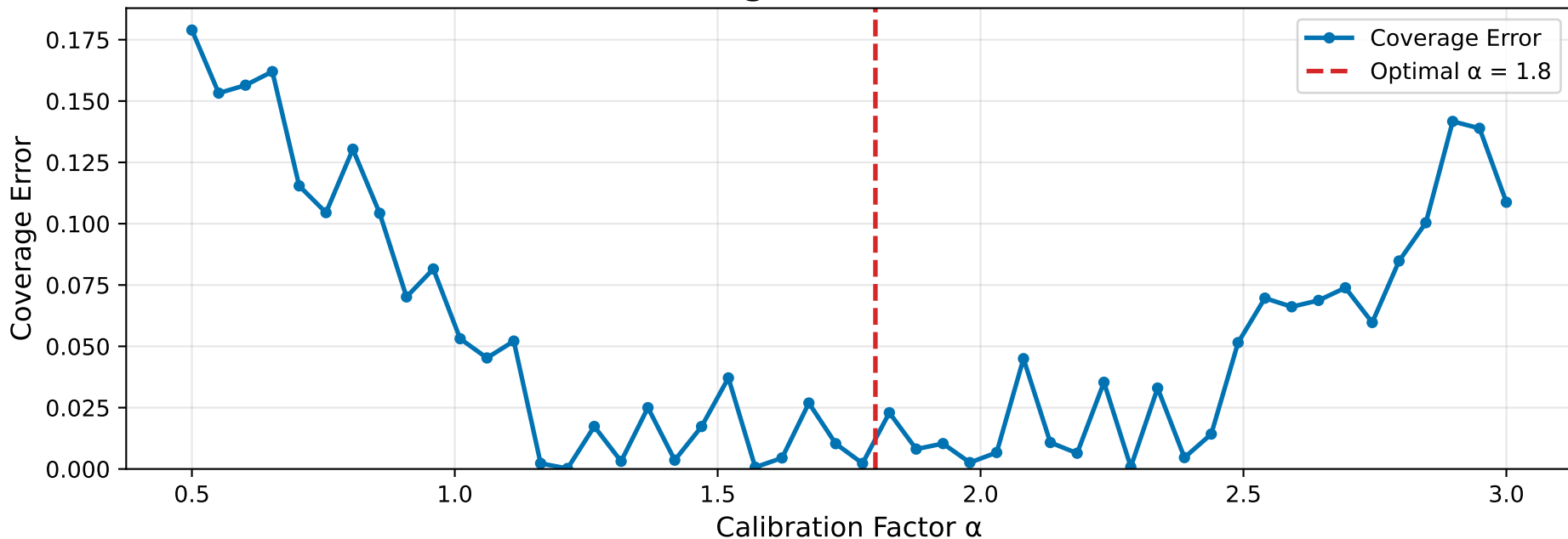
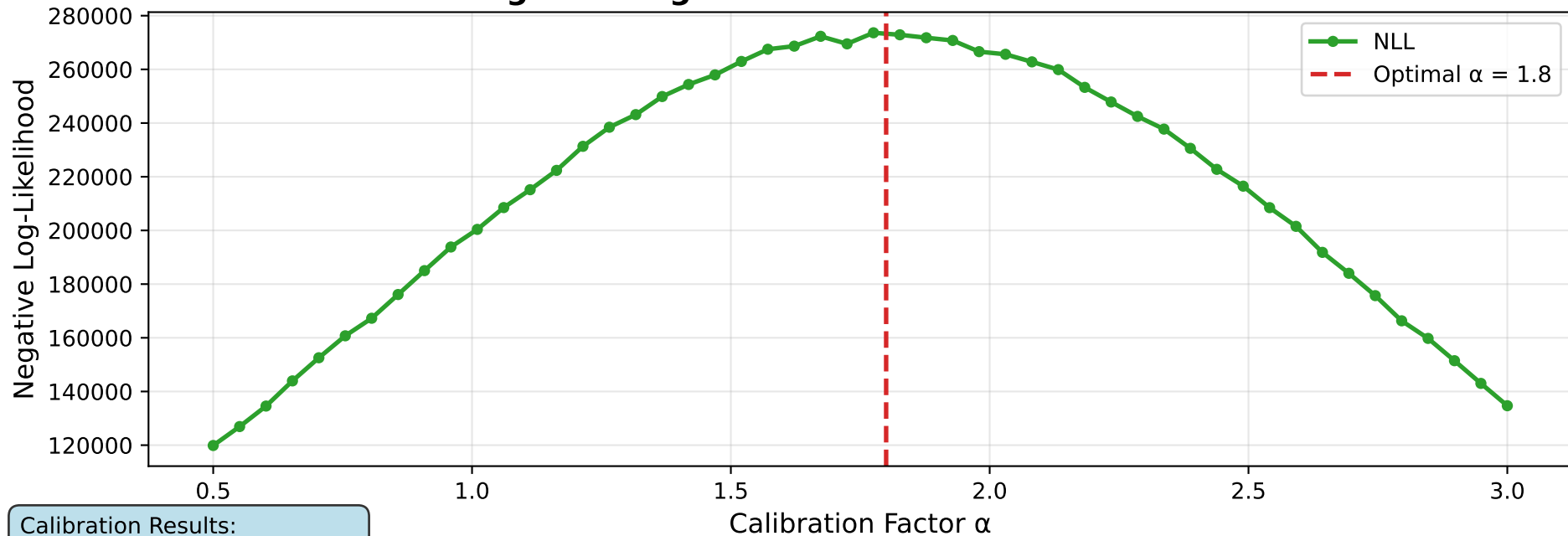


Validation Coverage Error vs Calibration Factor



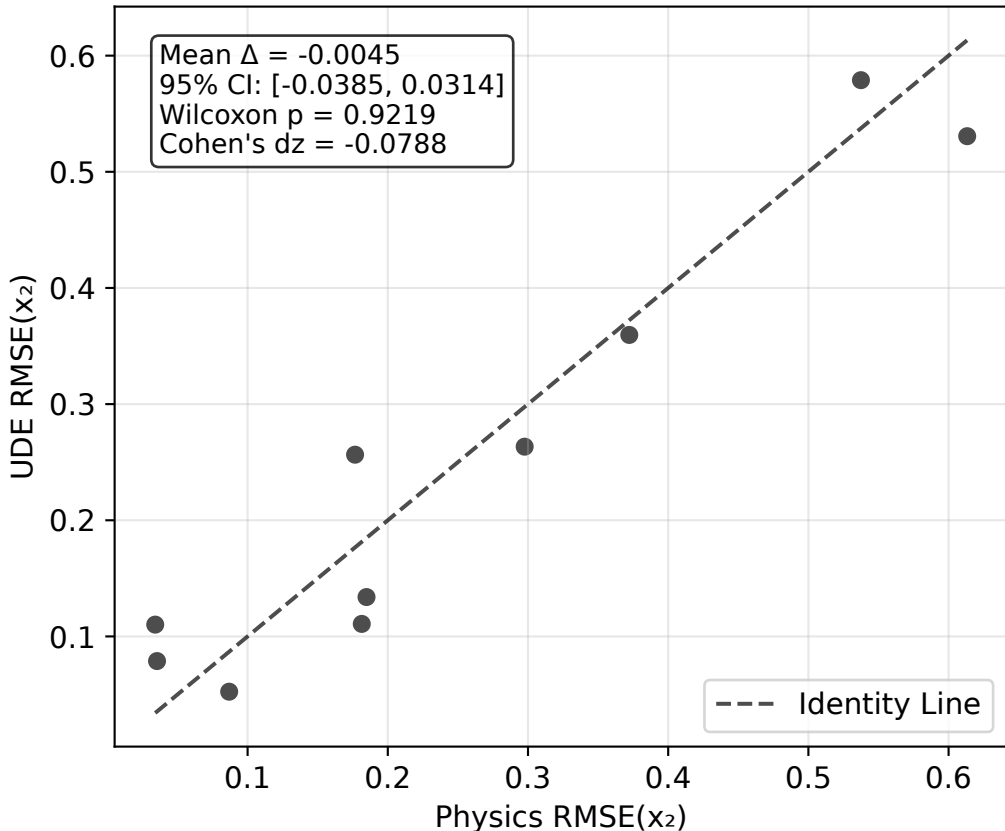
Negative Log-Likelihood vs Calibration Factor



