

We will have an opportunity for hands-on tool use!

To prepare:

1. Download RStudio: <https://posit.co/downloads/>
2. Download Shiny apps: <https://github.com/krometis/dataworks2026>
3. Install dependencies





ACQUISITION INNOVATION
RESEARCH CENTER



RECENT METHODOLOGICAL ADVANCES FOR INTEGRATED T&E

MINI-TUTORIAL, DATAWORKS 2026

Dr. Justin Krometis
April 22, 2026

The views, findings, conclusions, and recommendations expressed in this material are solely those of the author(s) and do not necessarily reflect the views or positions of the United States Government (including the Department of Defense (DoD) and any government personnel)

1. Background: What and Why
2. Methods
3. Case Study 1: Bayesian Reliability
4. Case Study 2: Changing Factors
5. Case Study 3: Test Planning
6. Tool Development
7. Future Outlook

BACKGROUND



ACQUISITION INNOVATION
RESEARCH CENTER



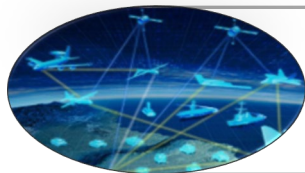
ARTIFICIAL INTELLIGENCE



INTEGRATED TESTING



DIGITAL ENGINEERING



KILL WEBS



Research Partners

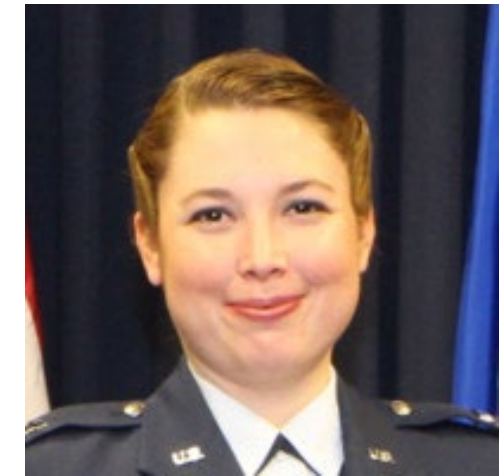




Dr. Laura Freeman
Virginia Tech



Corinne Stafford
STAT COE

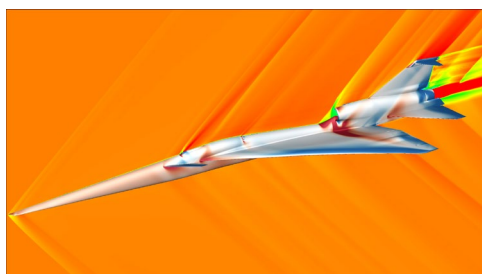


Dr. Victoria Sieck
U.S. Air Force

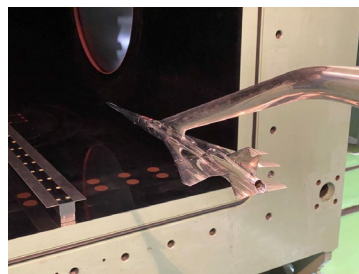
Using iterative sequential-progressive design of experiments coupled with evaluations that combine information from across the acquisition lifecycle.



Historical Data
SME Knowledge



Model Data



DT Data



OT Data



Operational
Performance

Leveraging all data enables better understanding of systems *earlier*...



...allowing for fewer or more optimal tests *later*



Alternative T&E Pathway (805(b)):

"(1)(A) design and execute a unified T&E strategy that **aligns DT and OT** to ... **build system understanding throughout the test program**"

"(1)(B)(v) automated analytics tools to **assess performance trends**, reliability, and maintenance needs"

"(1)(C)(iv) **integration of supporting or complementary data** from digital twins or other model-based systems engineering tools"

"(1)(D) define general test and evaluation objectives and data needs while **allowing detailed execution plans to evolve based on test results** and emerging requirements"

Role of DOT&E (805(c)):

"(1) provide independent evaluation of test data **across all phases of the program lifecycle**"

"(1)(B) evaluating whether the program collects and analyzes sufficient raw data, learns from test results at a **pace relevant to operational needs**, and converges on military effectiveness **based on data trends**"

"(1)(D) providing **continuous oversight through ongoing analysis** of test data"

Maximize Acquisition Flexibility:

"It is critical that teams leverage existing authoritative **data sources, including contractor data** and automated reporting mechanisms to assess program performance."

"scaling tools, policies, and practices to **maximize the use of modeling and simulation and automated testing** will allow the Department to accelerate and continually validate software."

"It is important that the Department **invest in test and evaluation (T&E) resources for validating digital technologies, including modeling and simulation, synthetic environments,** and managing performance/risks."

"the Department can **prioritize integrated testing to optimize test resources, avoid redundancies,** and utilize data-driven decisions to trace programs, performance, and testing"

Develop High-Performing Systems:

"The Department will integrate and scale adoption and investments in ...modeling and simulation environments for **virtual and constructive testing and training, and developmental integration and testing,** which will fuel rapid, iterative designs and technology insertion to maximize mission outcomes. These advanced techniques will **reduce the burden in test planning and execution** and the unplanned testing that results from unanticipated discovery"

"Integrate and **test new technology soonest to get key insights early** to shape weapon system development and fielding"

3.1(a):

“OT&E and LFT&E planning, execution, analysis, and reporting activities will use **the latest advances in science (e.g., design of experiments, statistical inference methods, or big data analytics)** to ... determine, with scientific rigor, the preliminary and final operational effectiveness, suitability, survivability, and lethality (as applicable) of DoD systems.”

3.1(c):

“Science and technology-based OT&E and LFT&E will **enable efficient use of data from multiple data sources** (e.g., contractor test (CT), developmental test (DT), operational test (OT), and live fire test (LFT) data or M&S results). Improved **sequential testing using Bayesian or similar inference methods** ... are critical to dynamically optimize the planning, execution, analysis, and reporting of integrated T&E, OT&E, and LFT&E across the acquisition life cycle.”

Consideration	Possible Levels (best case → worst case)
Safety	Minimal risk (not live projectile) → High risk (live fire)
Cost	\$ → \$\$\$
Resource Availability	Available → Partially available → Needs to be developed
Schedule	Easy & quick → Hard & extensive coordination
Historical operational performance data	Yes same factors → Yes but missing key factor(s) → None
Modeling and Simulation	Yes accurate and validated over time → Yes but not well understood and/or missing key factor(s) → None
Scale	Single component → Parallel systems → Series system

Consideration	Possible Levels (best case → worst case)
Safety	Minimal risk (not live projectile) → High risk (live fire)
Cost	\$ → \$\$\$
Resource Availability	Available → Partially available → Needs to be developed
Schedule	Easy & quick → Hard & extensive coordination
Historical operational performance data	Yes same factors →
Modeling and Simulation	Yes accurate and valid and/or missing key factors
Scale	Single component →

Variety of DoD programs means that a variety of analysis approaches may be appropriate!

Comprehensive Evaluations

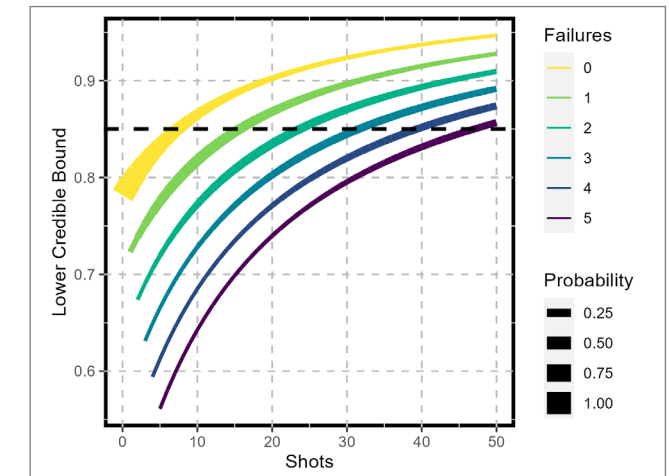
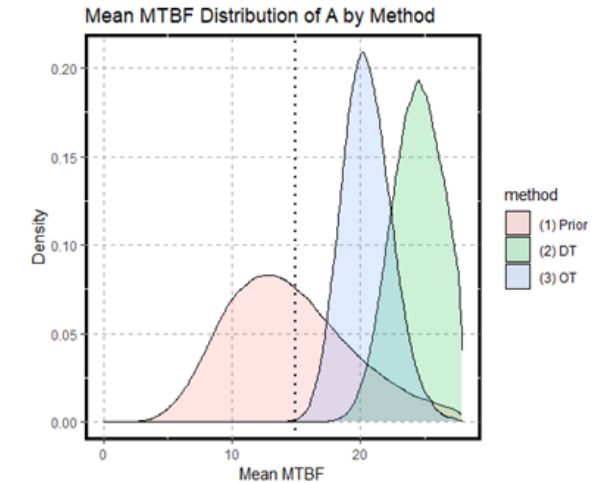
Estimating operational performance from data...

- ...from multiple phases of test
- ...of different fidelities/sizes
- ...with different experimental factors
- ...from digital representations of systems

Test Planning/Design

Given previous test results...

- ...where and how many tests should we run?
- ...when have we tested enough?



METHODS



ACQUISITION INNOVATION
RESEARCH CENTER

Fit a statistical model to all of the data

Naïve: Assume all data is equivalent and fit to all data equally

Blocking: Try to account for differences in data sources by adding source or phase-specific factors to the model

- Example: Add a shift parameter to account for possible biases in data sources

Relates the probability of a parameter value θ given data Y ($P(\theta|Y)$) to the probability of Y given θ and the probability of θ :

$$P(\theta|Y) \propto P(Y|\theta) P(\theta)$$

Relates the probability of a parameter value θ given data Y ($P(\theta|Y)$) to the probability of Y given θ and the probability of θ :

$$P(\theta|Y) \propto P(Y|\theta) \boxed{P(\theta)}$$

Prior

Relates the probability of a parameter value θ given data Y ($P(\theta|Y)$) to the probability of Y given θ and the probability of θ :

$$P(\theta|Y) \propto \underbrace{P(Y|\theta)}_{\text{Likelihood (Data)}} \underbrace{P(\theta)}_{\text{Prior}}$$

Relates the probability of a parameter value θ given data Y ($P(\theta|Y)$) to the probability of Y given θ and the probability of θ :

$$P(\theta|Y) \propto P(Y|\theta)P(\theta)$$

Posterior

Likelihood
(Data)

Prior

Relates the probability of a parameter value θ given data Y ($P(\theta|Y)$) to the probability of Y given θ and the probability of θ :

