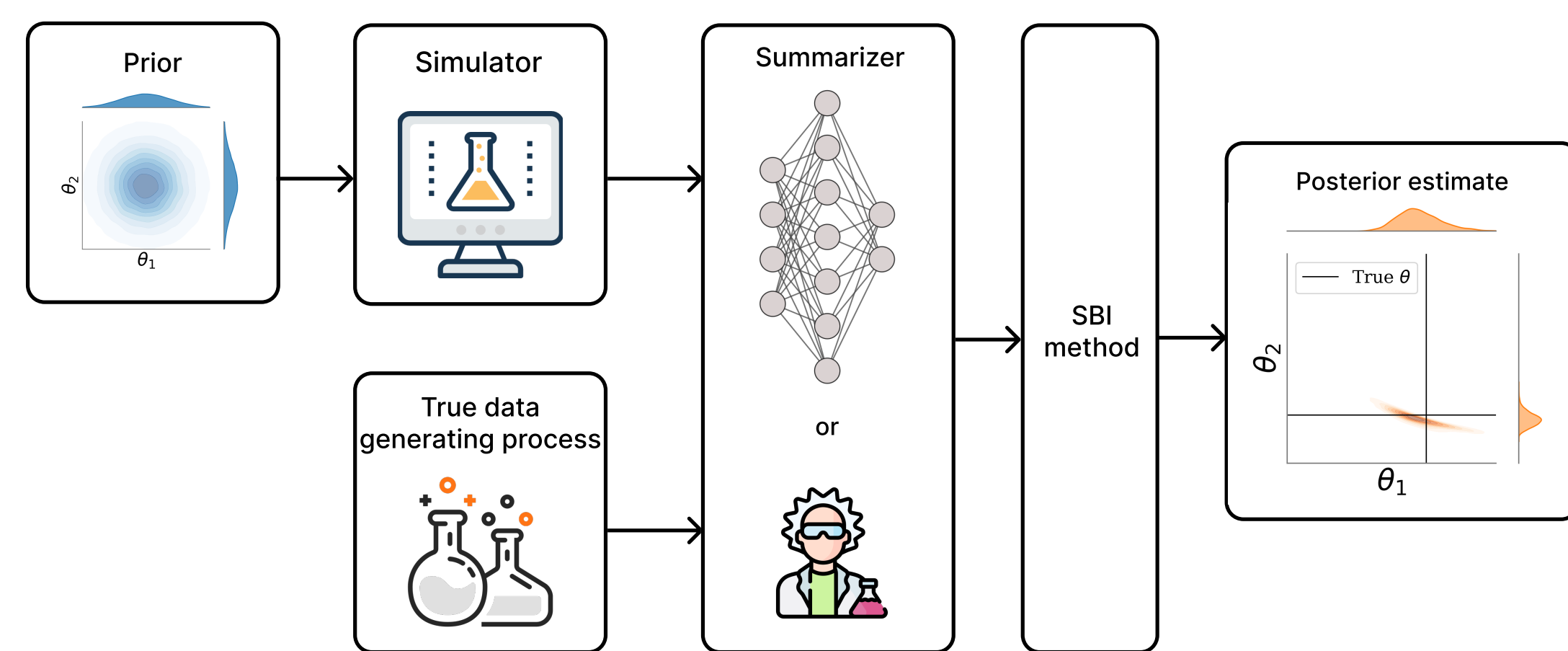


Learning Robust Statistics for Simulation-based Inference under Model Misspecification

Daolang Huang*, Ayush Bharti*, Amauri Souza, Luigi Acerbi, Samuel Kaski

Introduction

- Simulators are mechanistic models of real-world phenomenon that are used widely in many domains of science, medicine and engineering. However, their likelihood functions are intractable.
- Simulation-based inference (SBI) methods, such as approximate Bayesian computation (ABC) and neural posterior estimation (NPE), are used to fit such intractable models, which rely on simulating summary statistics.



