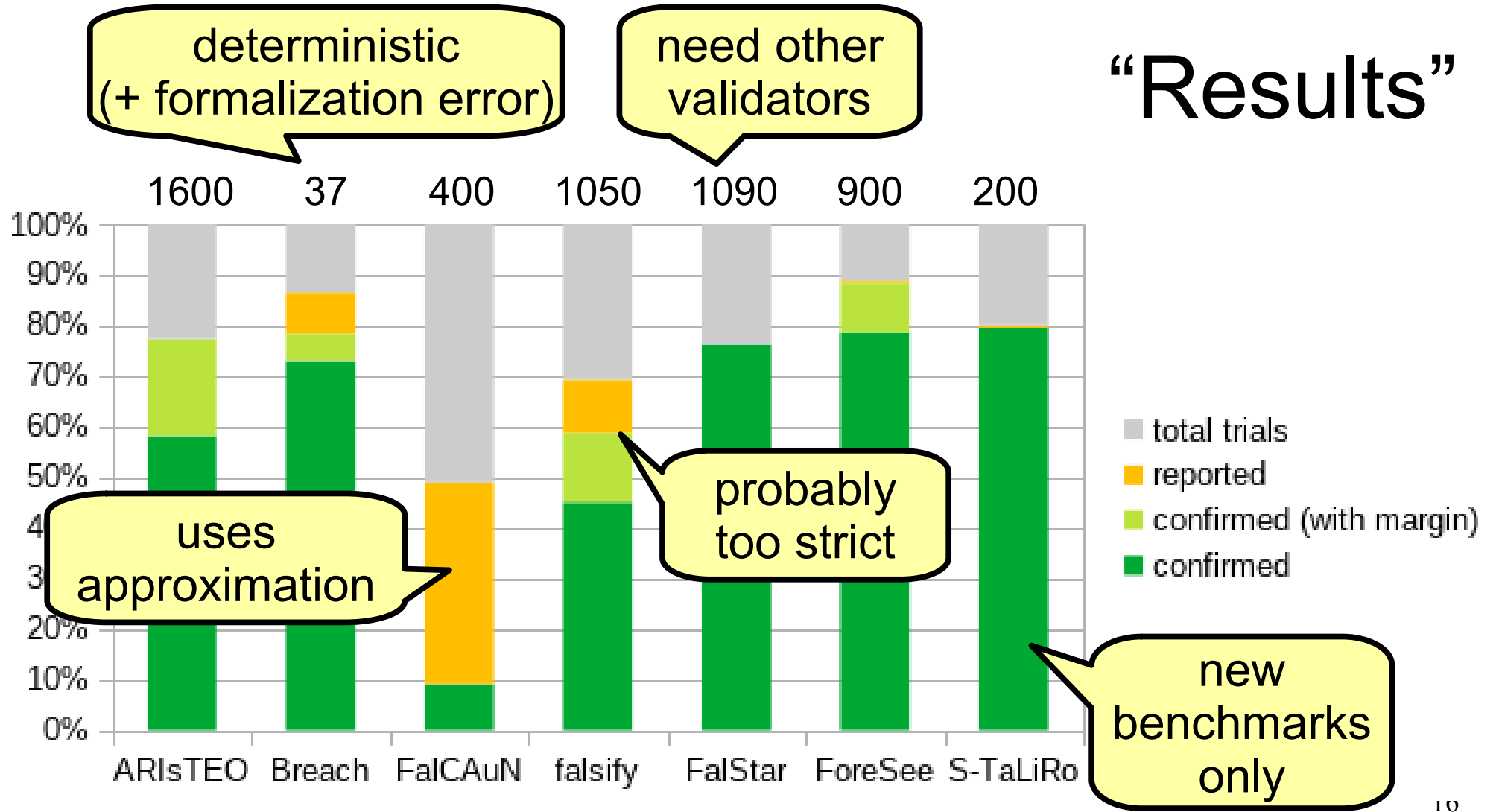
check manually:
does this signal falsify?



“Results”



“Results”



Conclusion & Outlook

- collection of results still ongoing
- both needed
 - Validation: reported results are plausible
 - Repeatability: performance is plausible
- next year
 - start much earlier, converge with reporting format
 - repeatability (on CodeOcean)
 - non-Matlab models (Zélus, Python), harder benchmarks