

Operationalizing Threats to MSR Studies by Simulation-Based Testing

Johannes Härtel and Ralf Lämmel

*University of Koblenz, Germany
Software Languages Team*

Motivation

- Empirical questions on SE can be answered by analyzing repository data.
- The methodology of such studies is often very complicated.
- ***How do we know that the methodology and answers are “correct”?***

Background

Some **existing ideas** to assure “correct” methodology and answers:

- Our **intuition** as a software developer helps us to judge answers.
- **Meta-analysis** checks the consistency of answers with previous work.
- We can **stick to** the methodology of **previous work**.
- **Model comparison** with a protection against under- and overfitting helps.

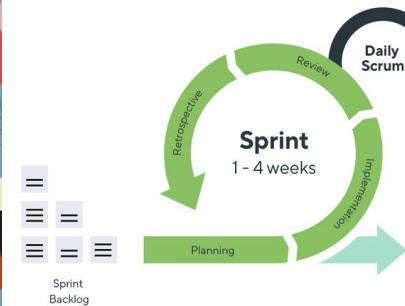
We will present a fresh idea on this topic.

We start with the typical methodology in MSR/ESE

(that you probably know...)

The typical methodology

Variables interesting for MSR/ESE.



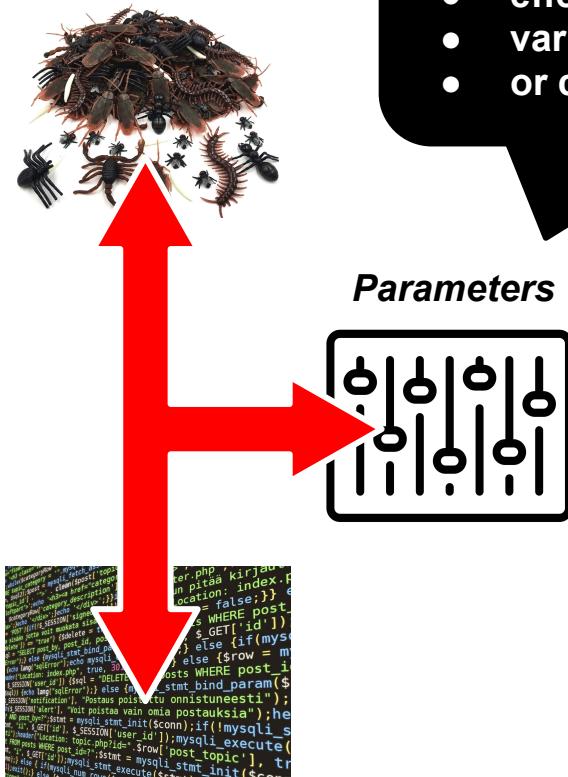
The typical methodology

Relationships between variables interesting for MSR/ESE.



The typical methodology

Models interesting for MSR/ESE.