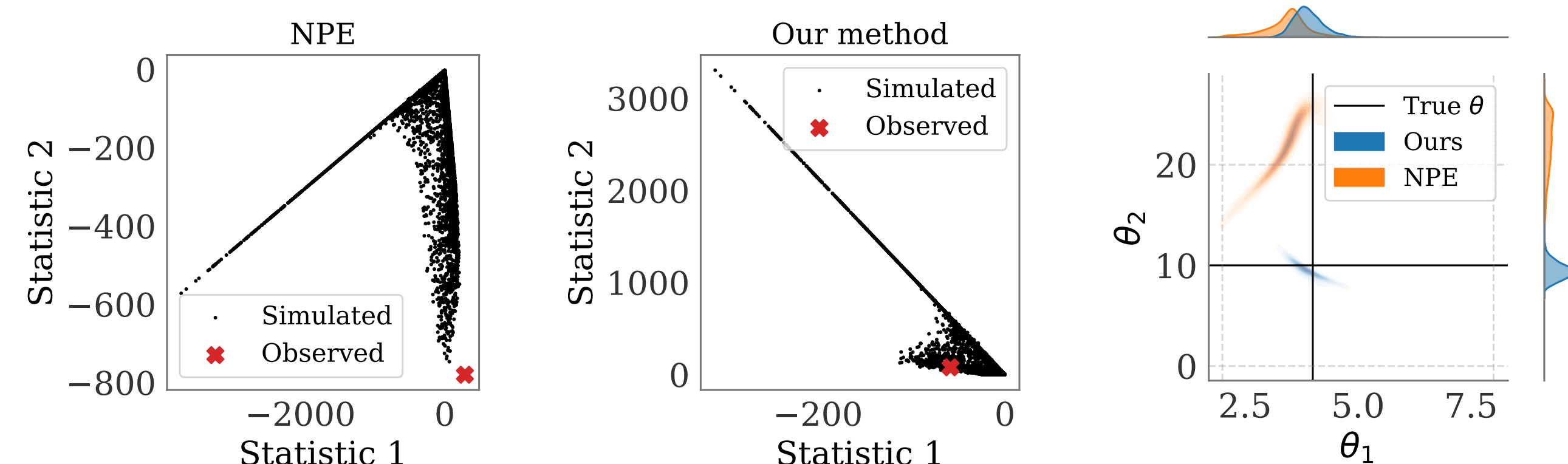
- However, **SBI methods are unreliable when the model is misspecified**, which occurs when the true data-generating process does not lie within the family of distributions defined by the model.

Method: SBI with robust statistics

Observation 1: Even if the model is misspecified, it may still be able to match the statistics, and hence, be “well-specified” in light of those statistics.

Observation 2: under misspecification, the model may be unable to match the observed statistic for any parameter, forcing the inference method to generalize outside its training distribution, causing shifts in the posterior.