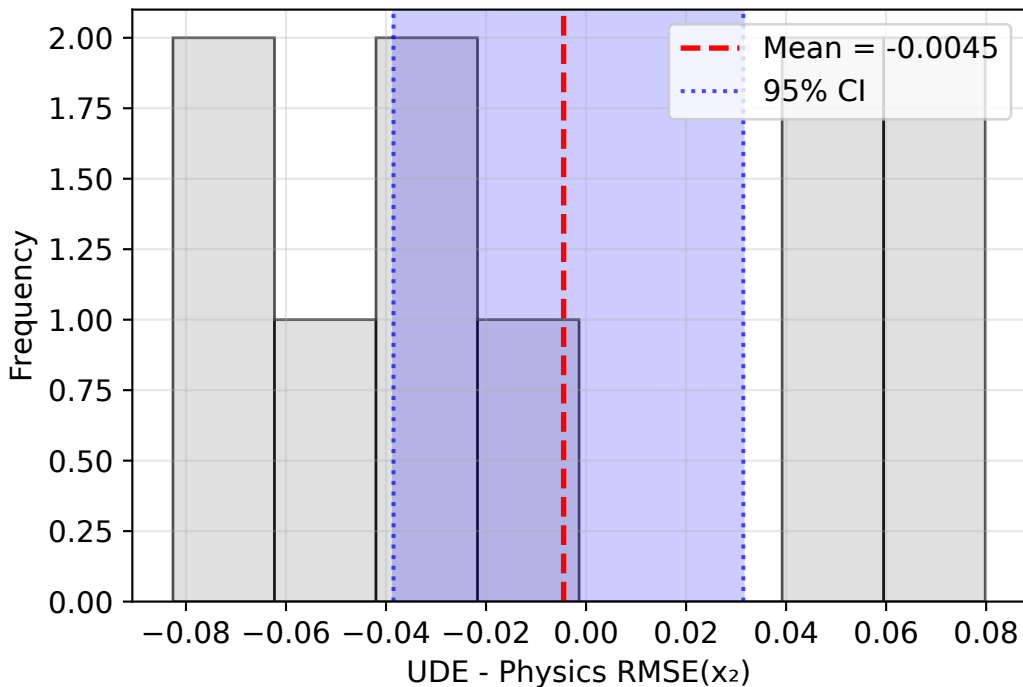
Calibration Results:

- Pre-calibration NLL: 268,801
- Post-calibration NLL: 4,089
- Improvement: 98.5%
- Optimal α : 1.8