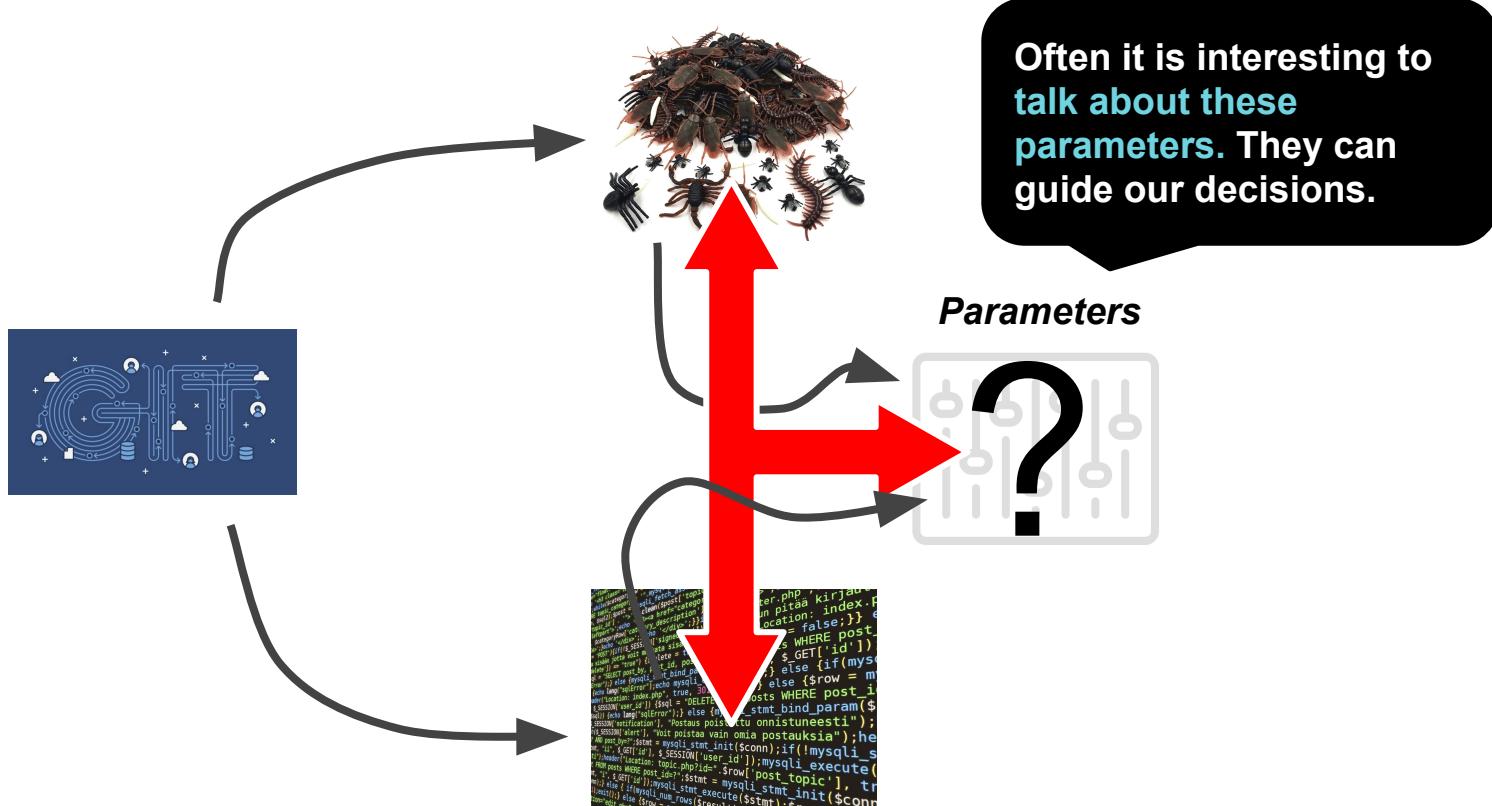


Parameters are variables too, e.g.,:

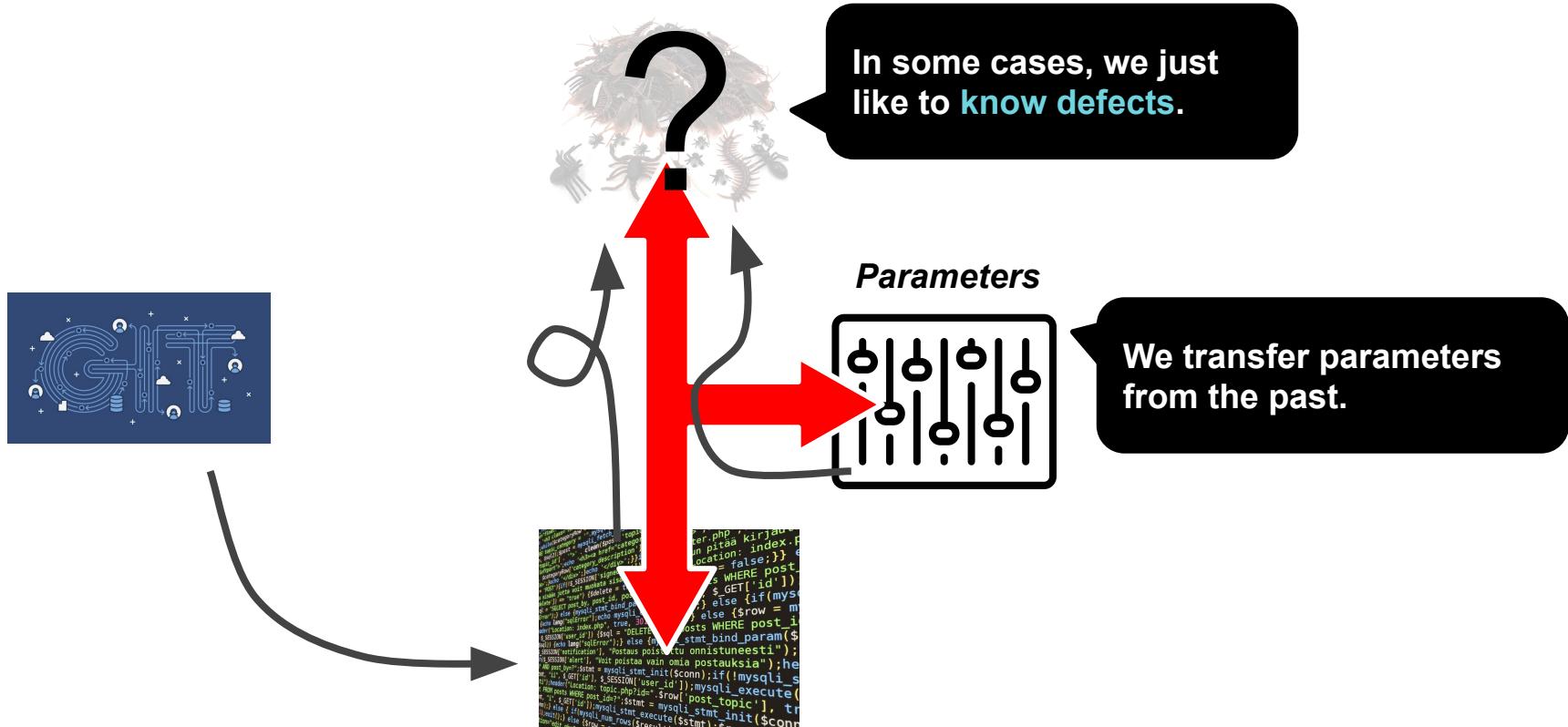
- effect strength,
- variance,
- or correlation.

What can we do with this model?