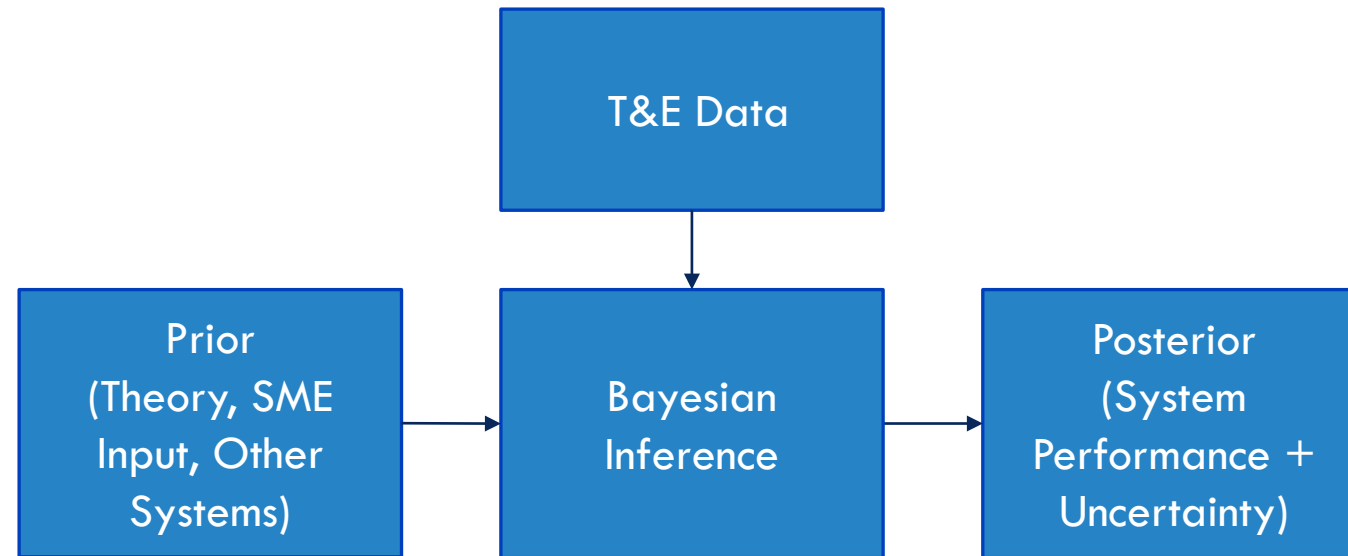
$$P(\theta|Y) \propto P(Y|\theta)P(\theta)$$

Posterior

Likelihood
(Data)

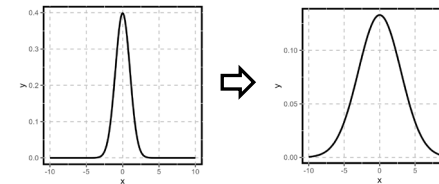
Prior

Models human learning: Understanding (prior) + Experience (likelihood) = Updated understanding (posterior)



SINGLE TEST PHASE

T&E Data



Prior
(Theory, SME
Input, Other
Systems)

Bayesian
Inference

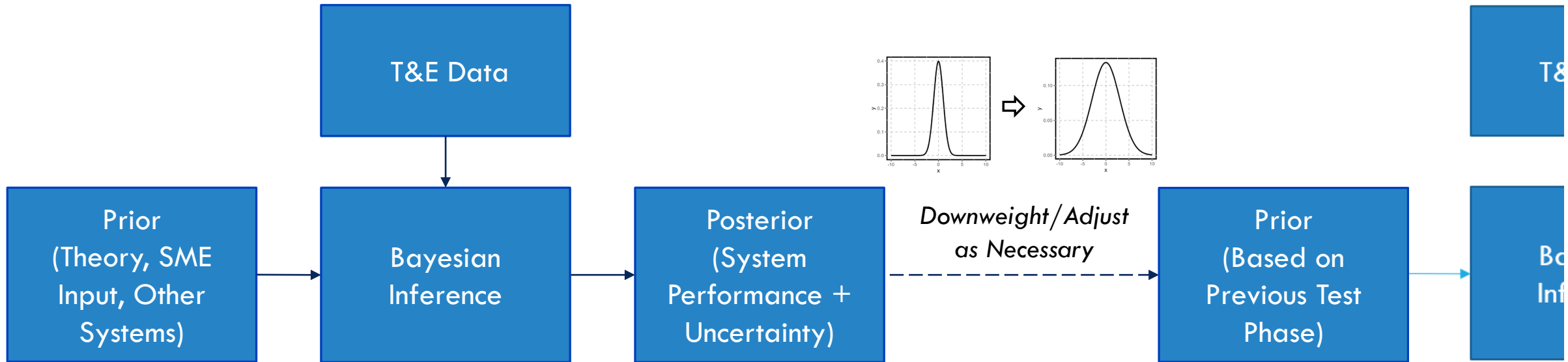
Posterior
(System
Performance +
Uncertainty)

*Downweight/Adjust
as Necessary*

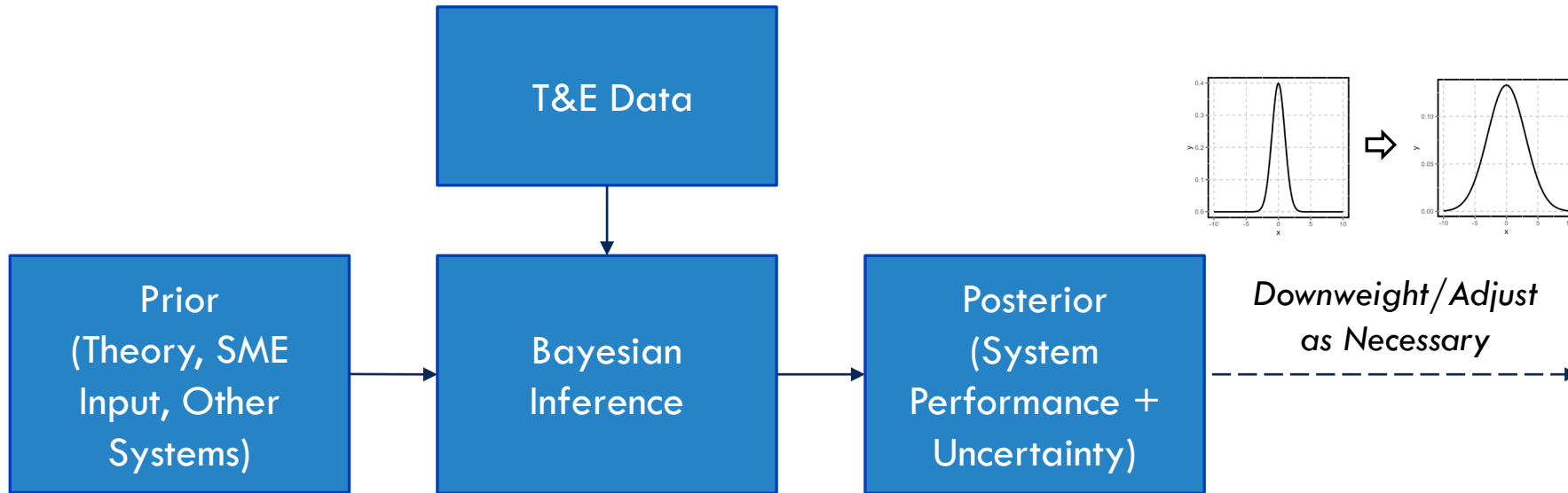
Prior
(Based on
Previous Test
Phase)