⇒ **If we pick statistics appropriately, we can be robust!**

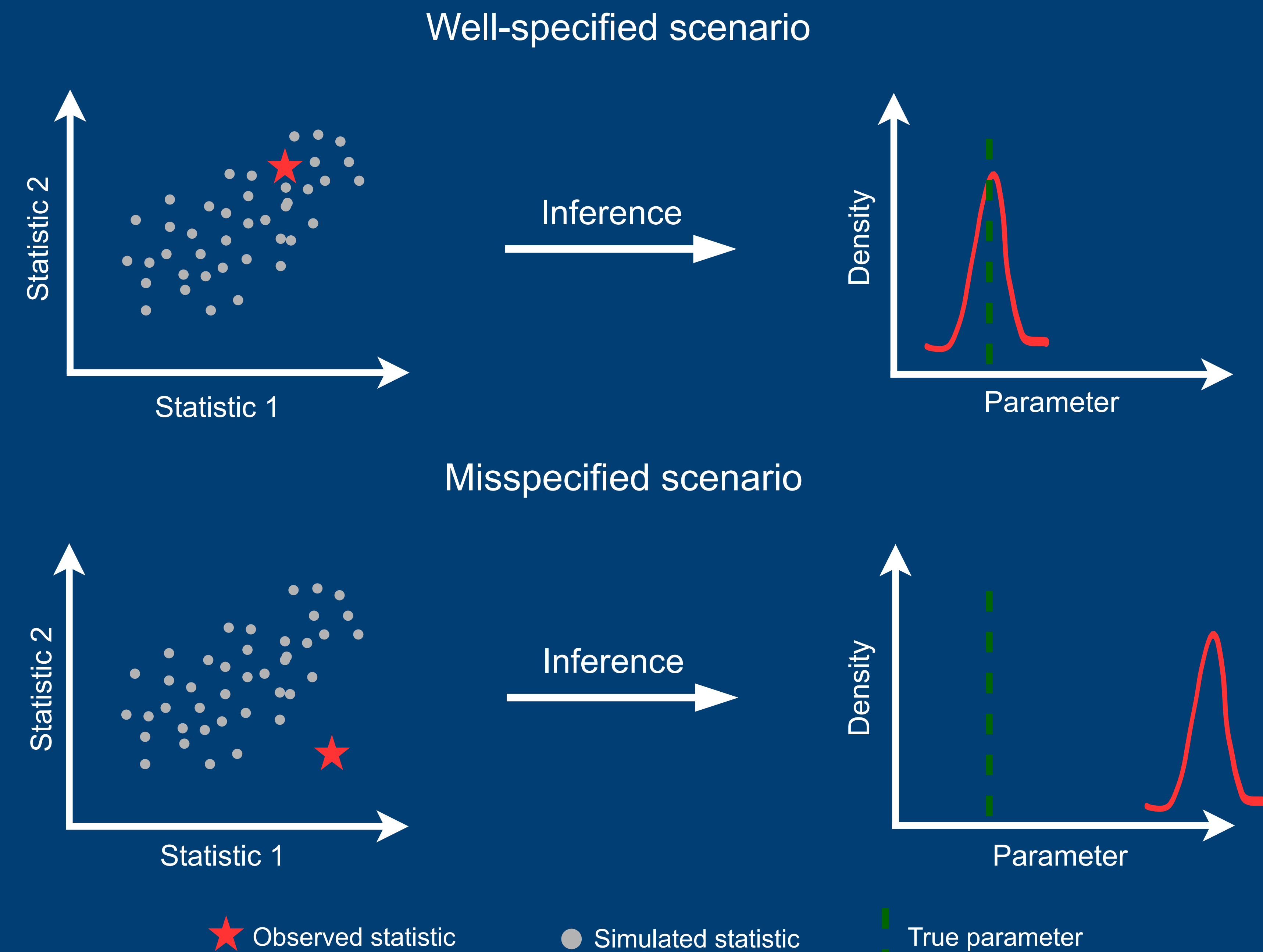


We introduce a **regularized loss function** that balances between learning statistics that are informative about the parameters, and penalizing those choices of statistics or features of the data that the model is unable to replicate.

Advantages of our method:

- Applies to all statistics-based SBI methods
- Performs reasonably well even when the model is well-specified