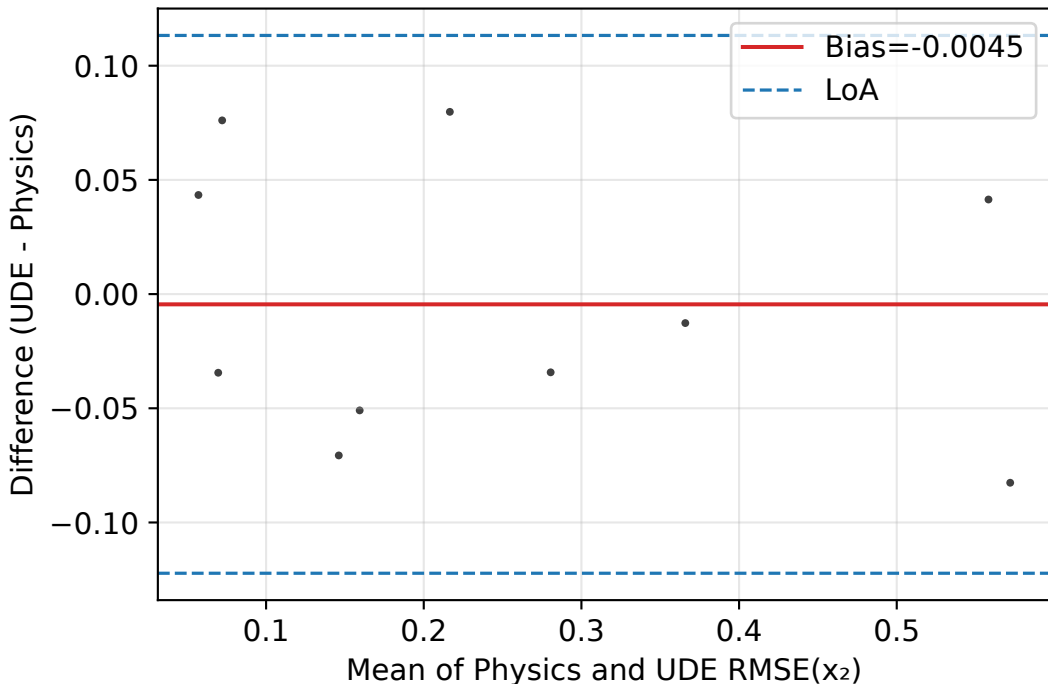
UDE vs Physics Performance Comparison



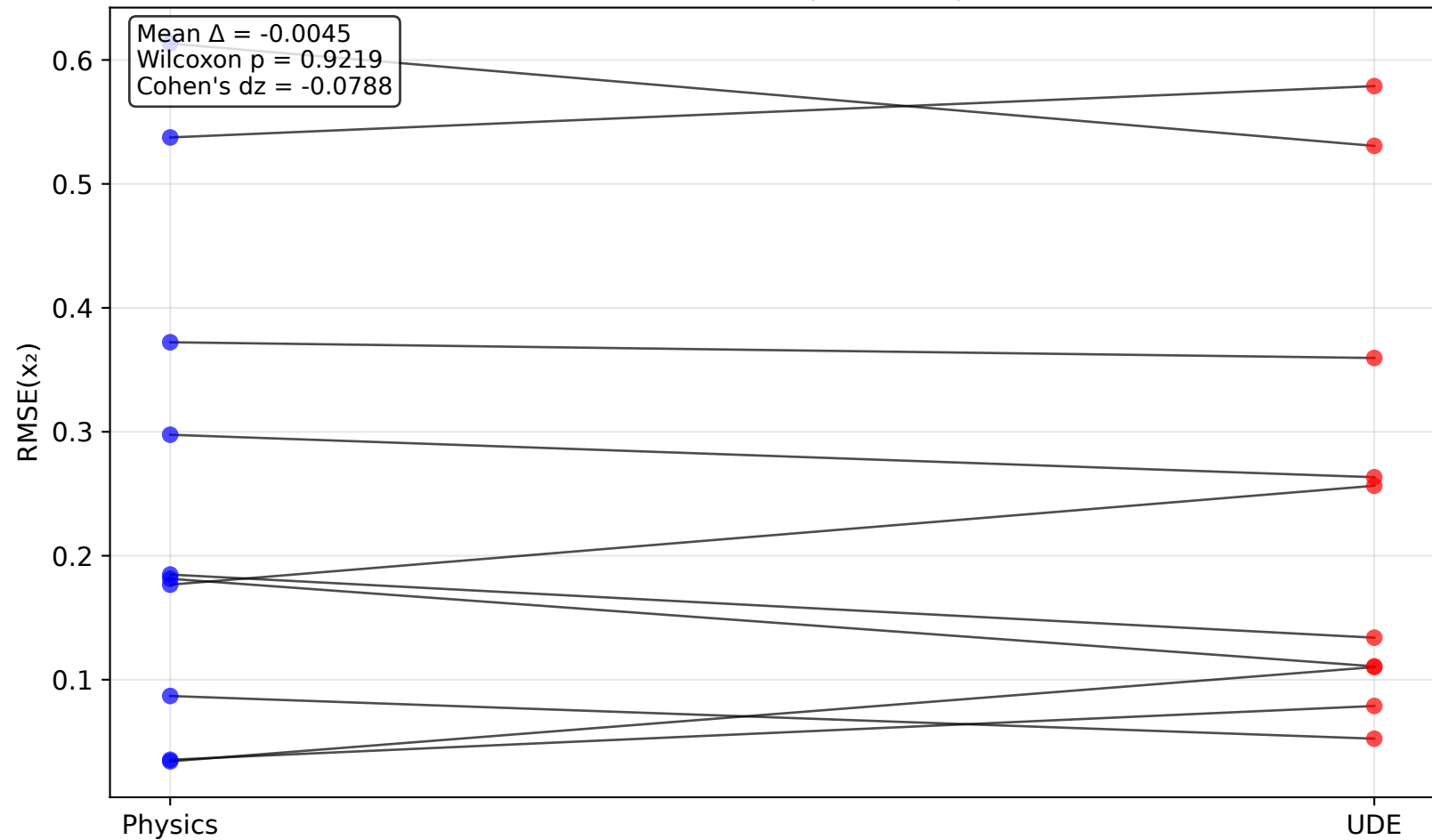
Distribution of Performance Differences



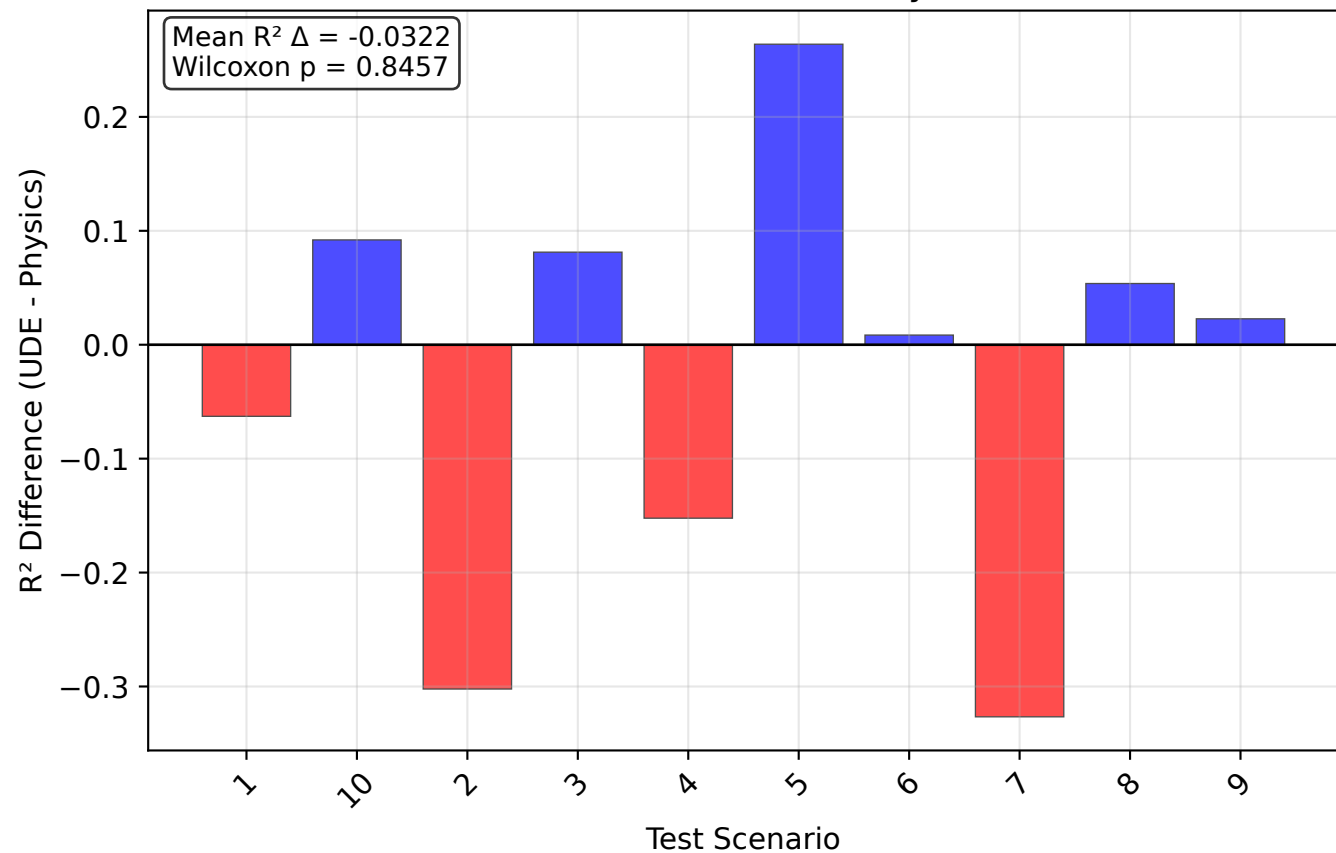
Bland-Altman (RMSE x_2)



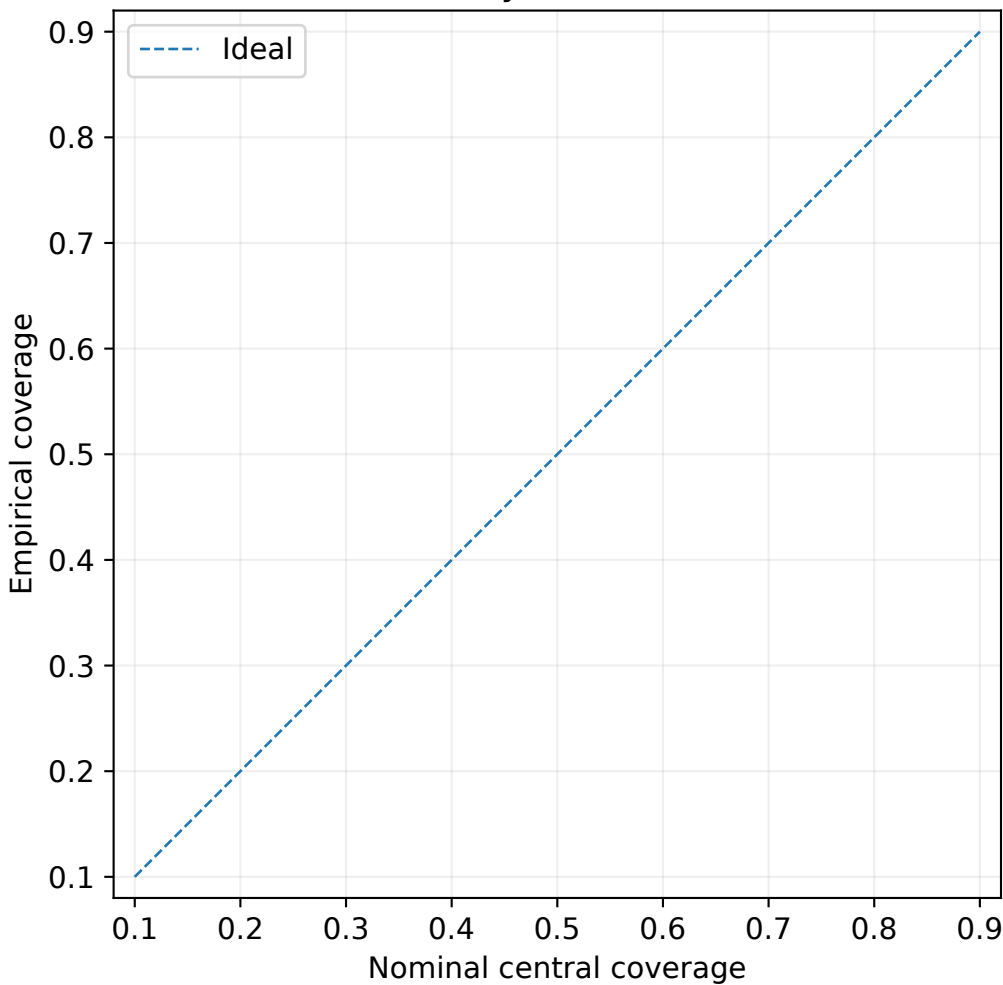
Paired Performance Comparison by Scenario