SINGLE TEST PHASE

NEXT TEST PHASE



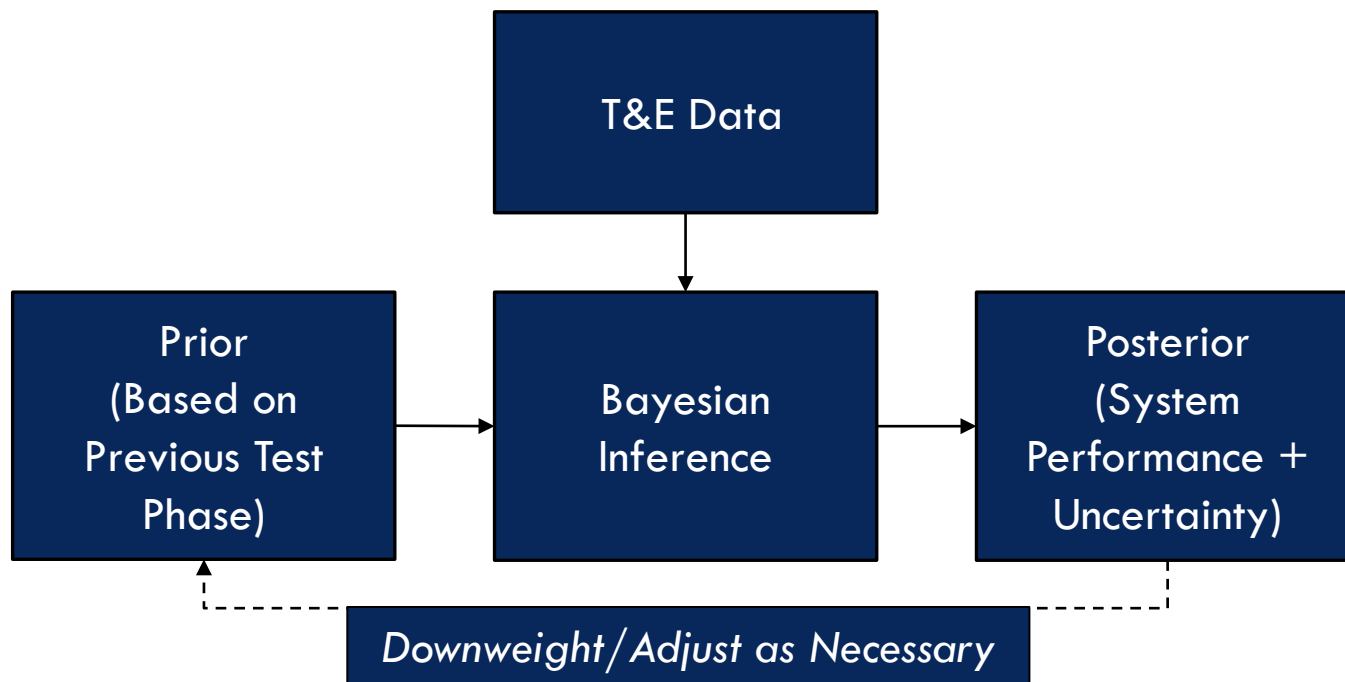
SINGLE TEST PHASE

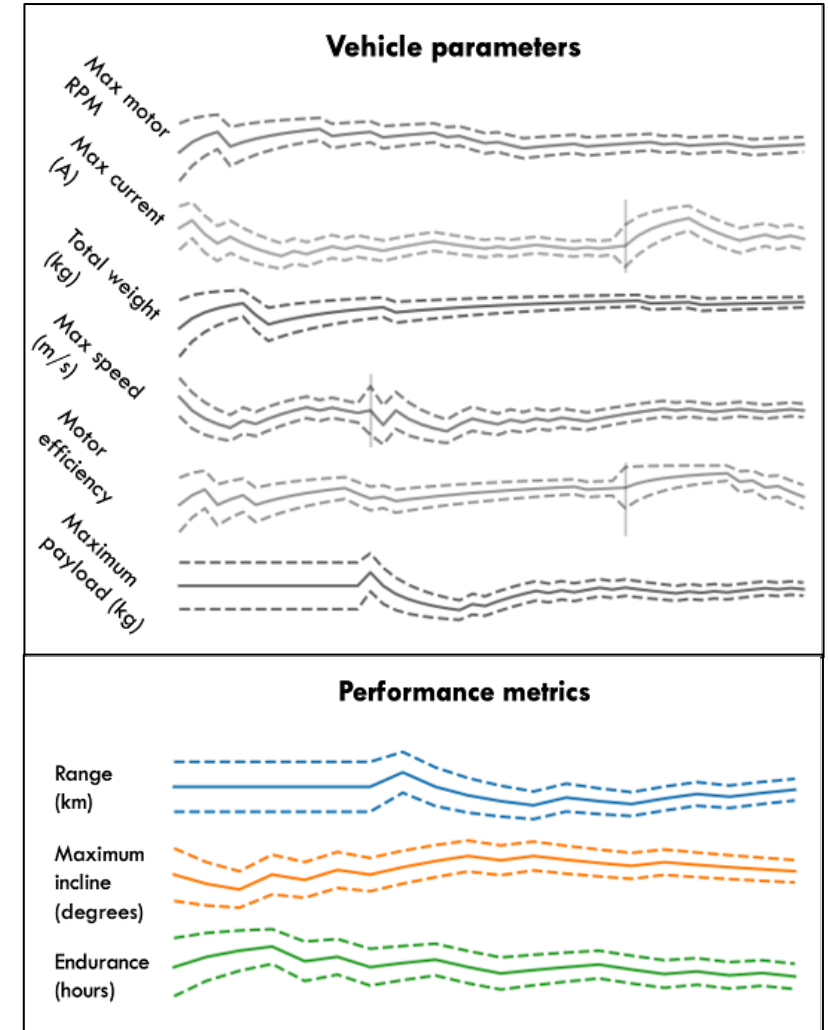
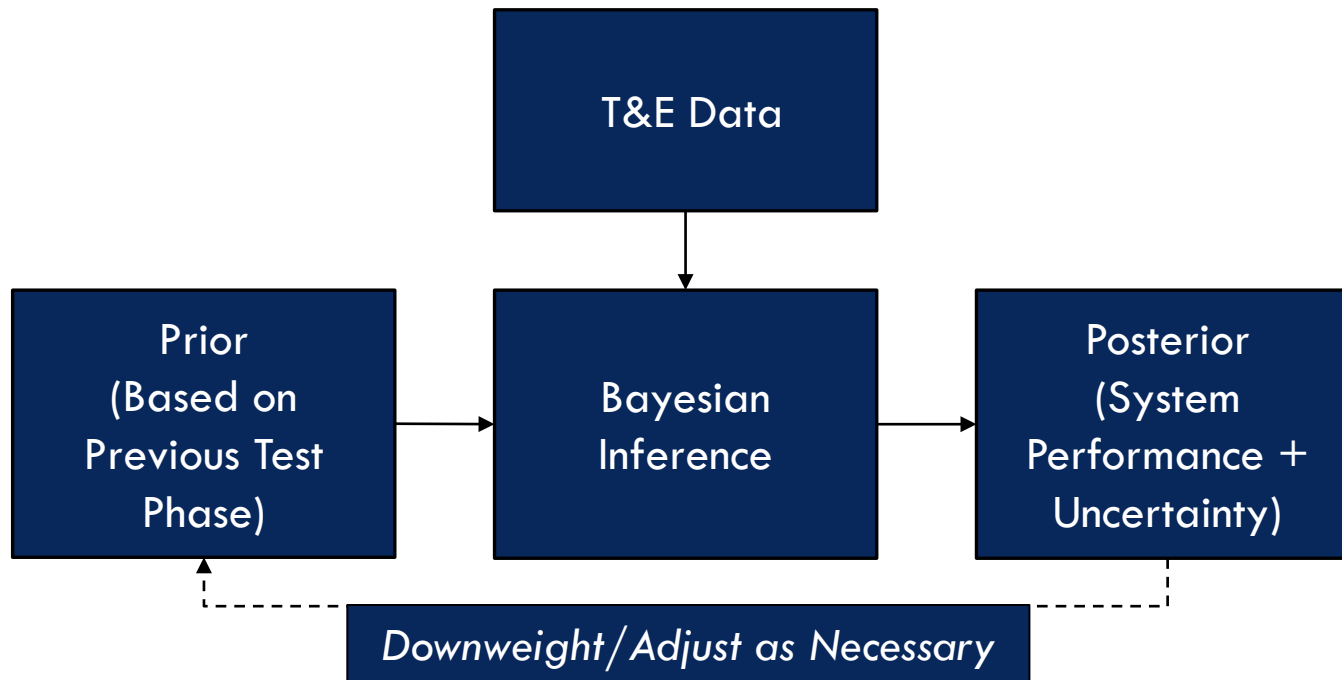


NEXT TEST PHASE

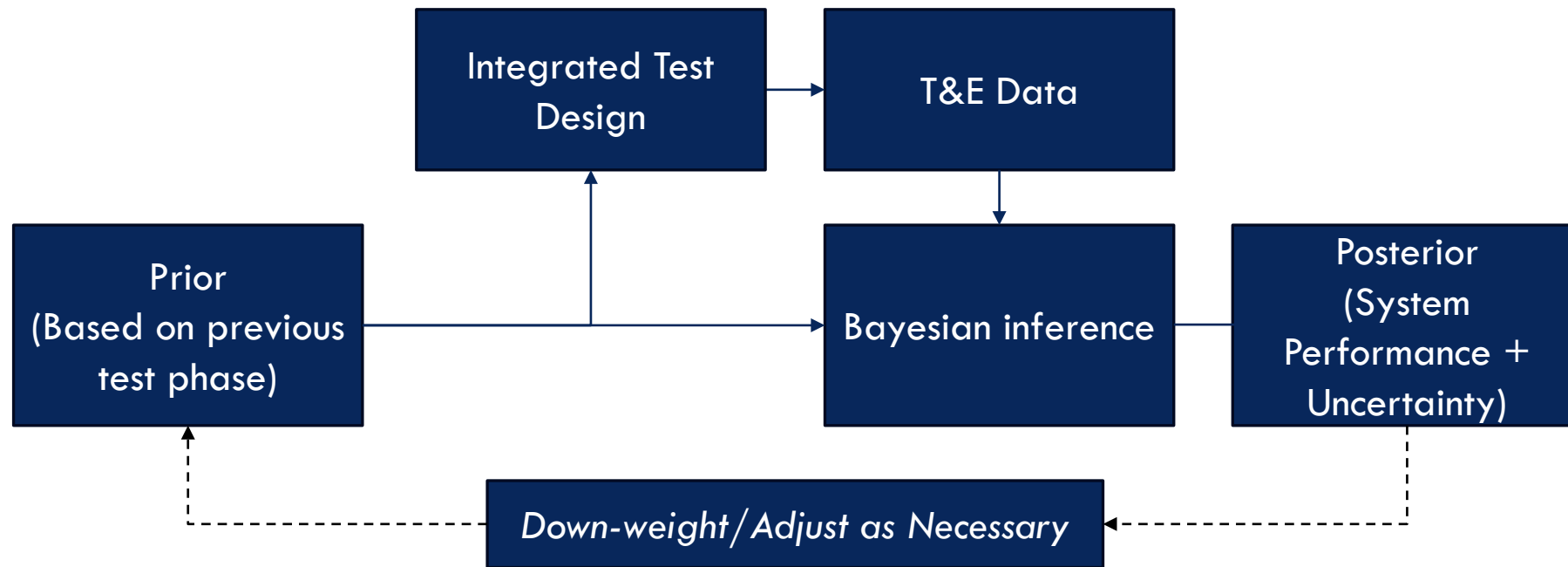
T&E
DATAWorks 2024 Tutorial on
Bayesian Methods







(Data and metrics are notional)



CASE STUDY 1: BAYESIAN RELIABILITY



ACQUISITION INNOVATION
RESEARCH CENTER

Publicly-released reliability dataset

- Metric: Miles before system abort (MBSA)
- Both DT and OT
- Eight variants

See:

- Dickinson et al. “Statistical methods for combining information: Stryker family of vehicles reliability case study.” 2015.
- Sieck et al. “A Framework for Using Priors in a Continuum of Testing.” 2024.
- Krometis et al. “A Comparison of Bayesian Methods for Integrated Test and Evaluation.” 2025.



Publicly-released reliability dataset

- Metric: Miles before system abort (MBSA)
- Both DT and OT
- Eight variants

See:

- Dickinson et al. “Statistical methods for combining information: Stryker family of vehicles reliability case study.” 2015.

- Sieck et al. “A Framework for Using Priors in a Continuum of Testing.” 2024.
- Krometis et al. “A Comparison of Bayesian Methods for Integrated Test and Evaluation.” 2025.



Informative Priors

Model reliability with Weibull distribution:

$$y_i | k, \tau_i \sim \text{Weibull}(k, \tau_i)$$

where:

- k is the *shape* parameter (failure change over time)
- τ_i is the *scale* parameter (units)

We infer k and τ_1, \dots, τ_8 from the data and then can compute the mean MBSA (MMBSA):

$$\bar{y}_i = \tau_i * \Gamma(1 + 1/k)$$

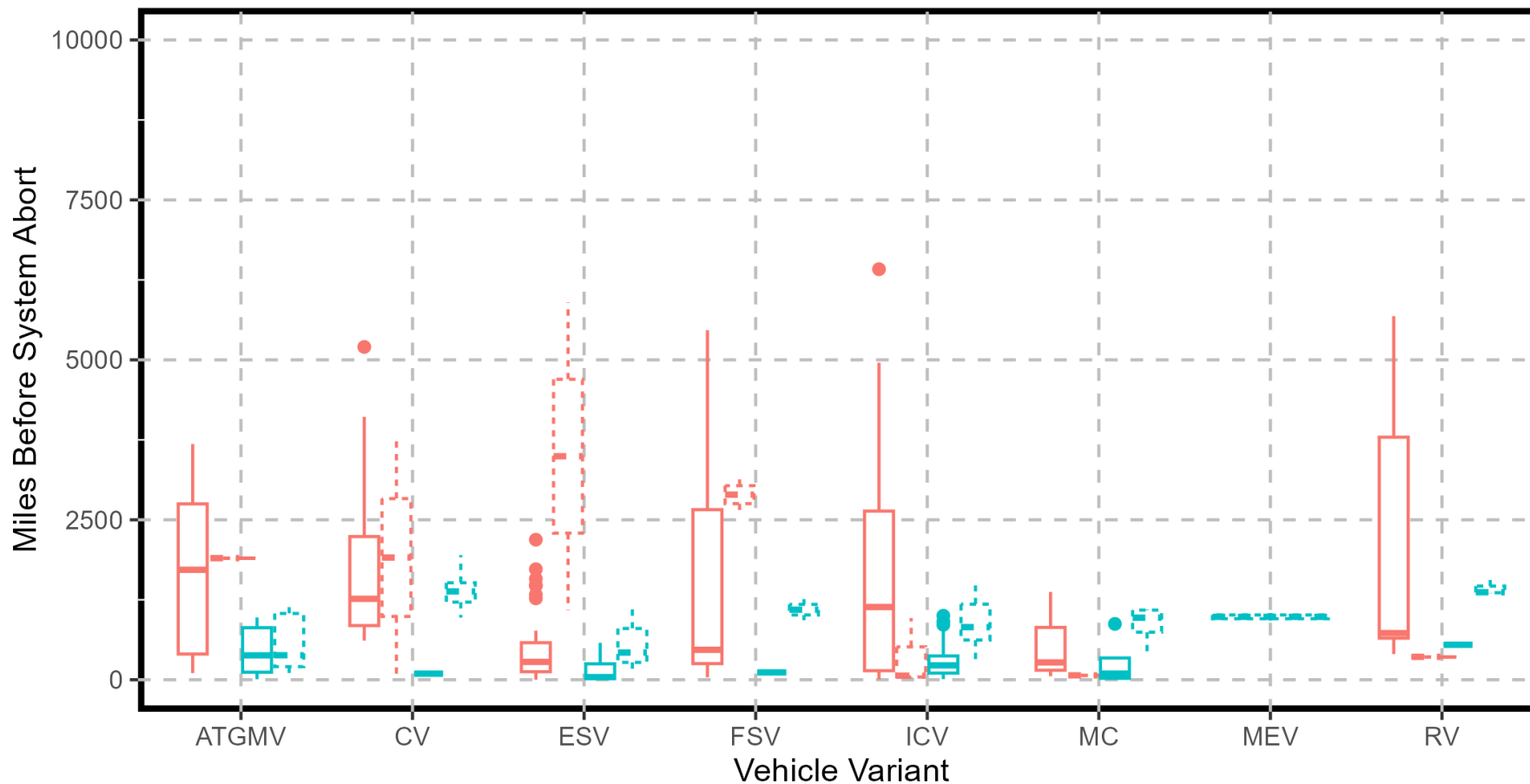


Uninformative – impart limited information via very vague priors:

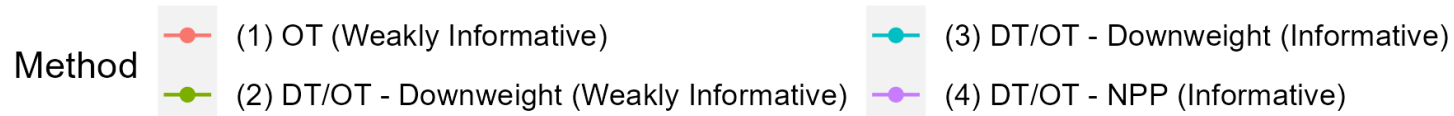
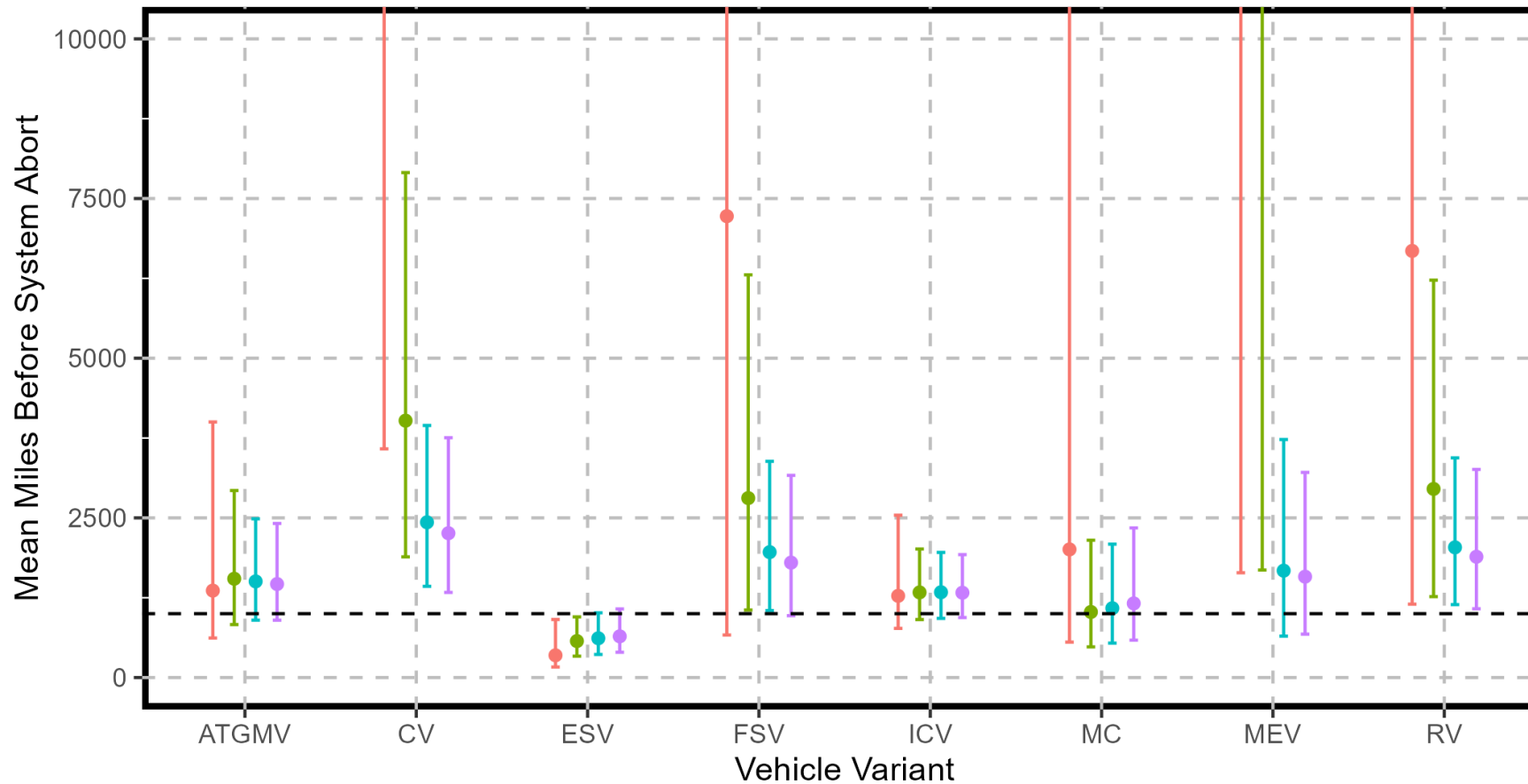
- Shape parameter k : Gamma(0.001,0.001)
- Scale parameter τ : Gamma(10^{-6} , 10^{-6})

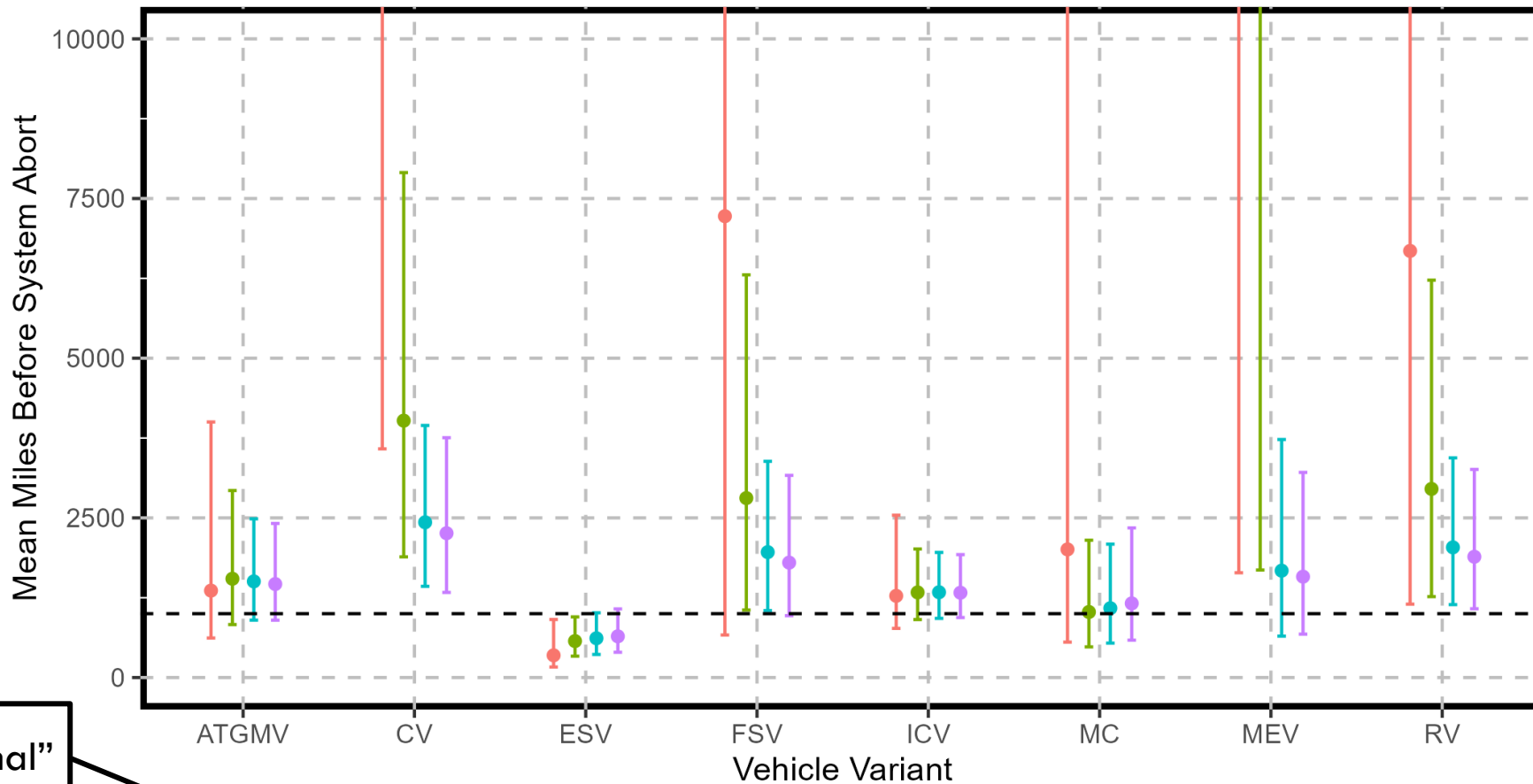
Informative – impart basic information:

- Shape parameter k : Gamma(4,4) (mean 1, std. dev. 0.5)
- Scale parameter τ : Gamma(4.6225,0.0043) (median 1,000, std. dev. 500)