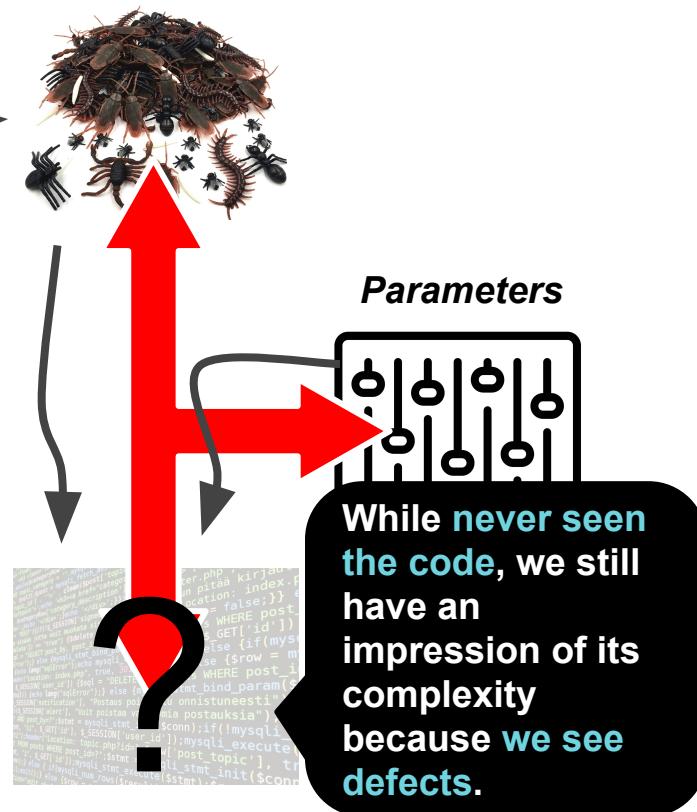
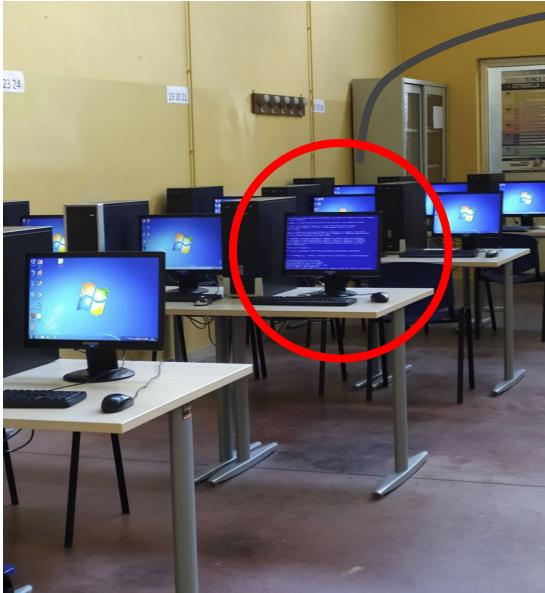
What can we do with this model?



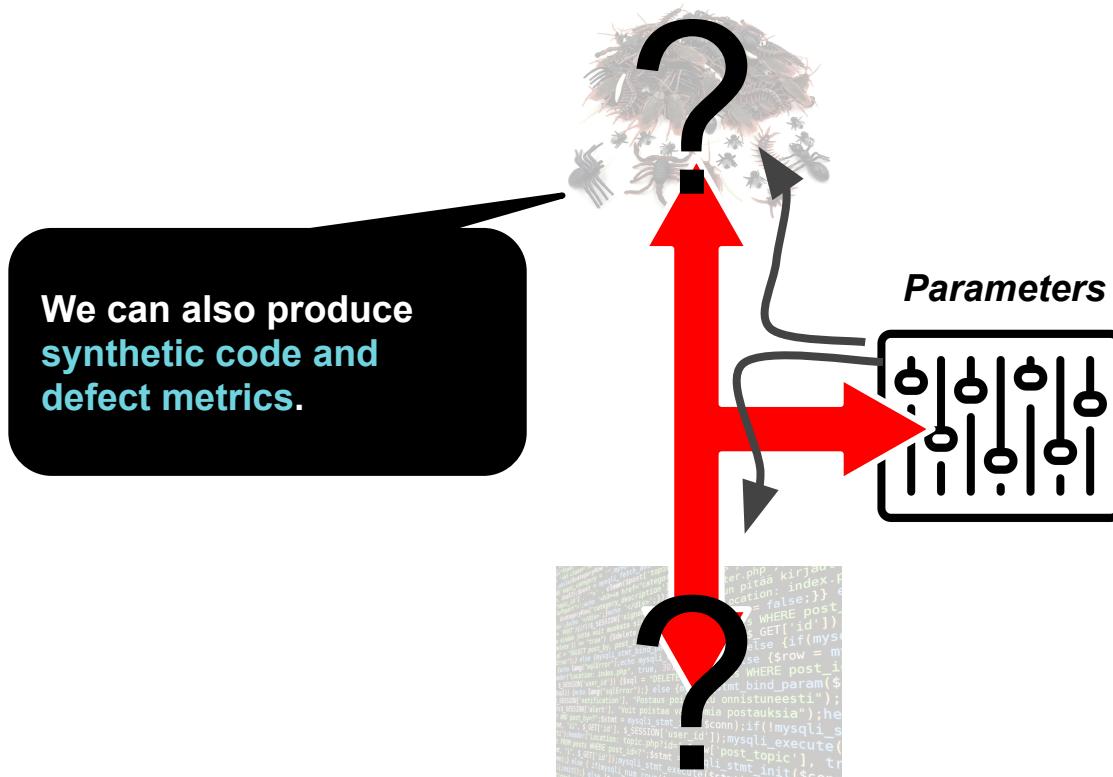
What can we do with this model?