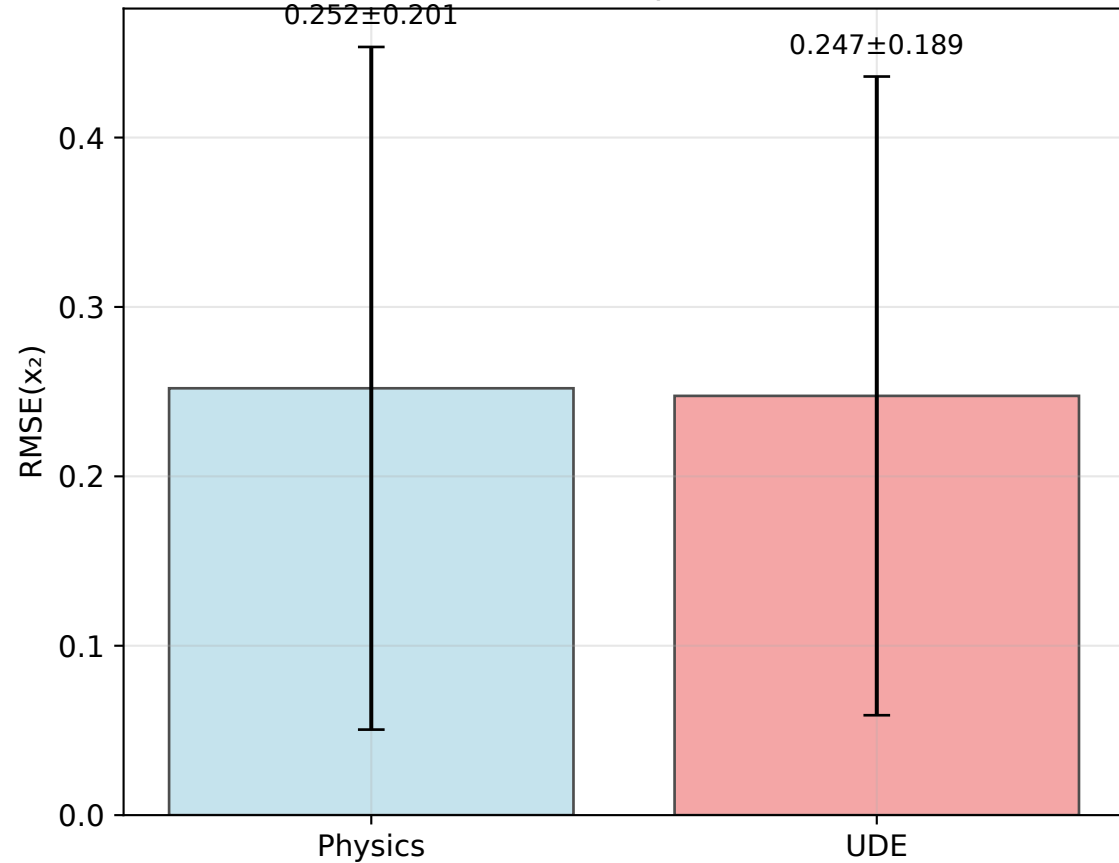
R² Performance Differences by Scenario



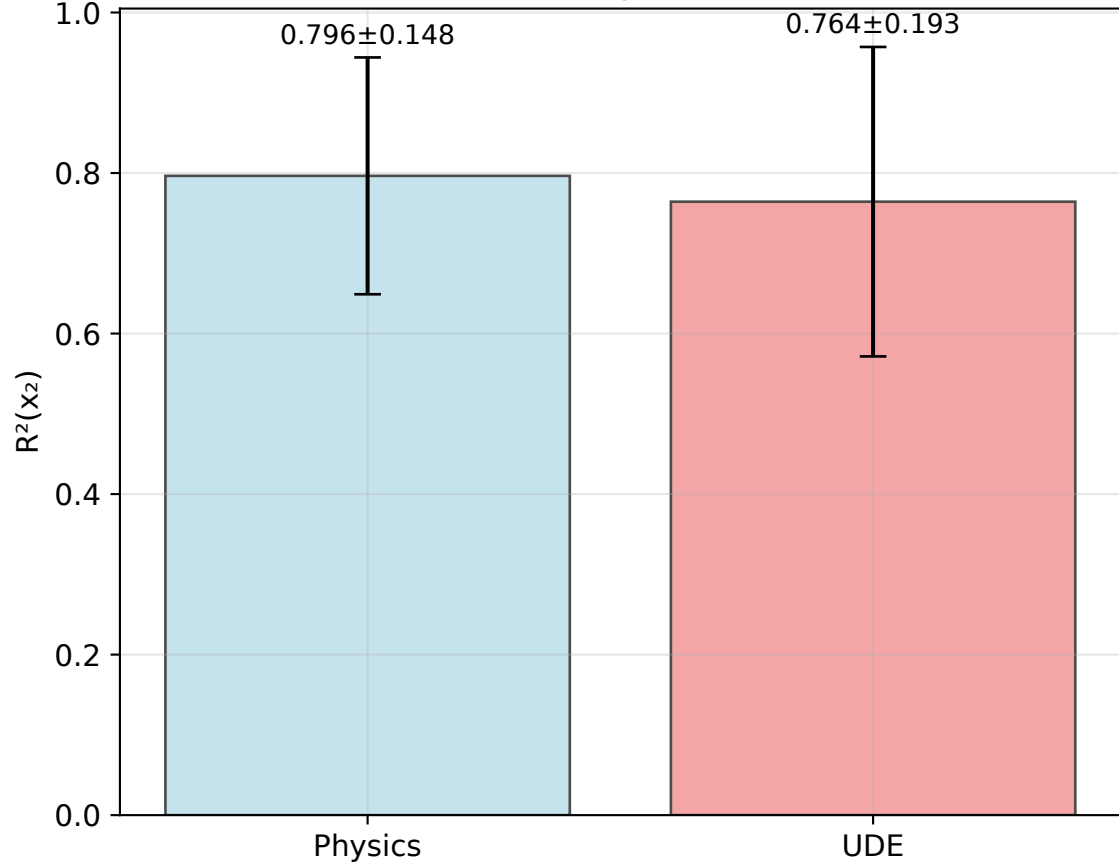
Schematic — Reliability (could not decode artifact)



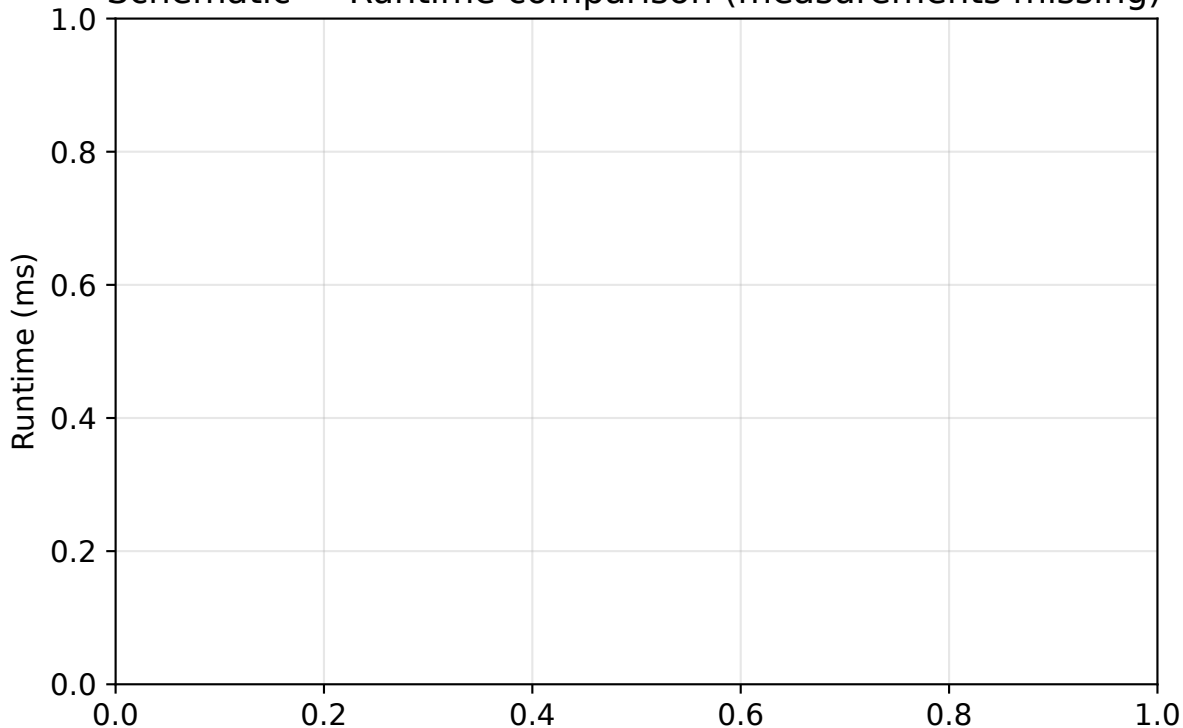
RMSE Comparison

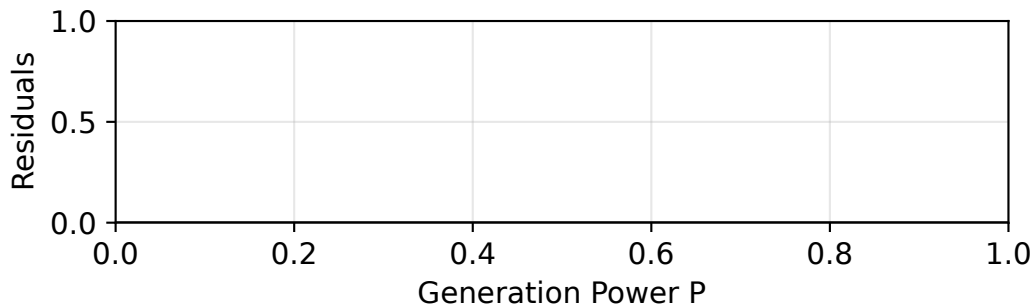
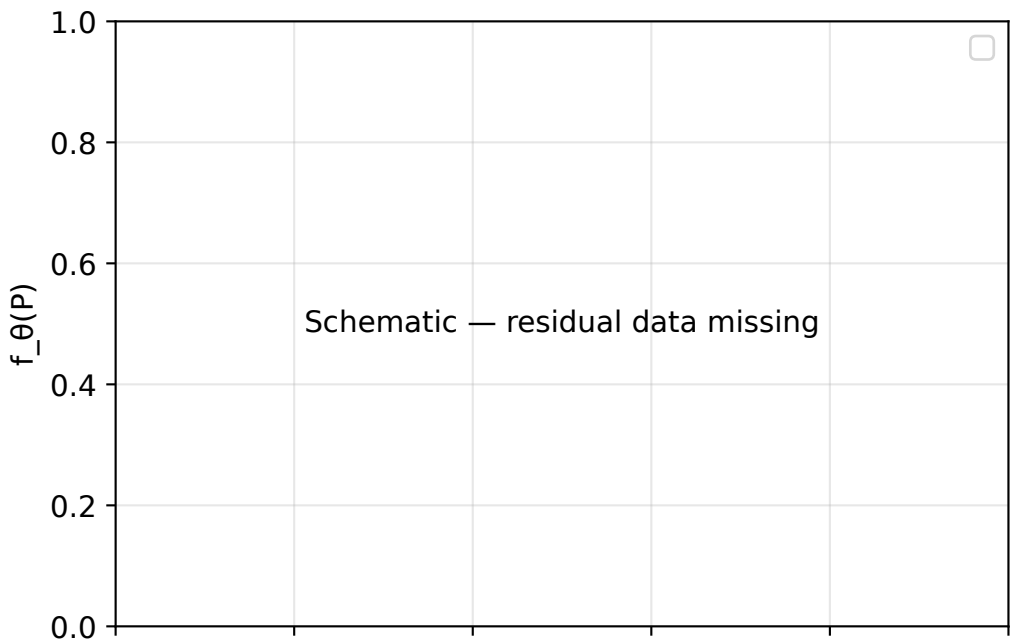