Right Censor 0 1 Test Phase DT OT

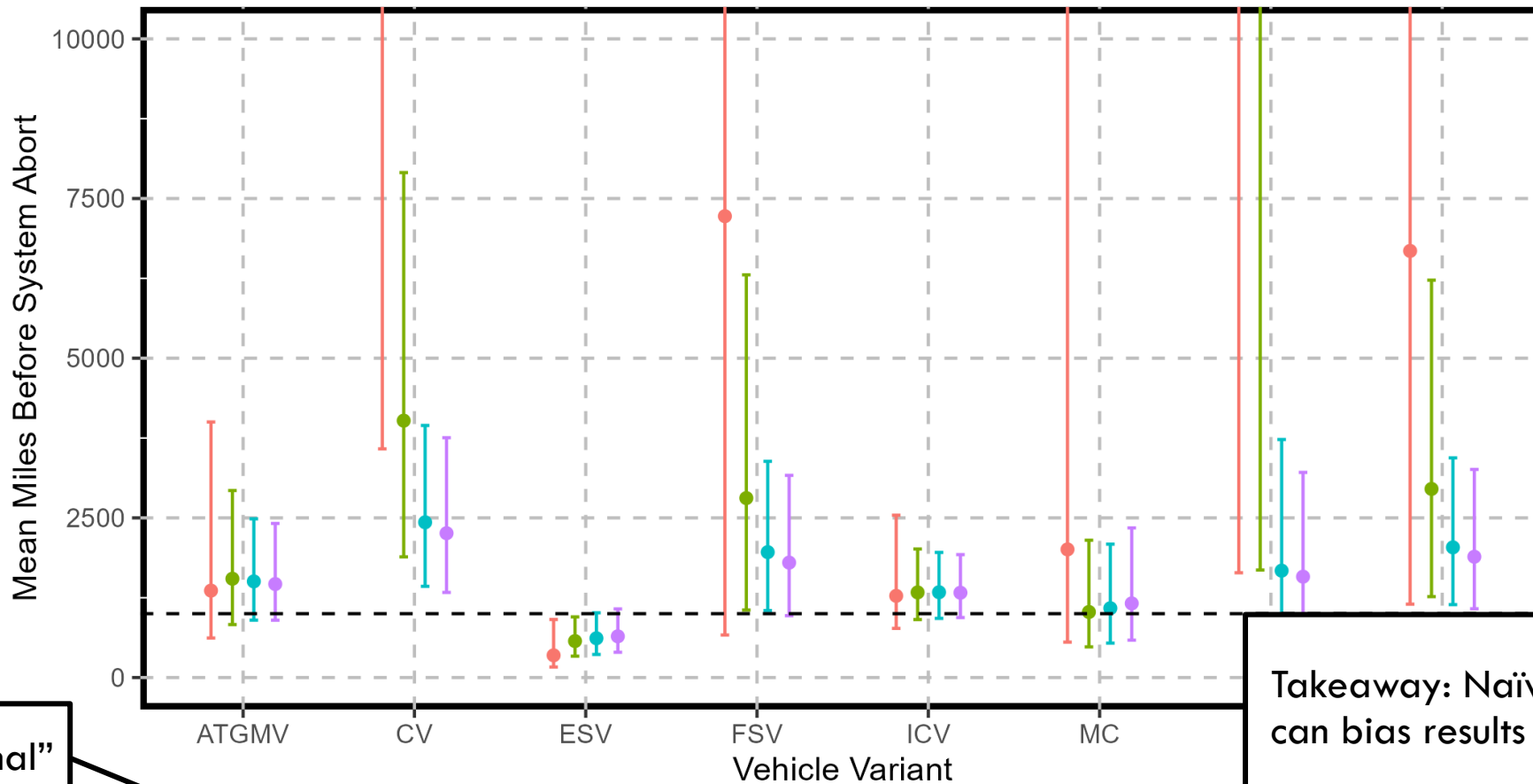




Analog of "traditional"

"Naïve" Bayes

- Method
- (1) OT (Weakly Informative)
 - (2) DT/OT - Downweight (Weakly Informative)
 - (3) DT/OT - Downweight (Informative)
 - (4) DT/OT - NPP (Informative)



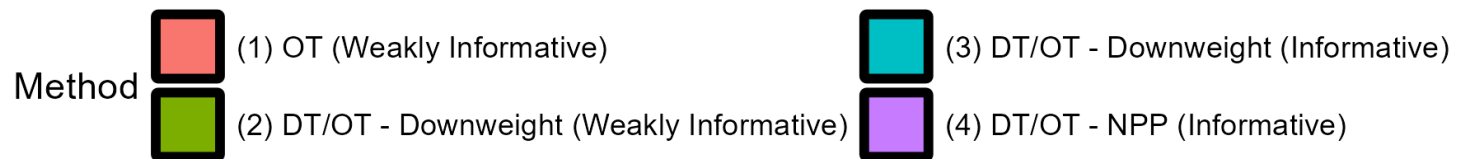
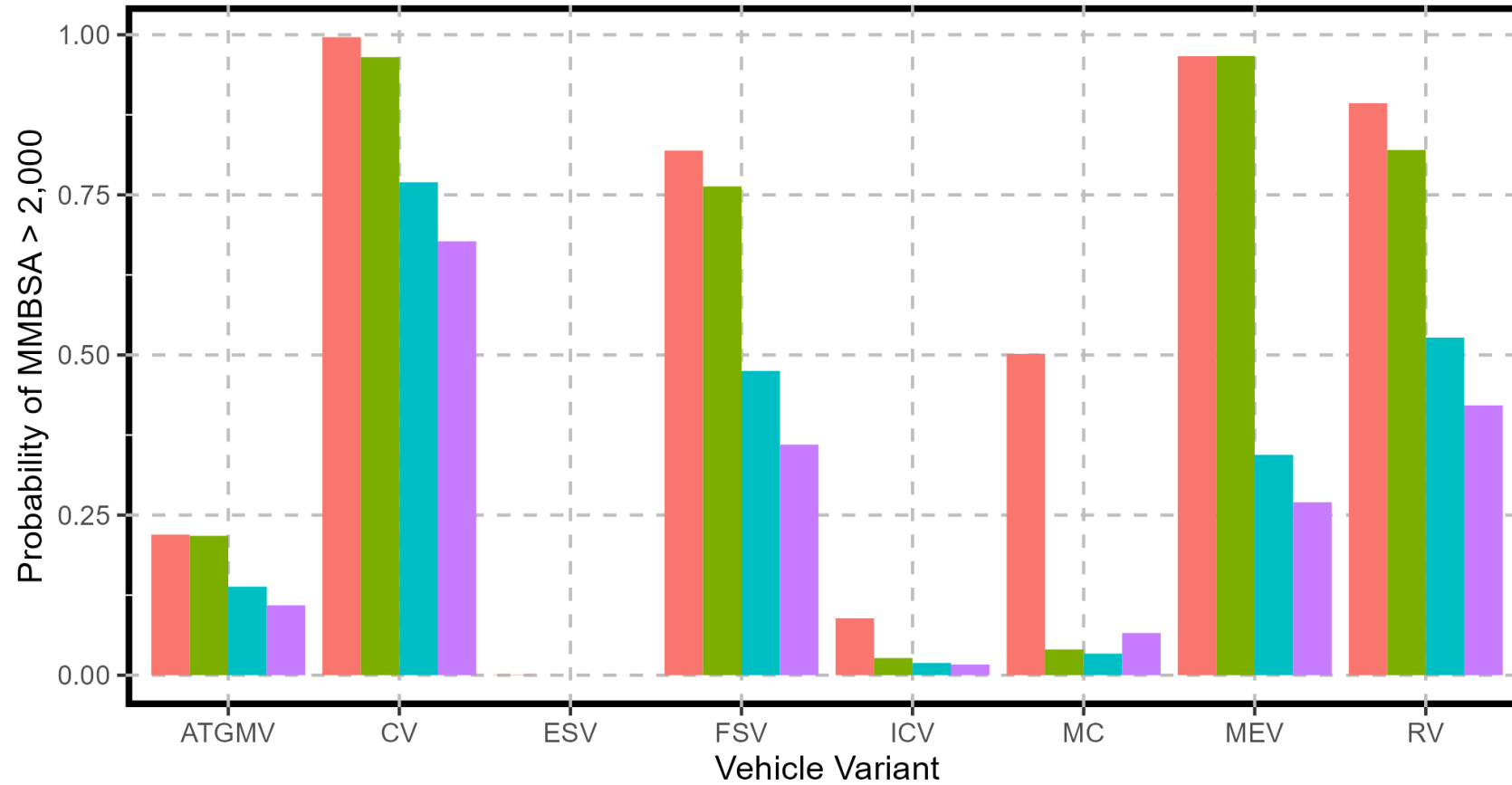
Analog of “traditional”

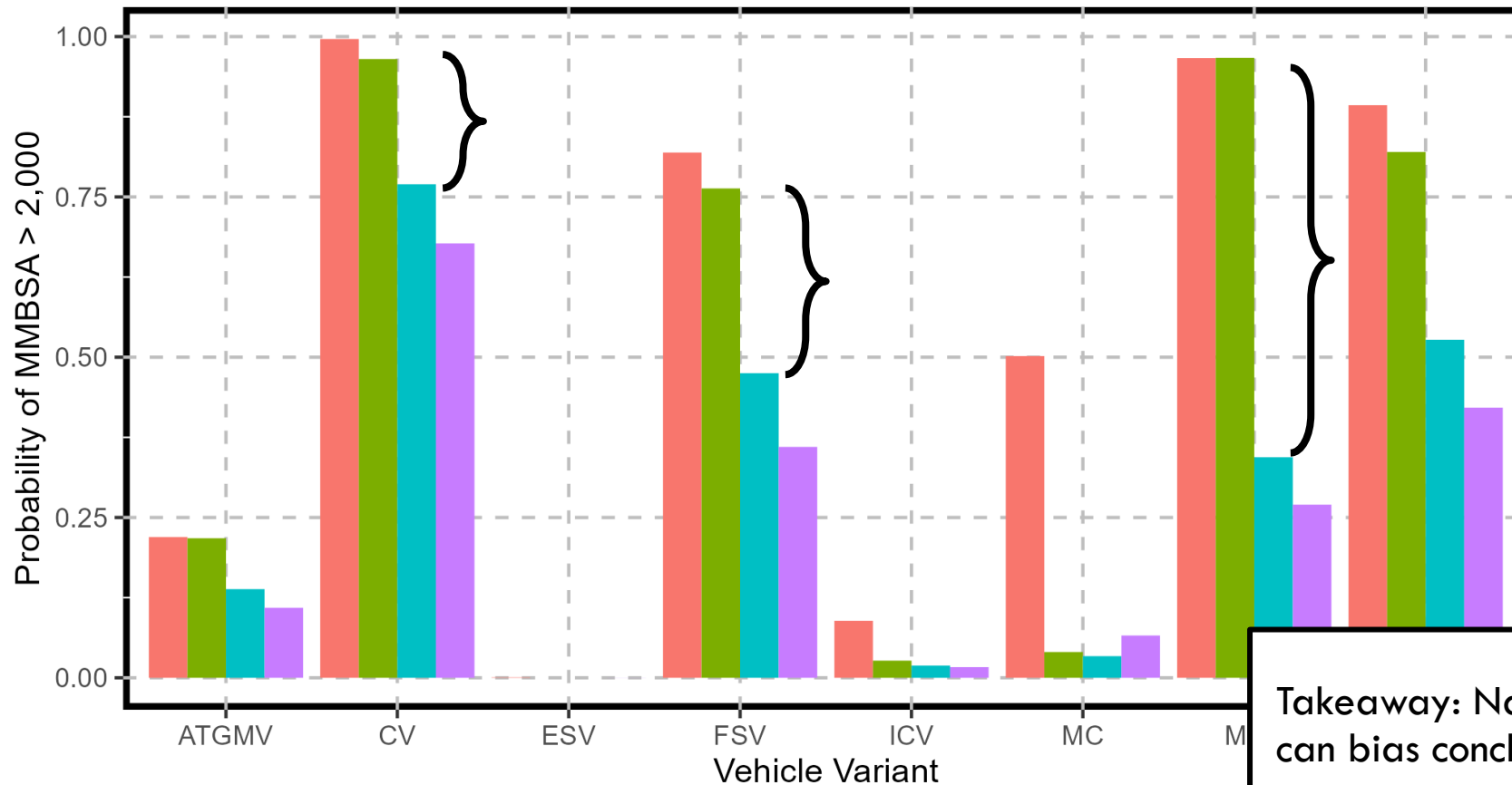
“Naïve” Bayes

Method

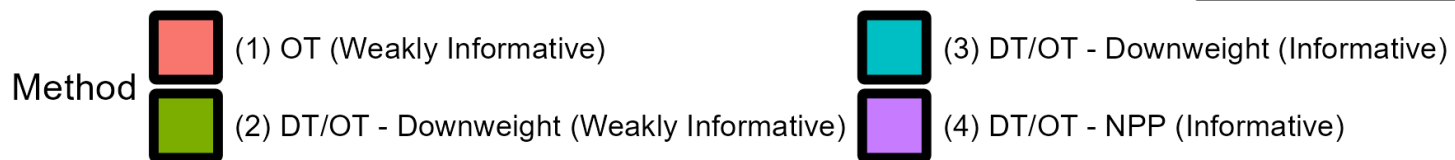
- (1) OT (Weakly Informative)
- (2) DT/OT - Downweight (Weakly Informative)
- (3) DT/OT - Downweight (Informative)
- (4) DT/OT - NPP (Informative)