What can we do with this model?



What can we do with this model?

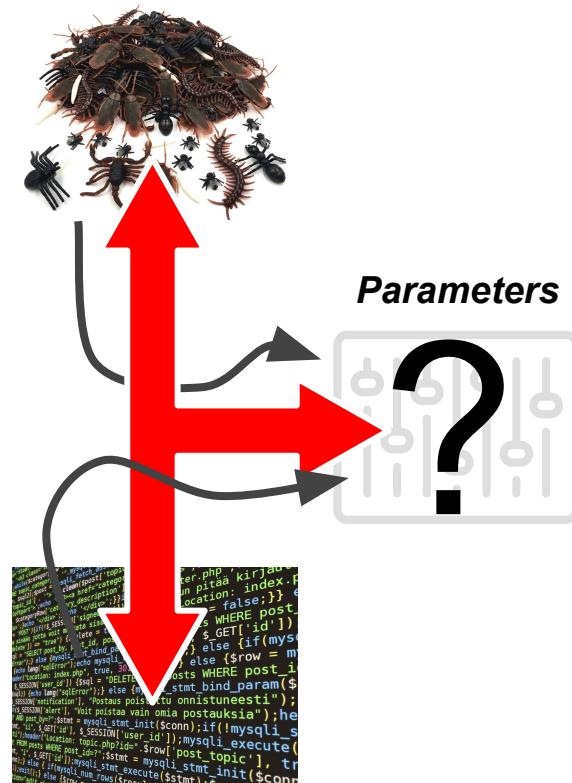
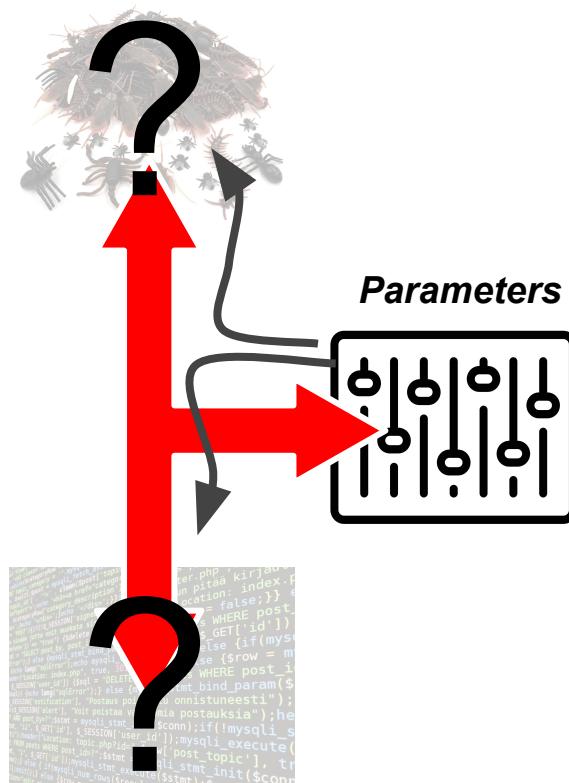


The starting point for Simulation-based Testing

(the topic of our paper)

Symmetry (0)