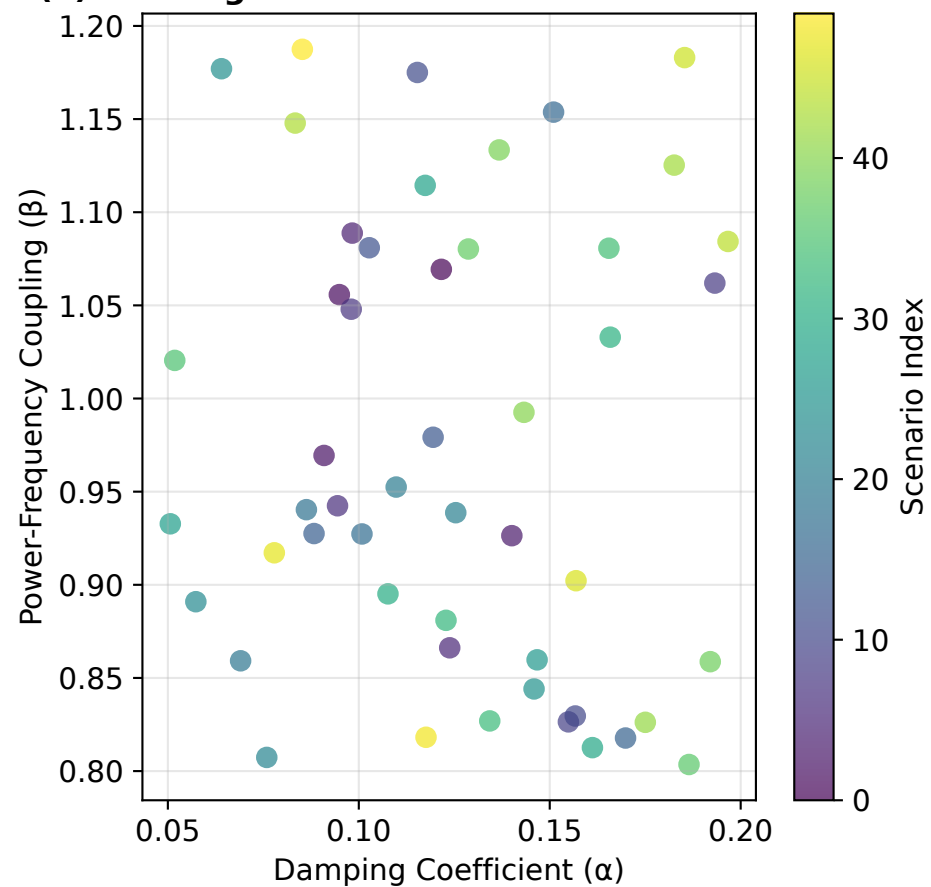
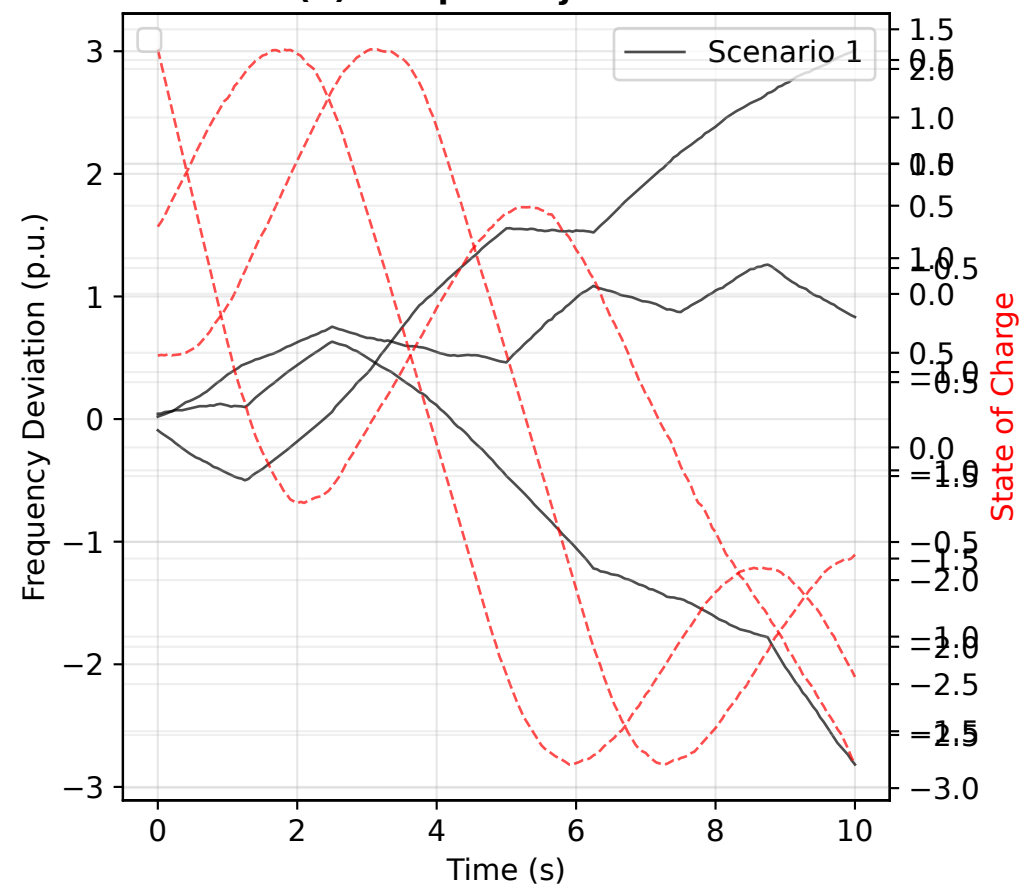
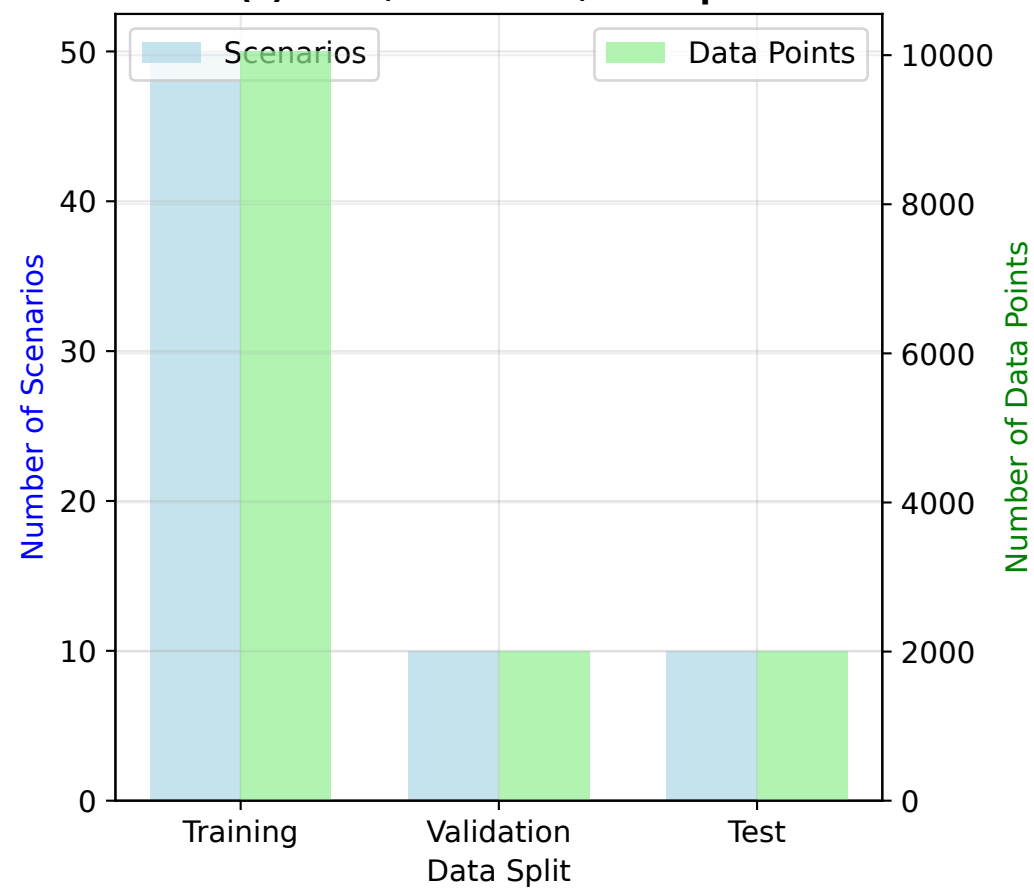
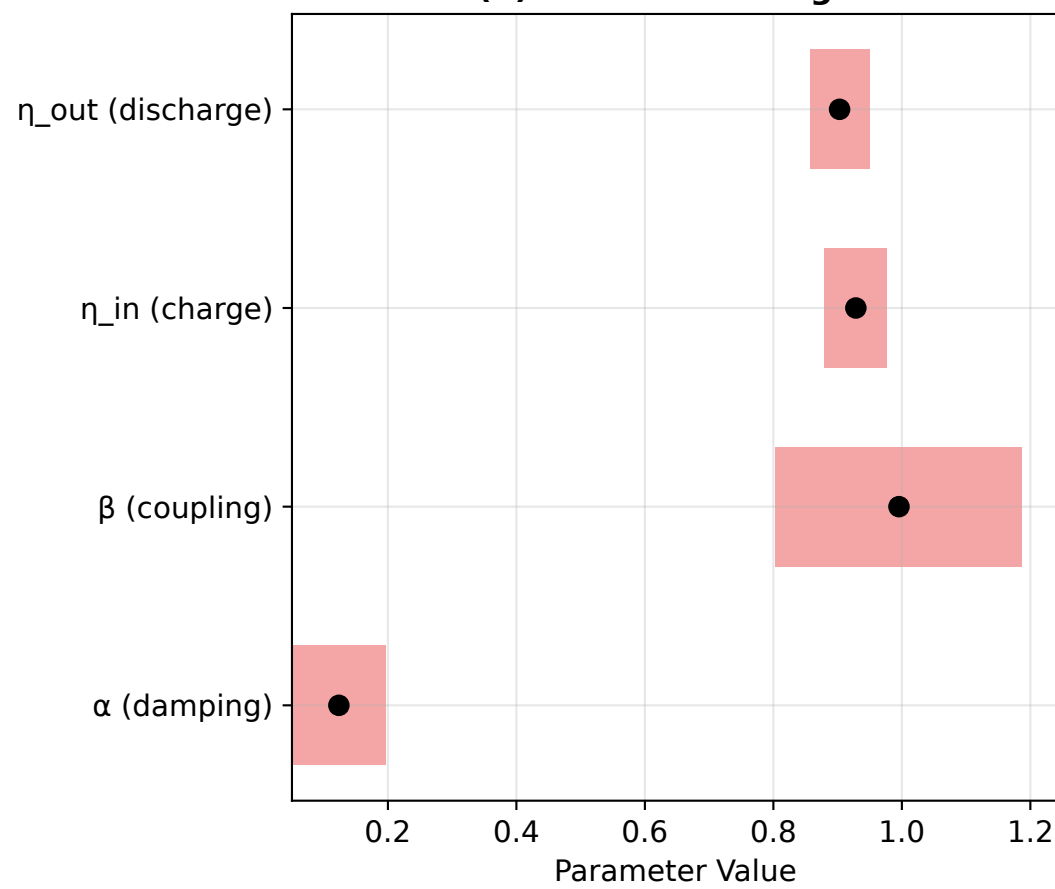
R² Comparison