Robustness to model misspecification in simulation-based inference is a statistics selection problem.



Under misspecification, observed statistic becomes an out-of-distribution sample.

Idea: choose statistics s.t. the observed statistic remains an in-distribution sample

Solution: learn parameters and statistics jointly using the regularized loss function

$$\text{proposed loss} = \text{usual loss} + \lambda \mathcal{D}(\text{simulated statistics}, \text{observed statistic})$$

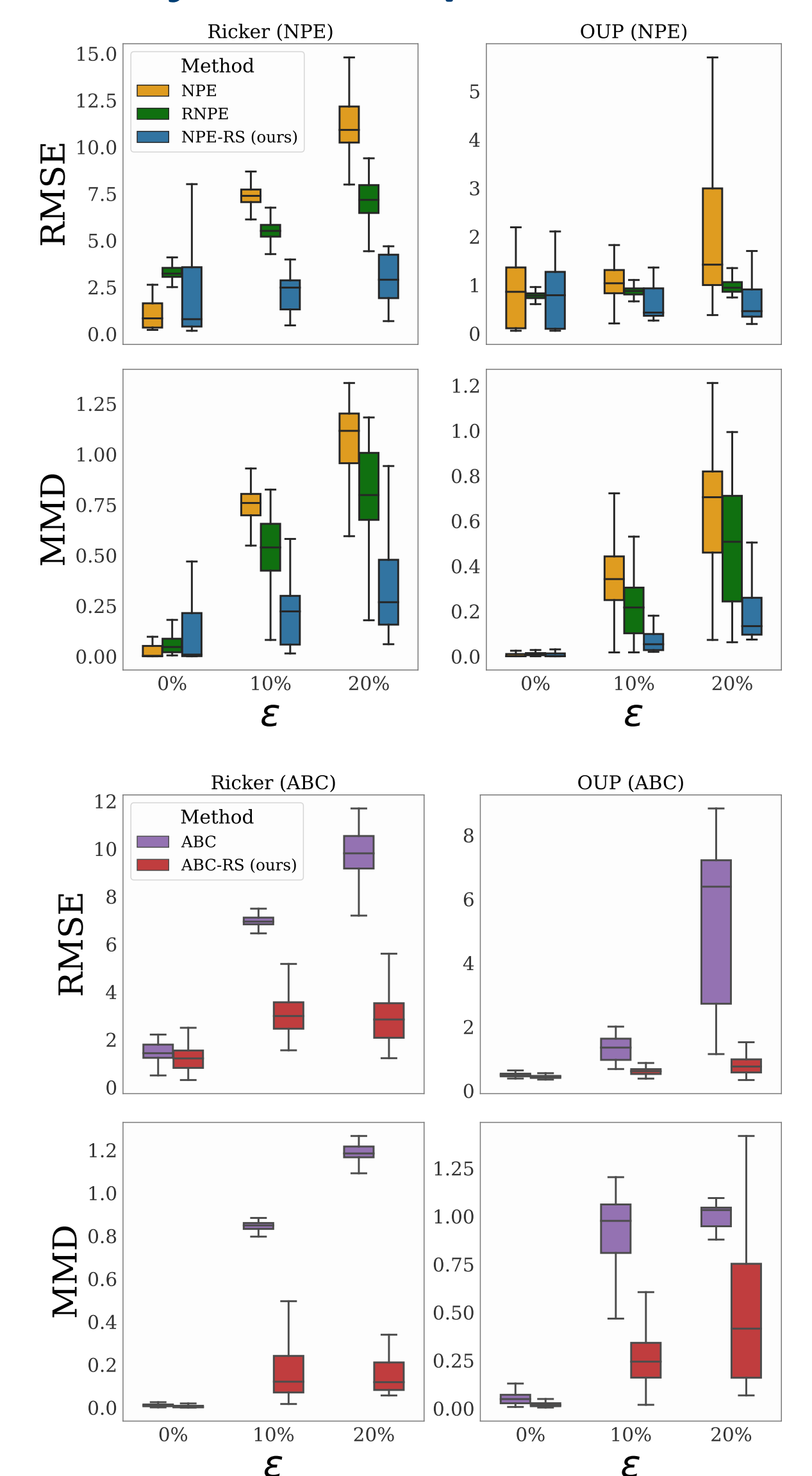
λ : encodes trade-off between efficiency and robustness