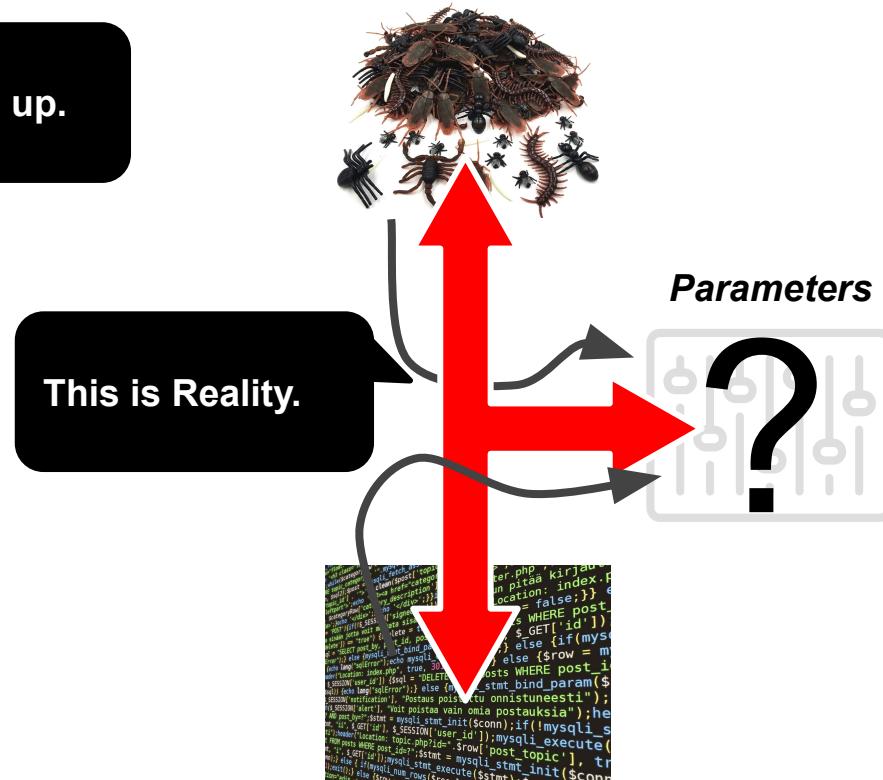
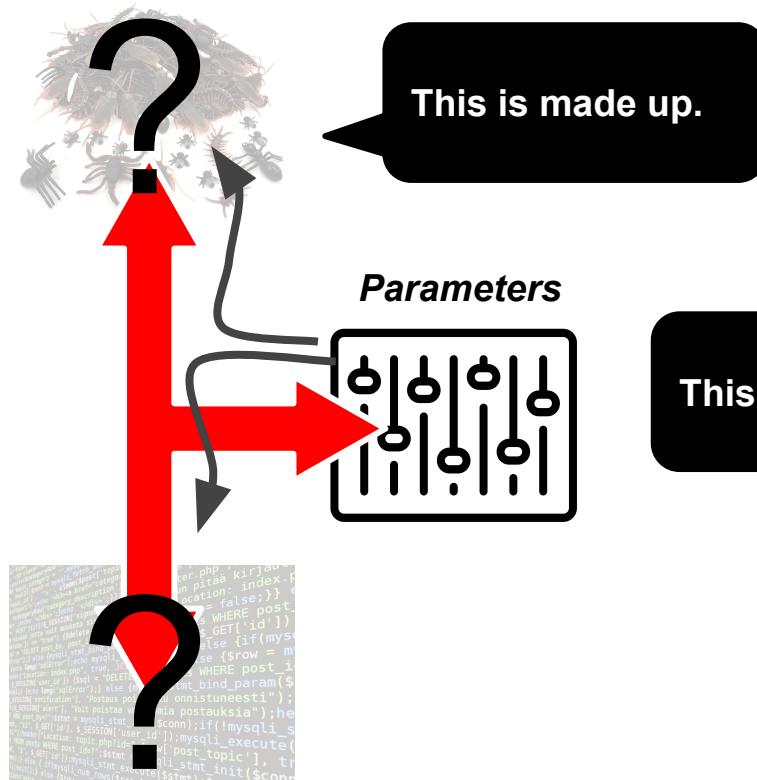
We expect a **correspondence** for those two workflows:



Symmetry (1)

We expect a **correspondence** for those two workflows, **and call them...**

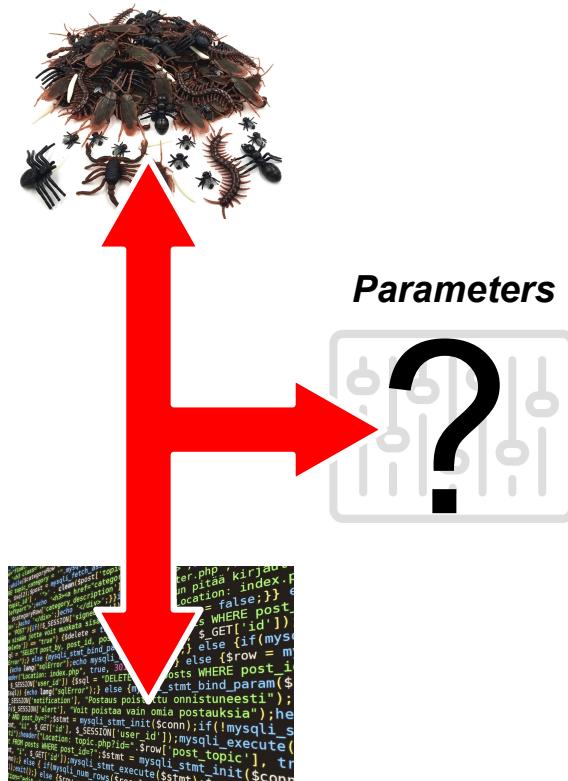
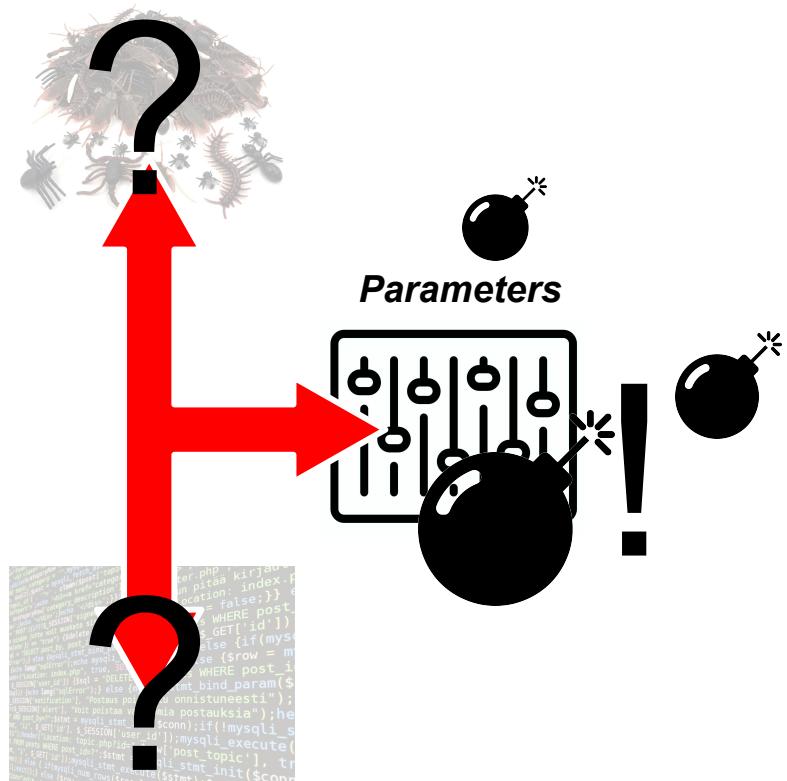
Simulation



Symmetry (2)

We set plausible but challenging **simulation** parameters.