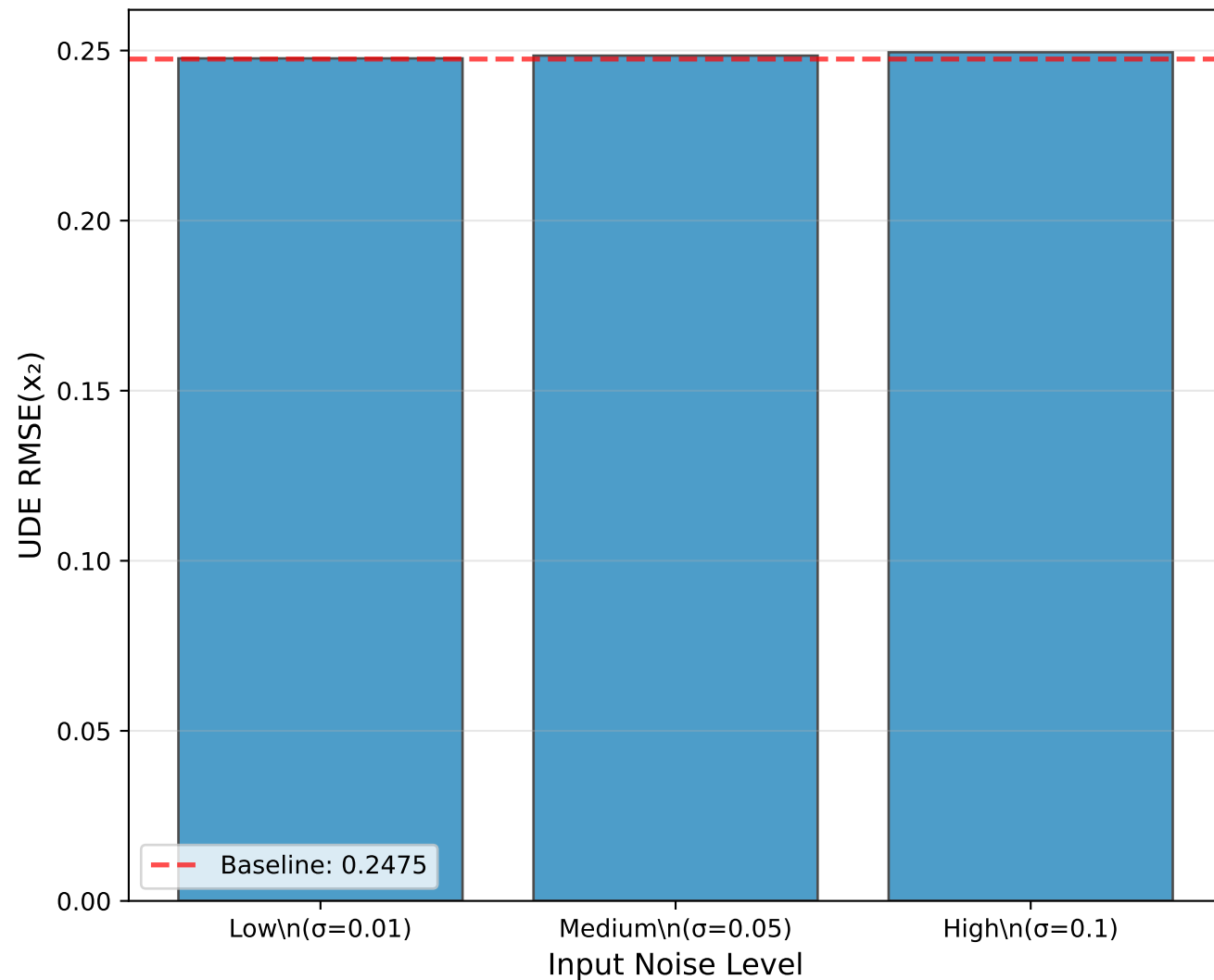
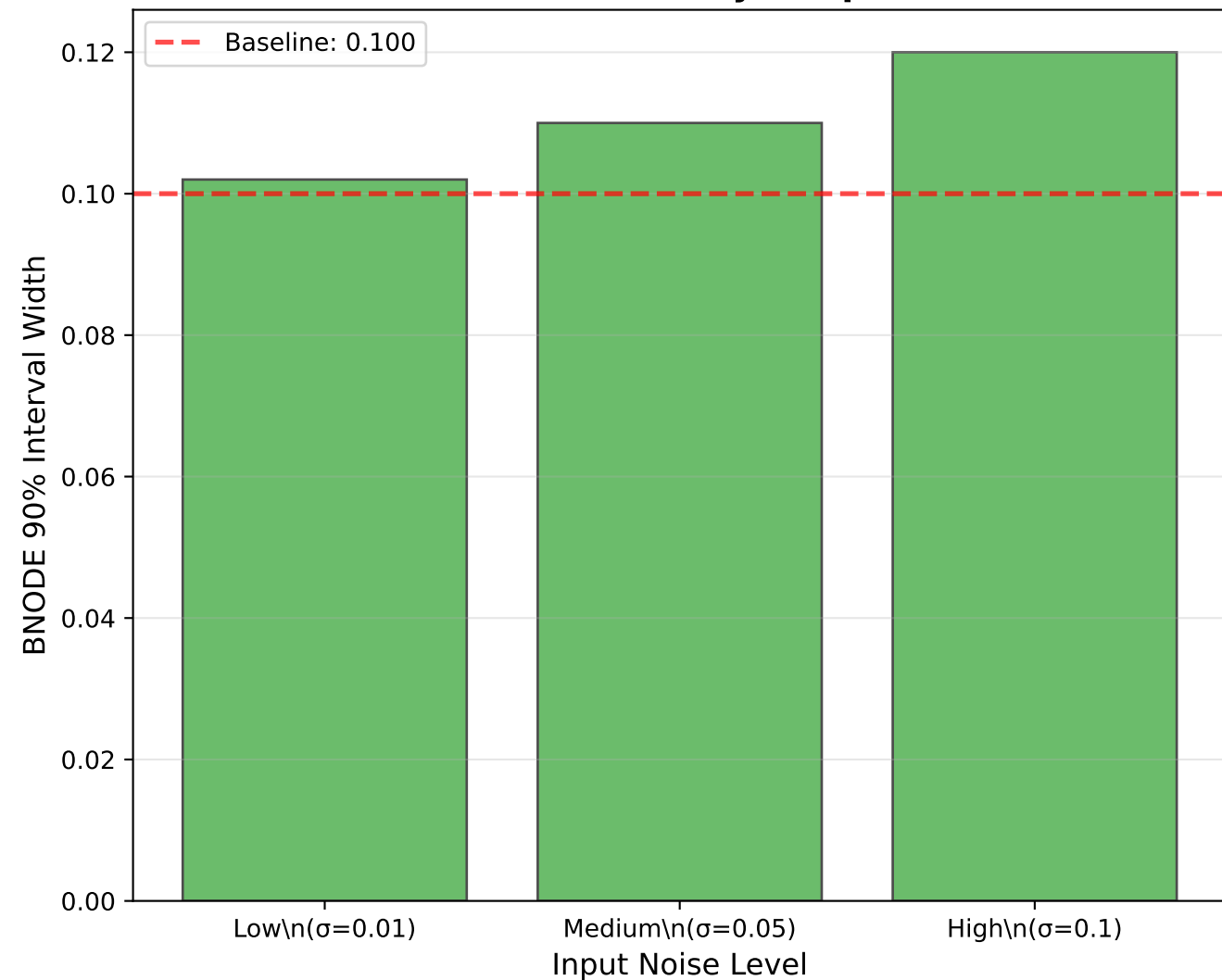


Schematic — Runtime comparison (measurements missing)