Takeaway: Naïve approaches can bias results

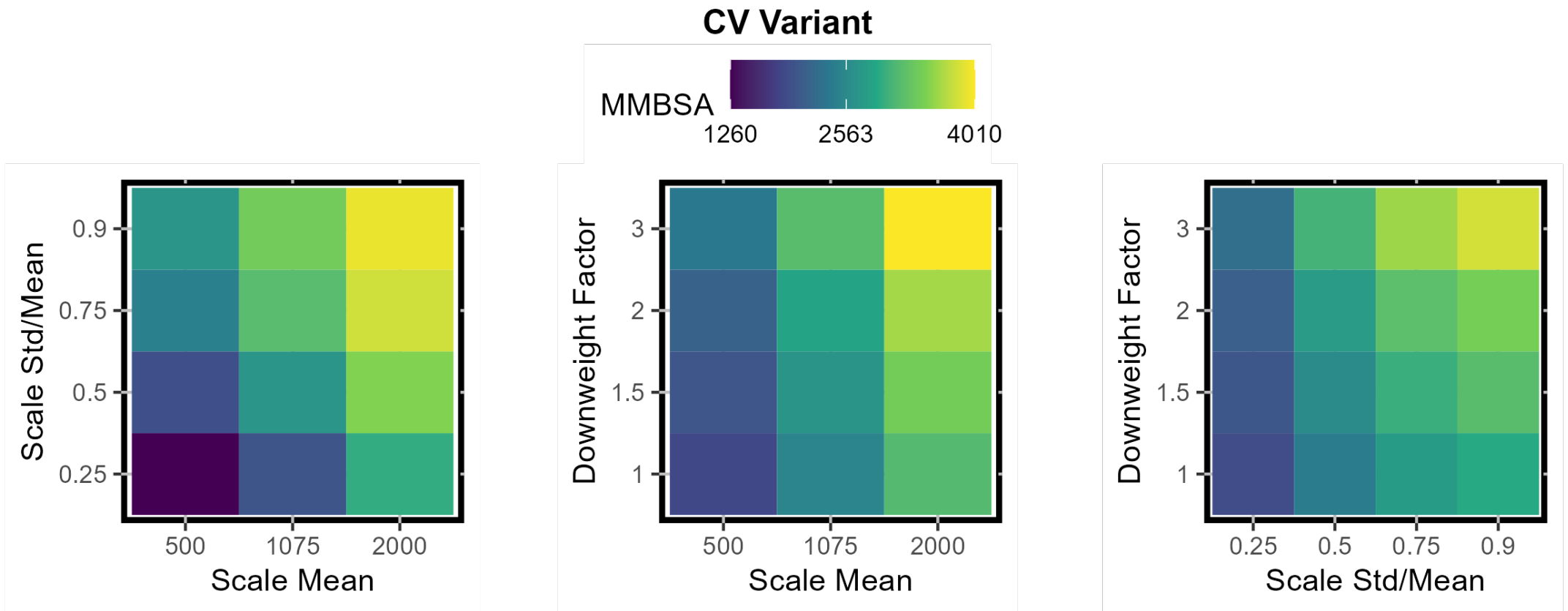




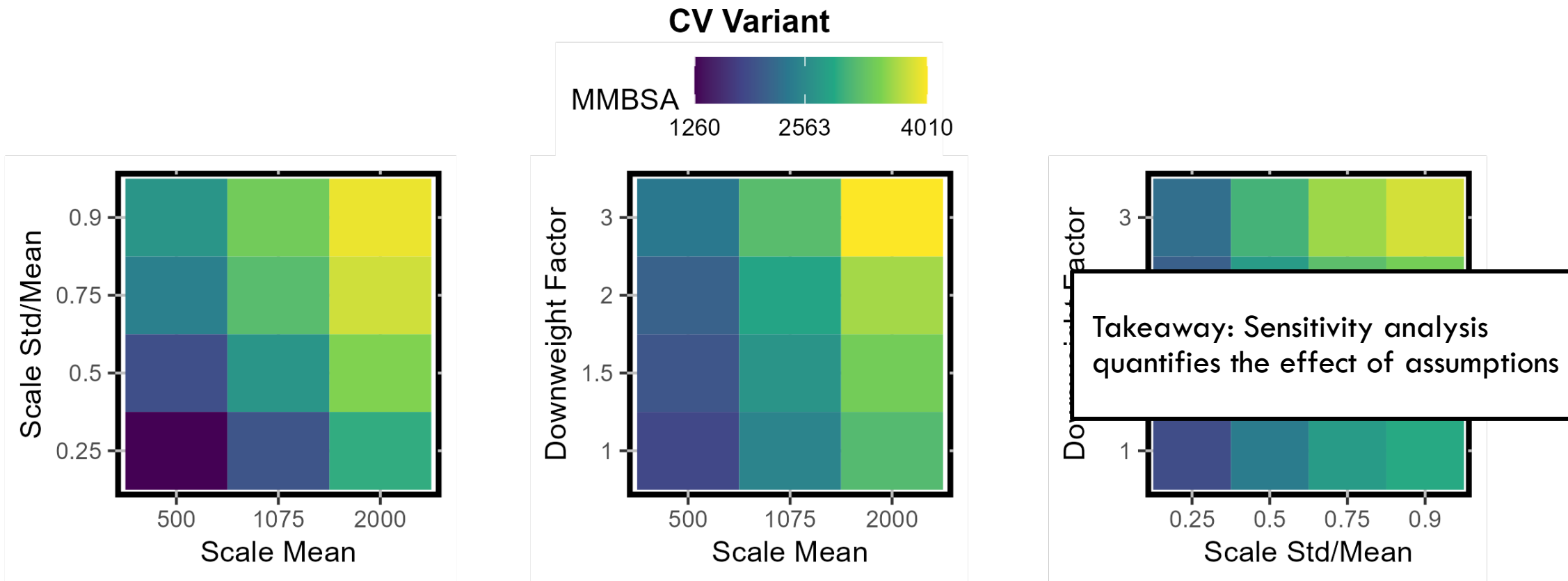
Takeaway: Naïve approaches can bias conclusions



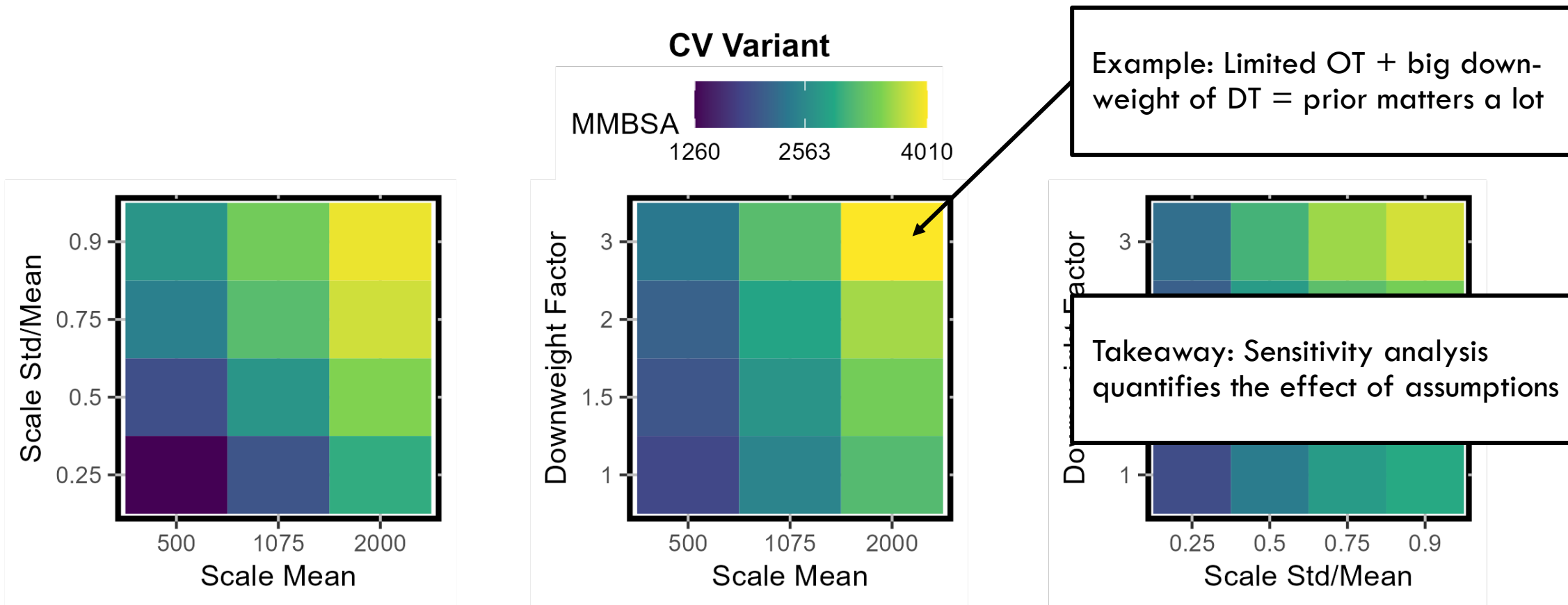
Key concern among Bayesian skeptics: How do modeling assumptions (including the prior) affect estimates?