Simulation

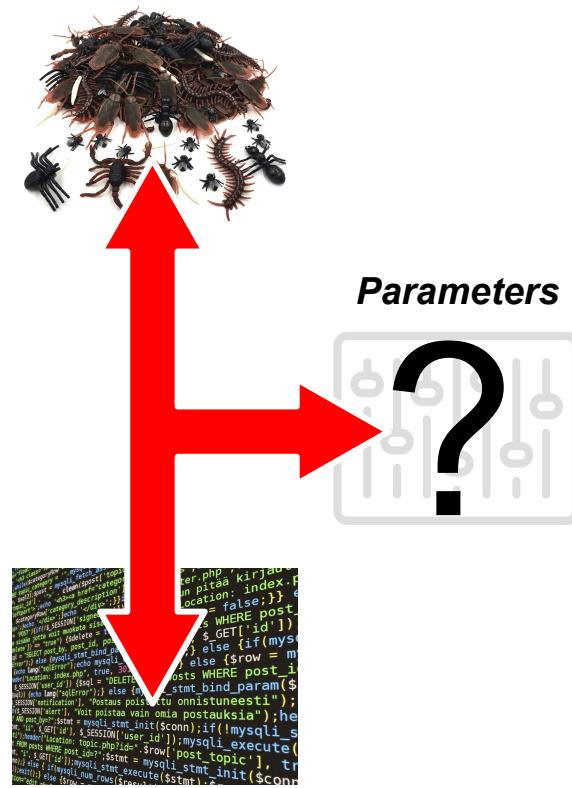
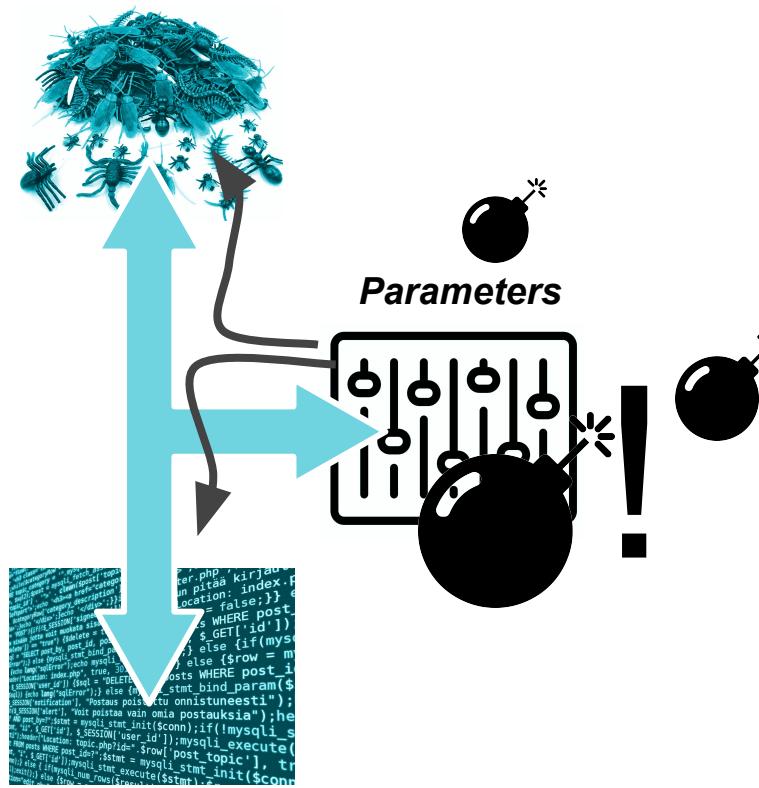


Original Methodology

Simulation

Symmetry (3)

We use relationships in revers to produce **artificial defects and code metrics**.