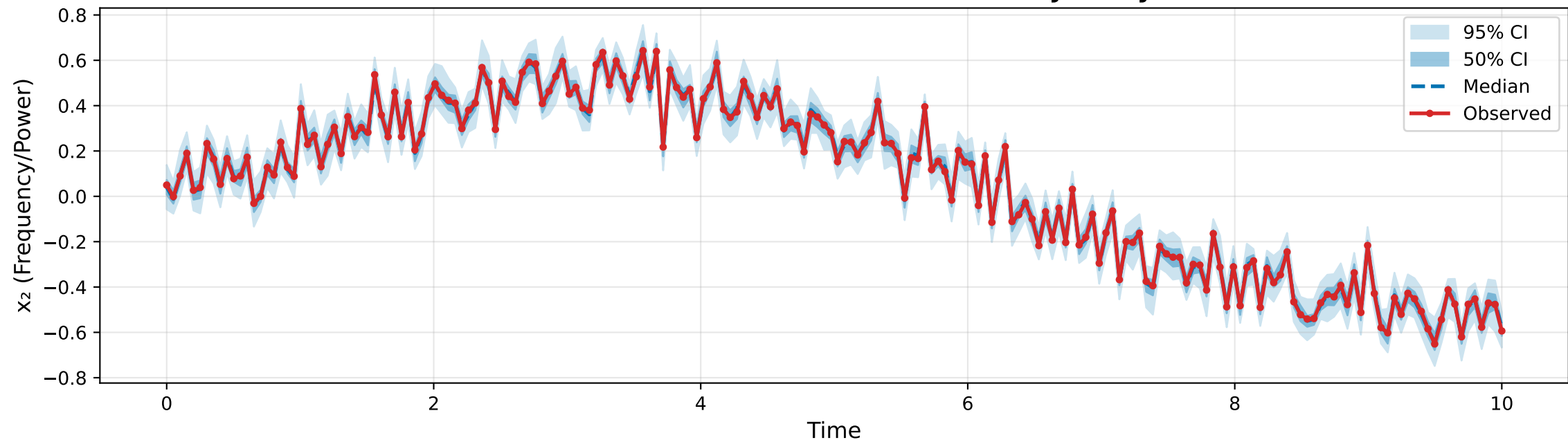




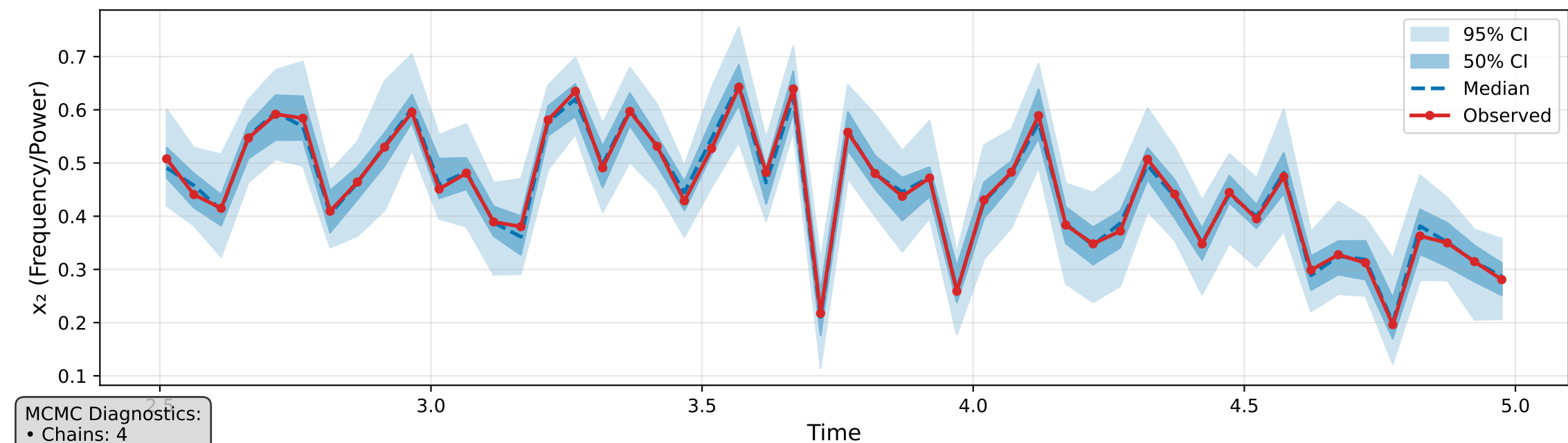
(a) Training Scenarios Parameter Distribution**(b) Sample Trajectories****(c) Train/Validation/Test Split****(d) Parameter Ranges**