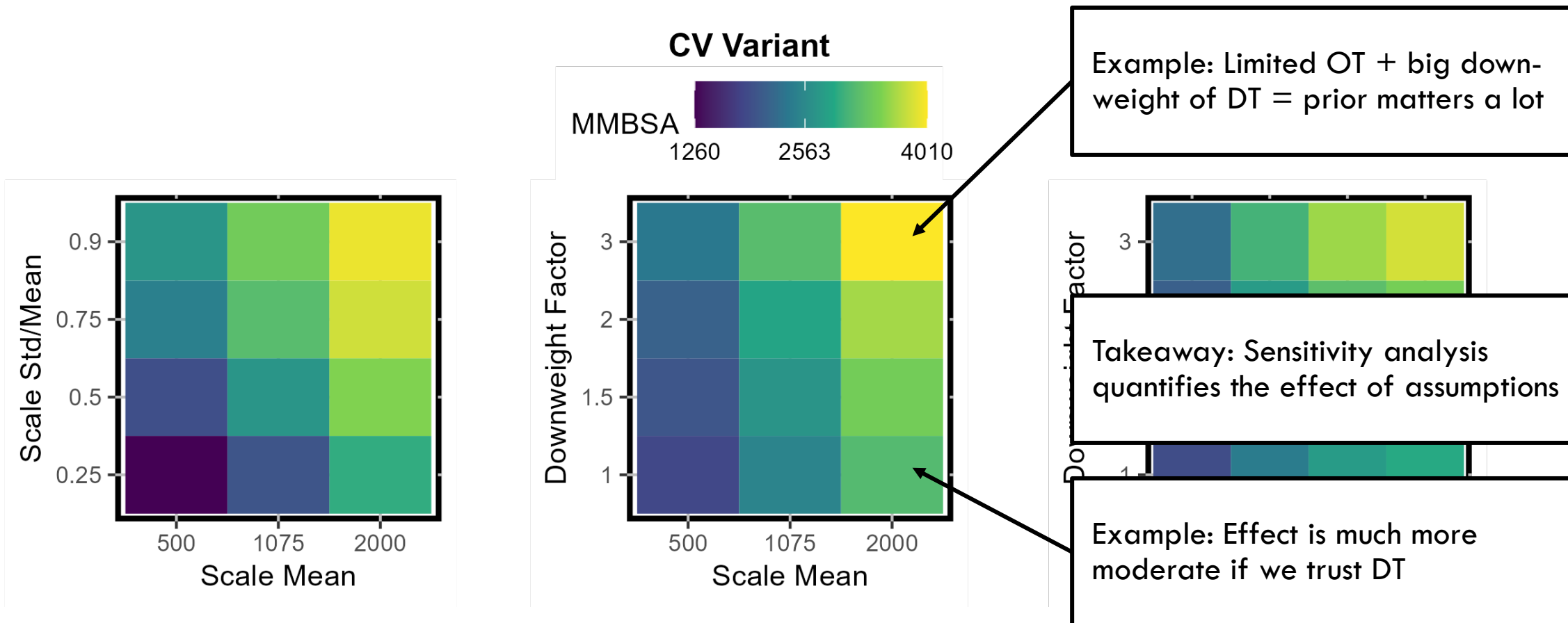
Key concern among Bayesian skeptics: How do modeling assumptions (including the prior) affect estimates?



Key concern among Bayesian skeptics: How do modeling assumptions (including the prior) affect estimates?



Key concern among Bayesian skeptics: How do modeling assumptions (including the prior) affect estimates?



Demonstrated the promise of Bayesian methods for integrating test data to improve estimates of operational performance

Highlighted that care should be taken in the development of assumptions, as attempts at objectivity can inadvertently bias results

CASE STUDY 2: CHANGING FACTORS



ACQUISITION INNOVATION
RESEARCH CENTER

IDA created¹ the following model of a synthetic counterfire radar:

$$Y = 79 - 6B + 4D - 7.5F + 5AF - 5.5BD + 4.5DF + 4D^2 - 9F^2$$

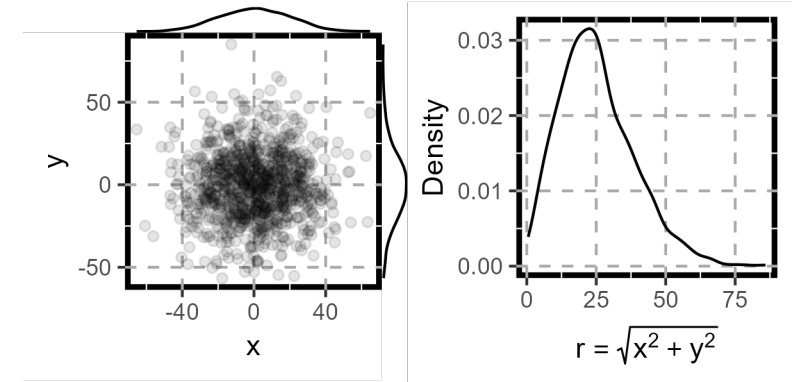
with the following factors:

Design Factor	Label	Type	Levels
Quadrant Elevation	<i>A</i>	Continuous	Low, High
Aspect Angle	<i>B</i>	Continuous	Incoming, Crossing
Munition Type	<i>C</i>	Categorical	Mortar, Rockets, Artillery
Shot Range	<i>D</i>	Continuous	Low, High
Operating Mode	<i>E</i>	Categorical	90, 360
Radar to Weapon Range	<i>F</i>	Continuous	Low, High

¹Ahrens, Monica, Rebecca Medlin, Keyla Pagán-Rivera, and John W. Dennis. "Case Study on Applying Sequential Analyses in Operational Testing." *Quality Engineering* 35, no. 3 (December 12, 2022): 534–45. <https://doi.org/10.1080/08982112.2022.2146510>.

Scenario	Rationale	Formula
Operations (“Real Life”)	Most complicated – full model	$79 - 6B + 4D - 7.5F - 5.5B * D + 4.5D * F + 5A * F + 4D^2 - 9F^2$
Operational Testing (OT)	Less fidelity than operations – drop quadratic terms	$79 - 6B + 4D - 7.5F - 5.5B * D + 4.5D * F + 5A * F$
Developmental Testing (DT)	Drop Quadrant Elevation (A)	$79 - 6B + 4D - 7.5F - 5.5B * D + 4.5D * F$
Modeling & Simulation (M&S)	Drop Radar to Weapon Range (F)	$79 - 6B + 4D - 5.5B * D$

Assume location error is normally distributed in two dimensions: Rayleigh distribution



Models give distribution mean, which can then be used to generate data



Can compare:

- Analysis methods
- Design of experiments techniques (test designs via skpr package)

For problems with:

- Different numbers of test phases/data sources
- Varying data sizes, e.g., trials and reps by phase
- Evolving test factors
- Shifts/biases in test data (e.g., in M&S data)
- Different error/noise in measurements

Five methods considered:

- **Frequentist:**
 - Using OT data only
 - All data, Blocking: With shift factors added for M&S and DT
 - All data, Without blocking: No shift factors for M&S and DT
- **Bayesian informative priors w/ downweighting:**
 - Resetting intercept uncertainty to prior value
 - Doubling intercept uncertainty

Key metric: RMSE between Rayleigh means

- Fitted model vs. **Operational model**
- Computed on full factorial dataset generated using the Operational model

Five methods considered:

- **Frequentist:**

- Using OT data only
- All data, Blocking: With shift factors added for M&S and DT
- All data, Without blocking: No shift factors for M&S and DT

- **Bayesian informative priors w/ downweighting:**

- Resetting intercept uncertainty to prior value
- Doubling intercept uncertainty

Key metric: RMSE between Rayleigh means

- Fitted model vs. **Operational model**
- Computed on full factorial dataset generated using the Operational model

Benefit of working with synthetic data: We know the “truth”!

Consider scenarios where OT is limited and M&S is quite a bit larger than DT

Intuition:

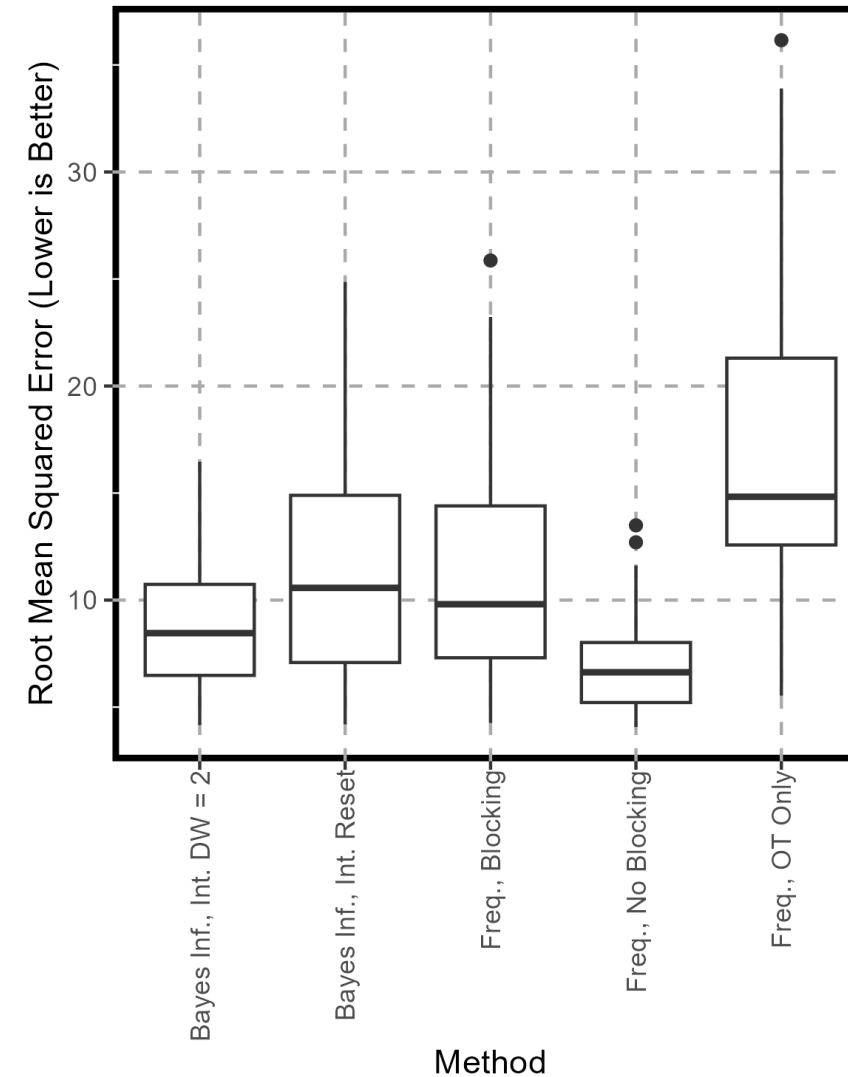
- Integrated testing should provide a benefit
- Challenge of managing different data sizes and changing test factors

Parameter	Value
M&S Trials	Full factorial (9 trials)
M&S Reps	100
DT Trials	10, 20, 40
DT Reps	5, 10
OT Trials	10, 20
OT Reps	1, 2
DT Optimality	D
OT Optimality	D

(Additional assumption: No more than 200 DT datapoints.)