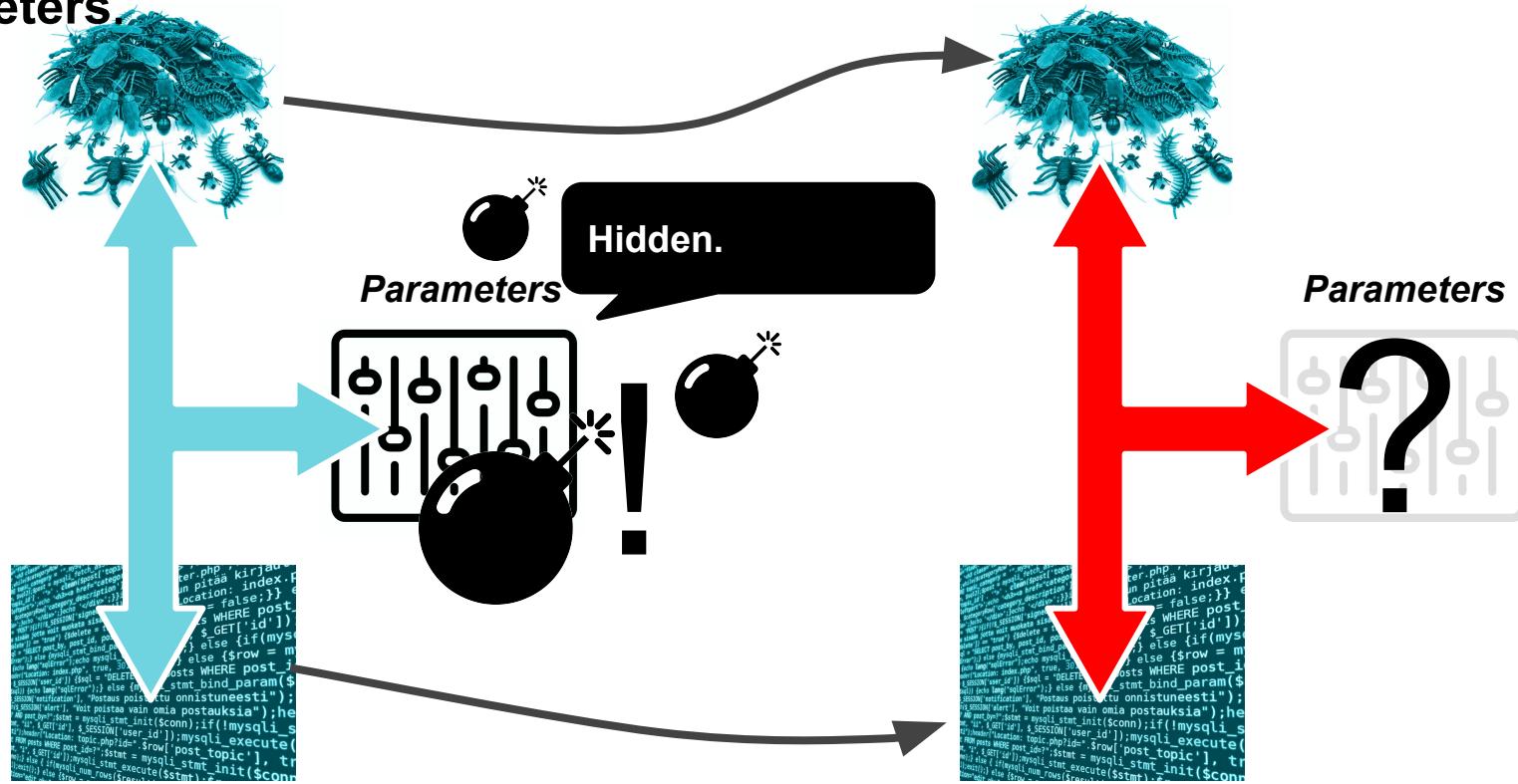


Johannes Härtel — University of Koblenz — johanneshaertel@uni-koblenz.de

Simulation

Symmetry (4)

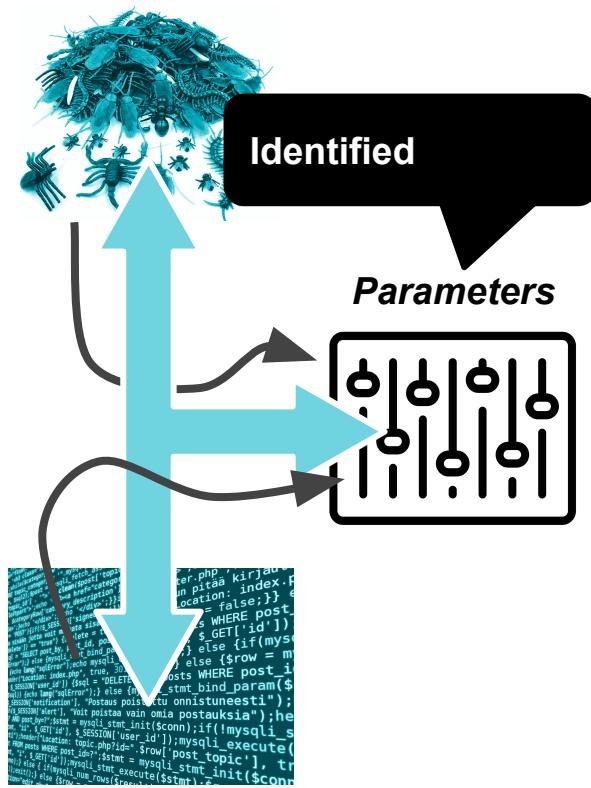
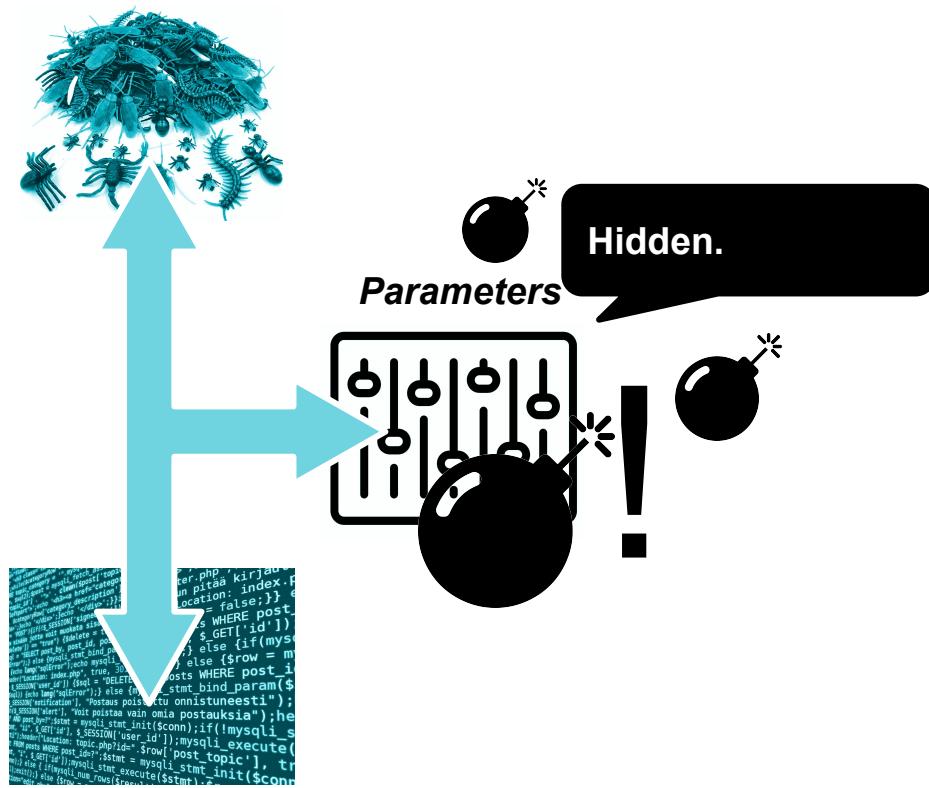
We substitute the real variables with the synthetic variables, except the parameters.