\mathcal{D} : maximum mean discrepancy

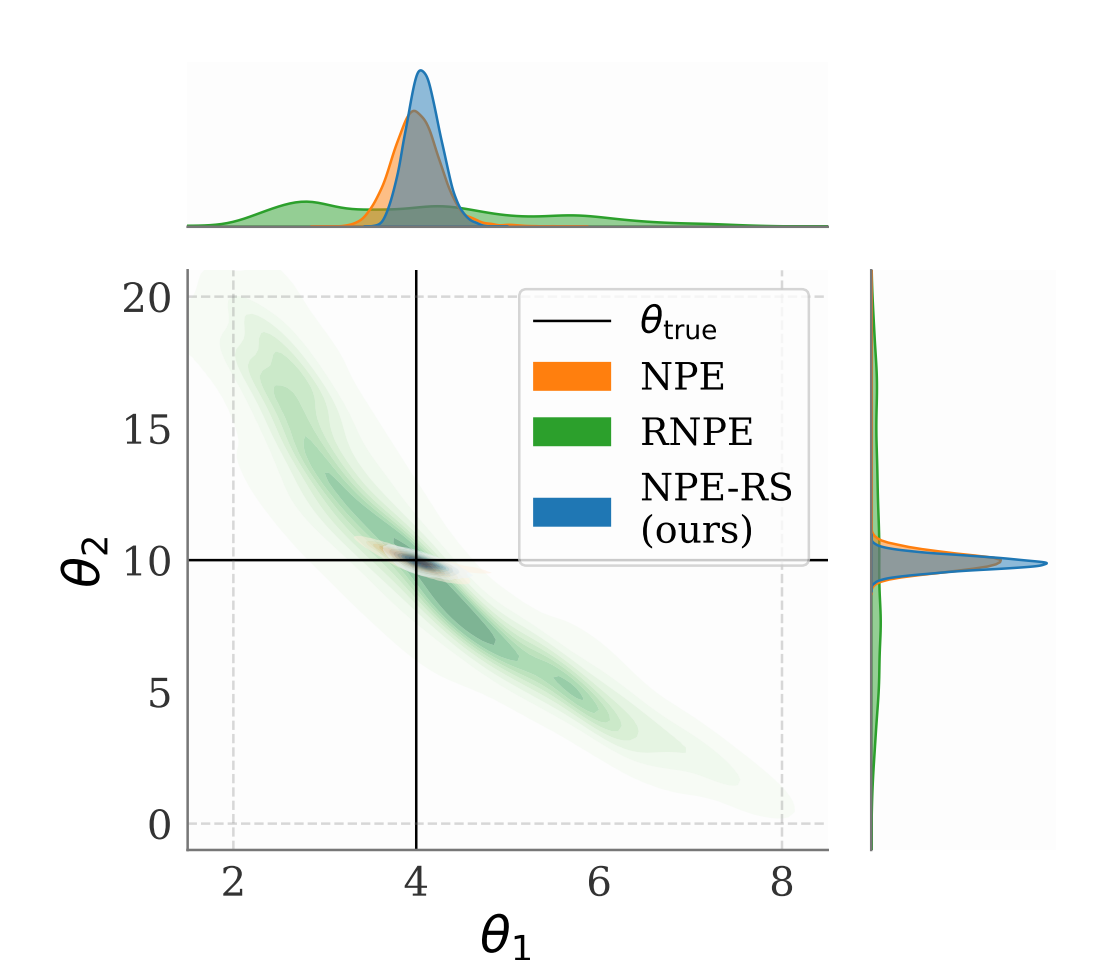


Take a picture to
download the full paper

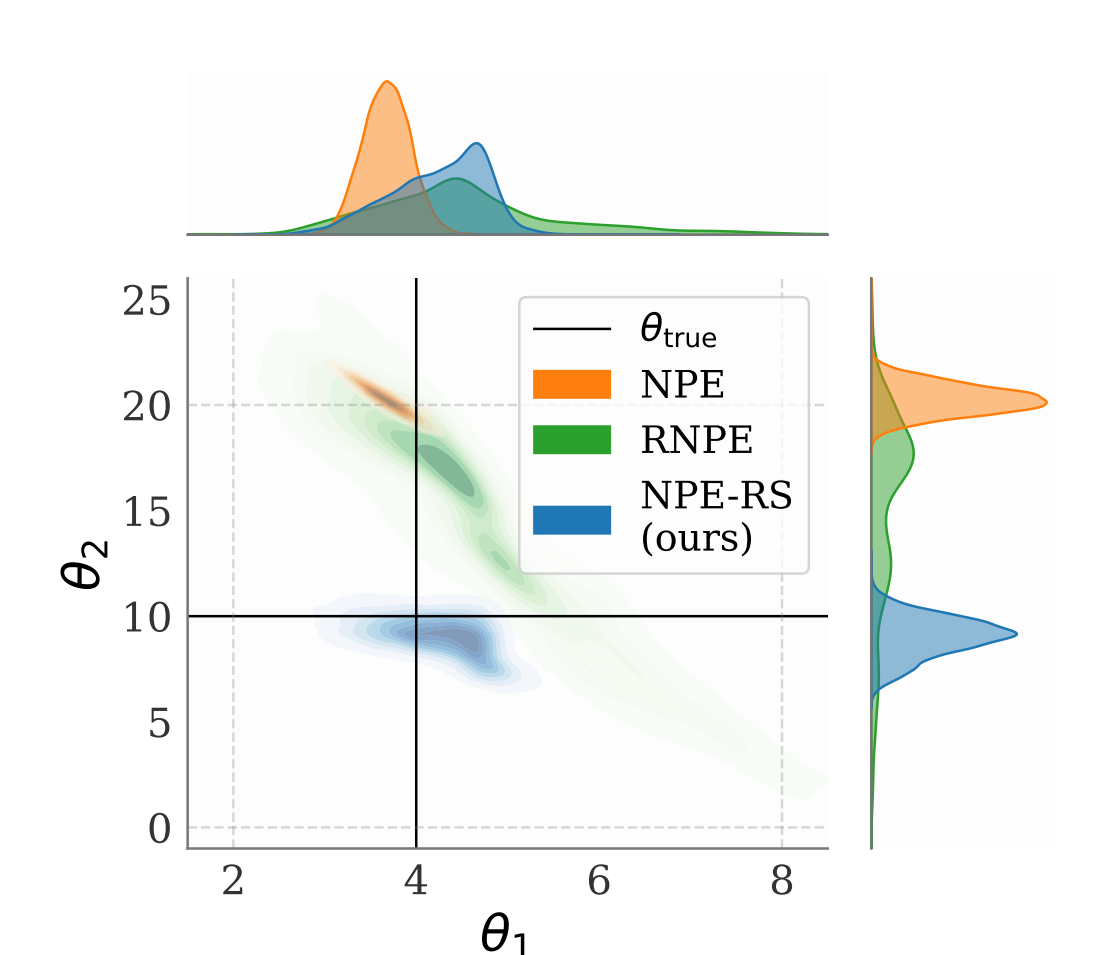
Synthetic experiments



Well-specified



Misspecified (epsilon=10%)



Real-world experiment