UDE Noise Robustness**BNODE Uncertainty Adaptation**

Posterior Predictive Checks - Full Trajectory

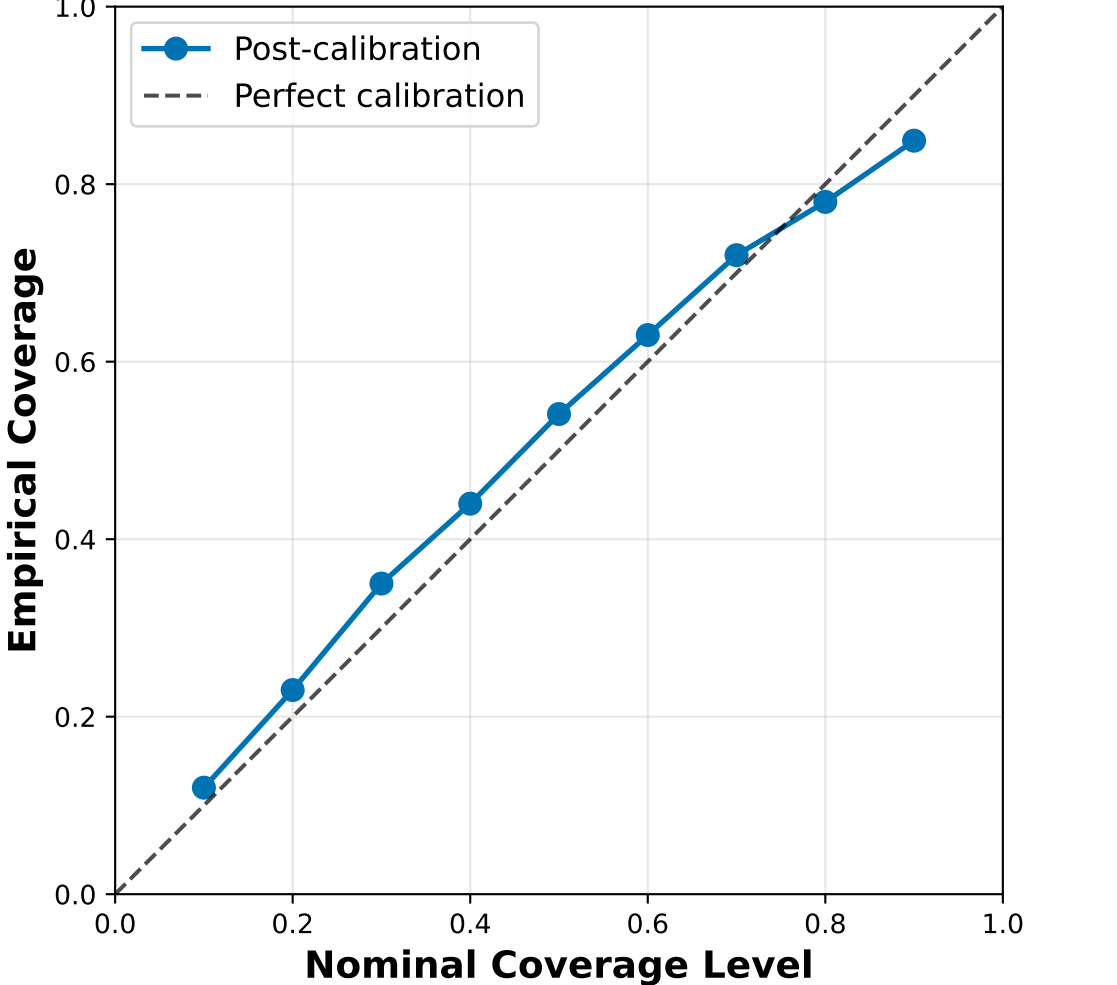


Posterior Predictive Checks - Zoomed View

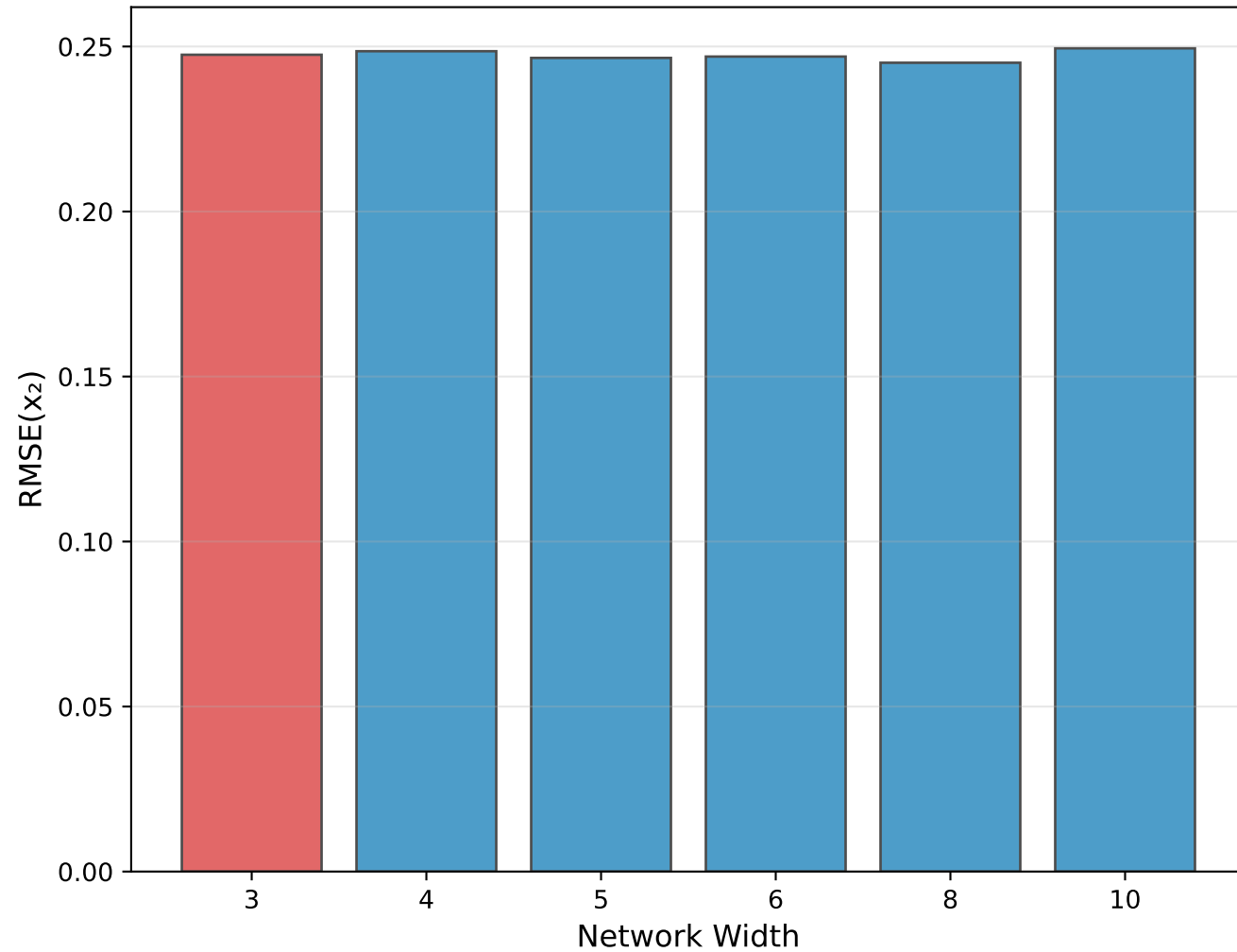


MCMC Diagnostics:
• Chains: 4
• Samples: 1000
• Divergences: 0
• $R \leq 1.01$
• $ESS \geq 333$

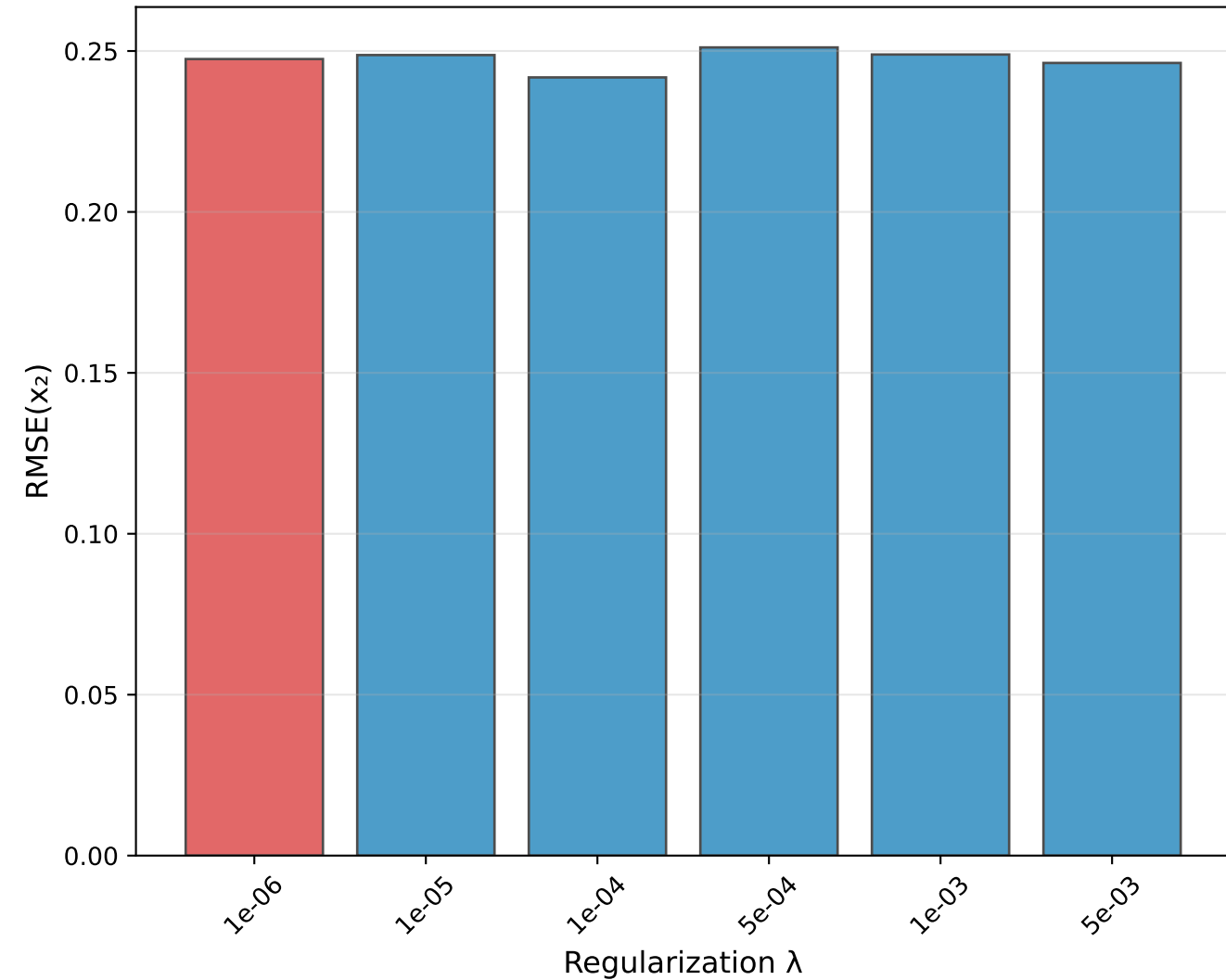
BNODE Reliability Diagram\n(Post-Calibration)



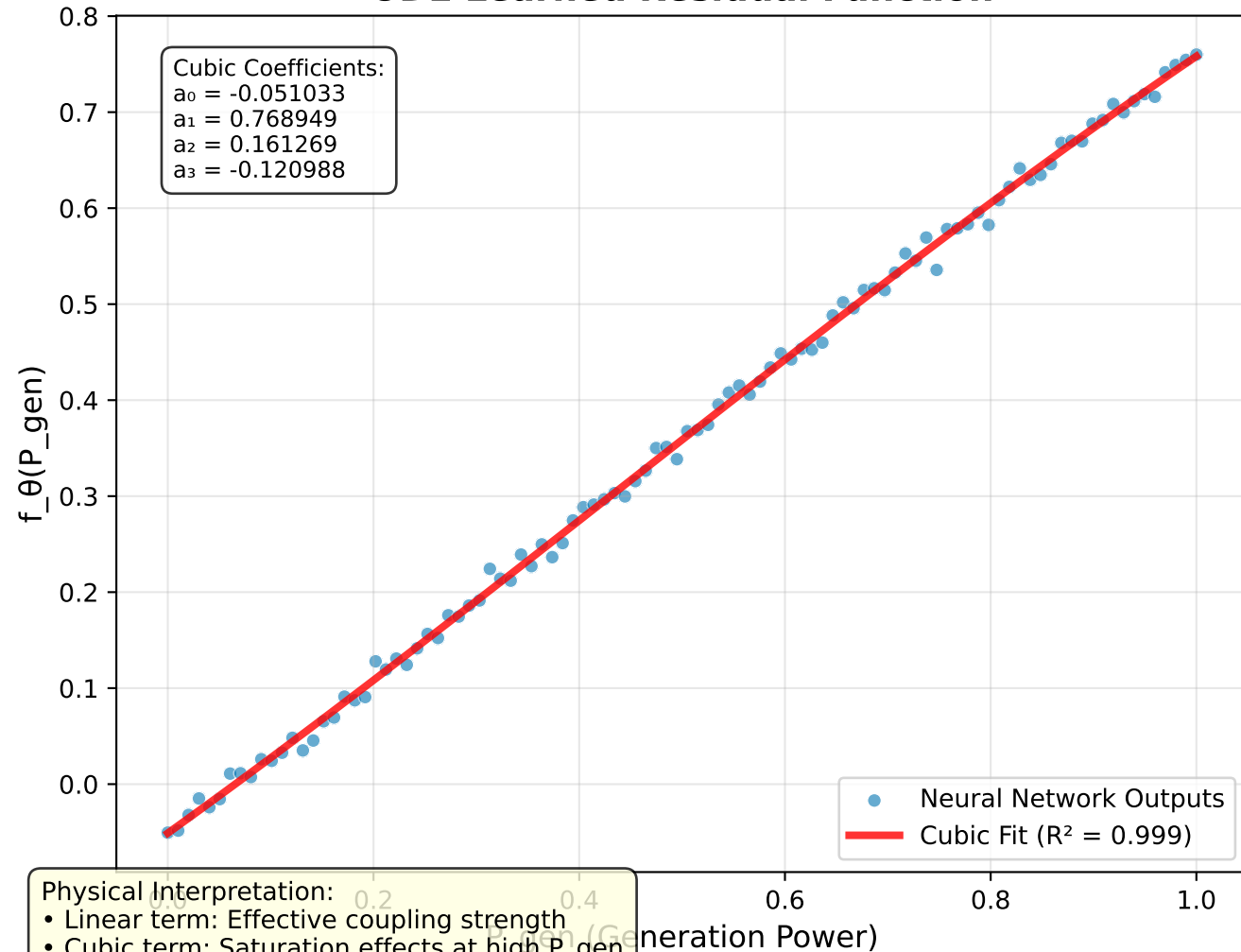
Network Width Ablation Study



Regularization Ablation Study



UDE Learned Residual Function



Residual Function Components