Simulation

Symmetry (5)

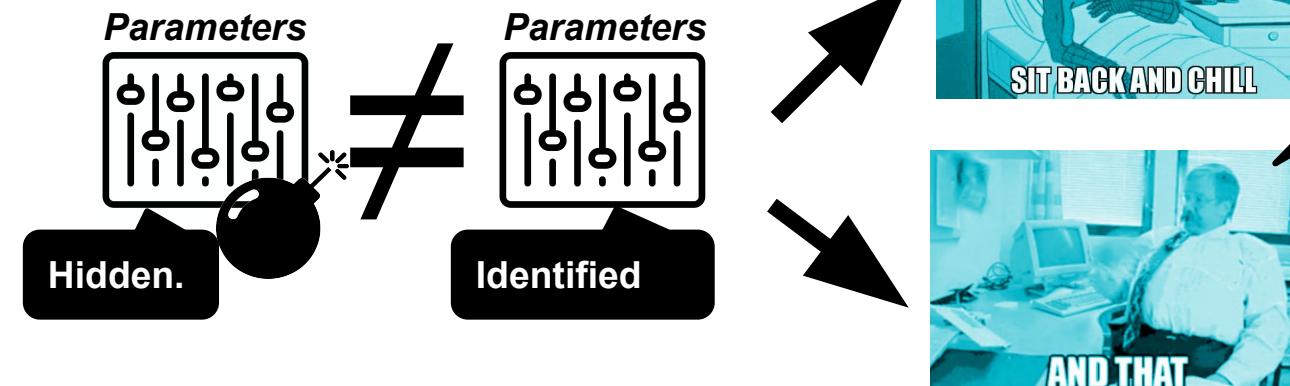
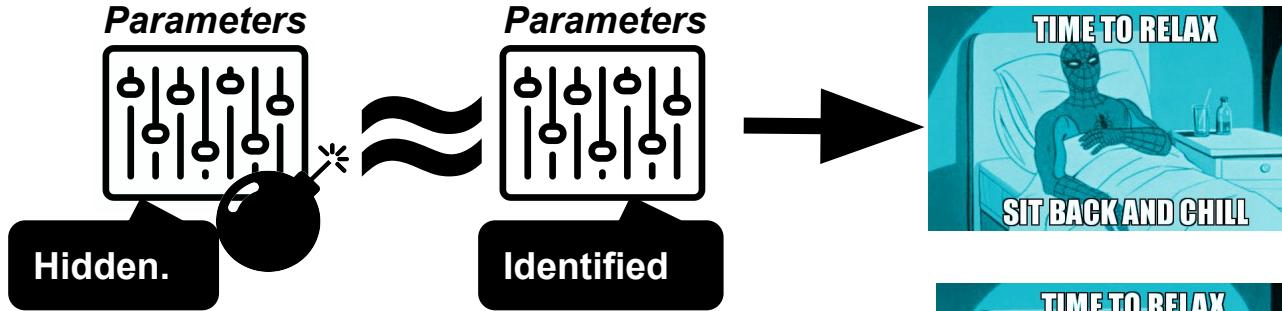
We run the original methodology to identify parameters again.



Symmetry (6)

We can check correspondence.

Nice to know that results are **correct**, I will relax now.



It is **a threat** to the study, I will relax, but list it in the paper.

I will **admit impossibility** and report on this.

Applying Simulation-based Testing (our evaluation)

What kind of bombs do we have?



Dependent observations



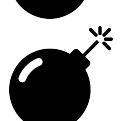
Causation vs. prediction



Control of variables



Correlated variables



... (needs to be continued by you)

We used simulation-based testing to pinpoint threats in [real studies](#) (as an evaluation).

Try the simulations that we provide online.

<https://github.com/topleet/MSR2022>