Takeaways:

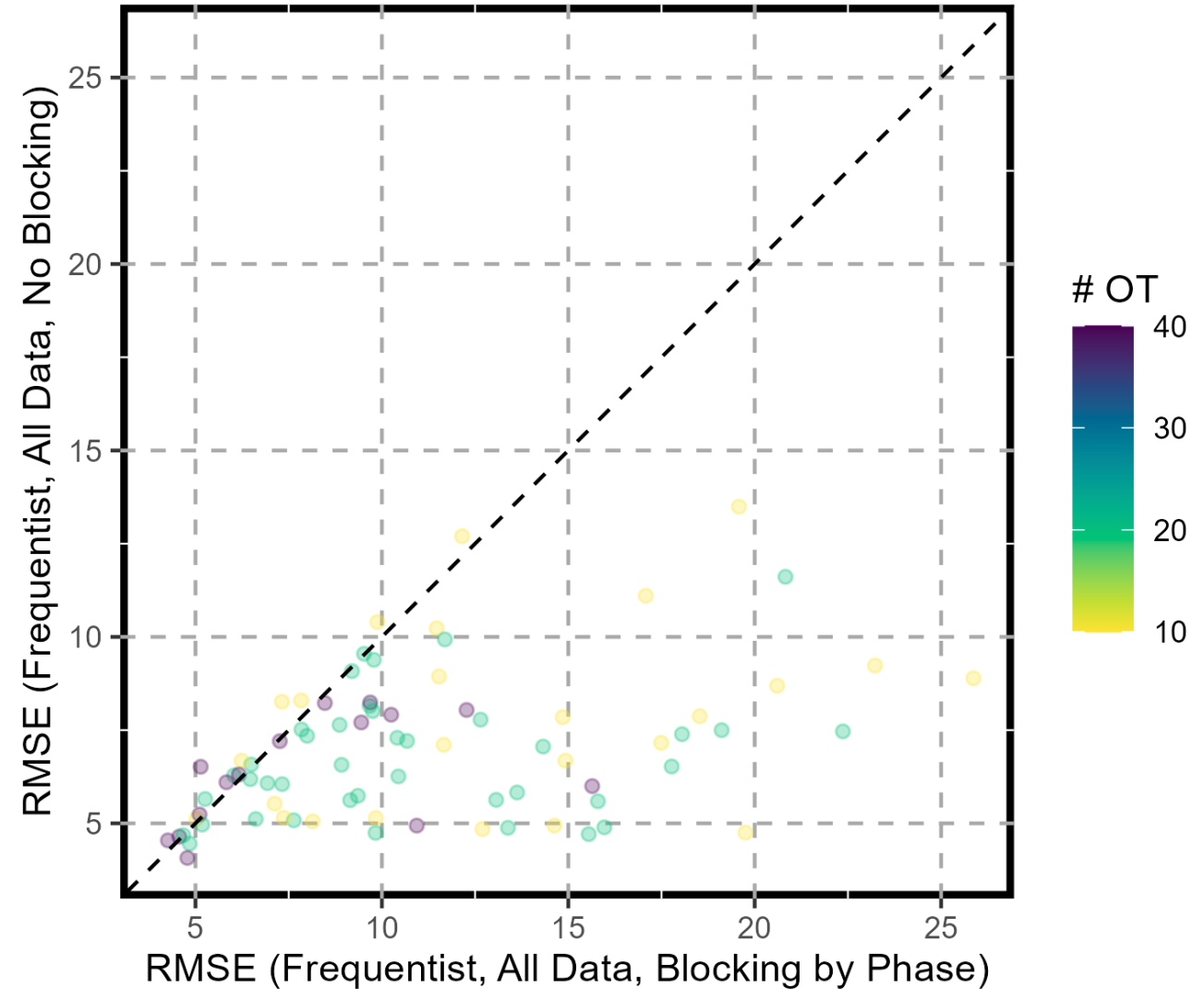
- Integrated evaluation helps provide better models
- Without blocking seems to do a little better



Error is lower for integrated model *without* blocking

Takeaways:

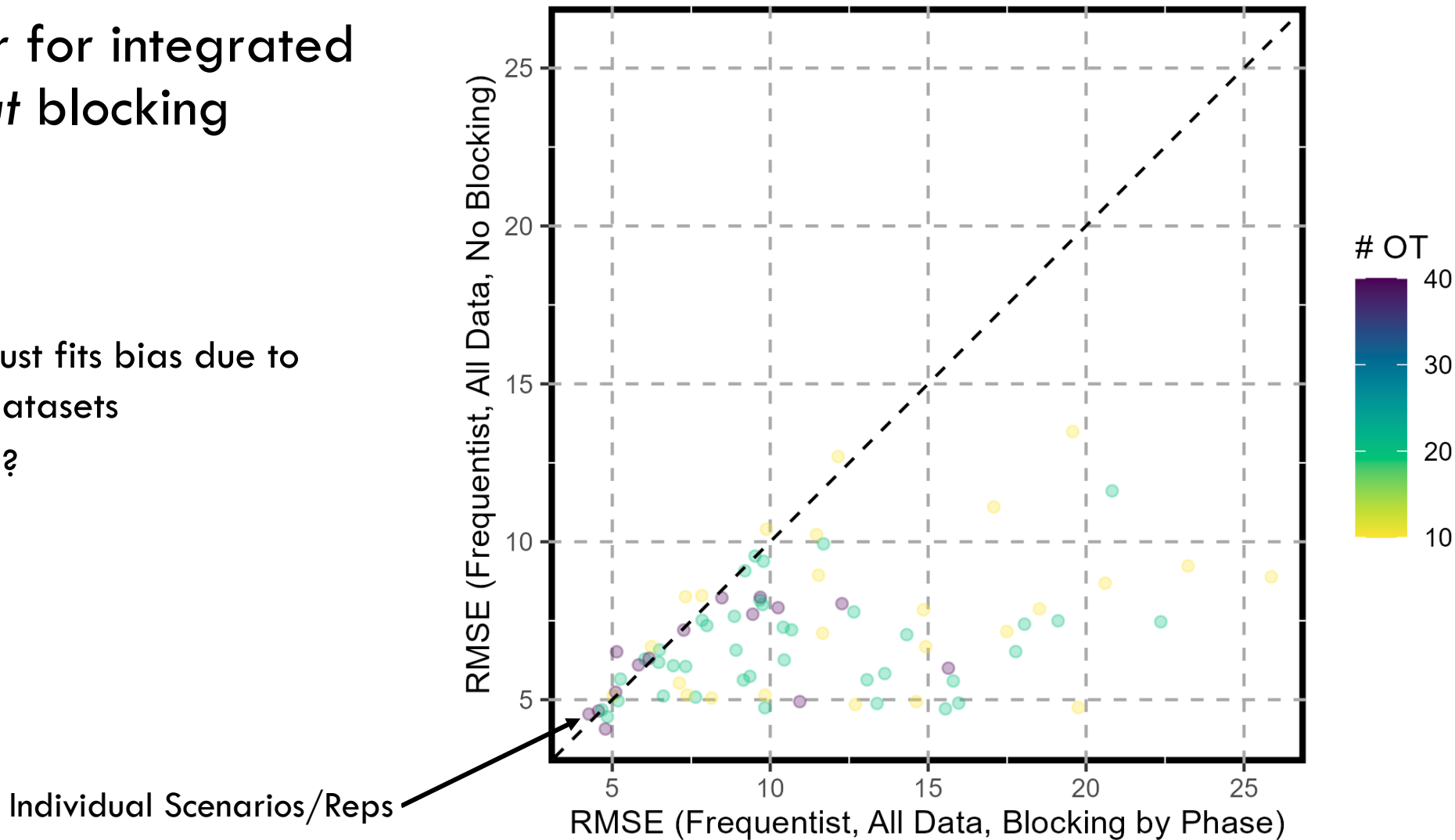
- Blocking shifts just fits bias due to noise in small datasets
- Blocking is bad?



Error is lower for integrated model *without* blocking

Takeaways:

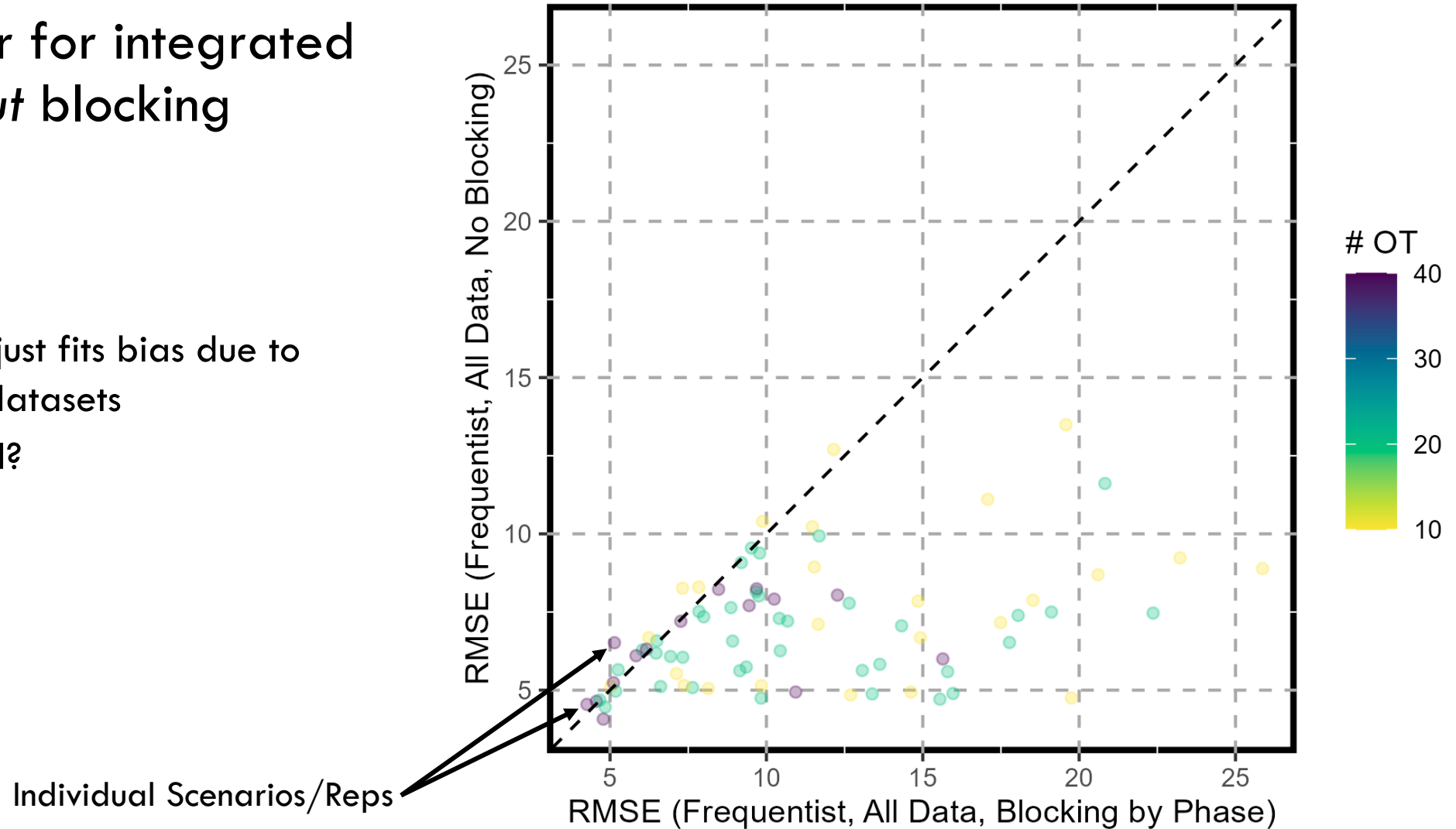
- Blocking shifts just fits bias due to noise in small datasets
- Blocking is bad?



Error is lower for integrated model *without* blocking

Takeaways:

- Blocking shifts just fits bias due to noise in small datasets
- Blocking is bad?



Consider scenarios where OT is limited and M&S is quite a bit larger than DT

- Add random bias to M&S, DT model intercepts

	M&S	DT	OT
Intercept	99.6	86.0	79

Intuition:

- Integrated testing should provide a benefit
- Biases in some of the data might change results?

Parameter	Value
M&S Trials	Full factorial (9 trials)
M&S Reps	100
DT Trials	10, 20, 40
DT Reps	5, 10
OT Trials	10, 20
OT Reps	1, 2
DT Optimality	D
OT Optimality	D

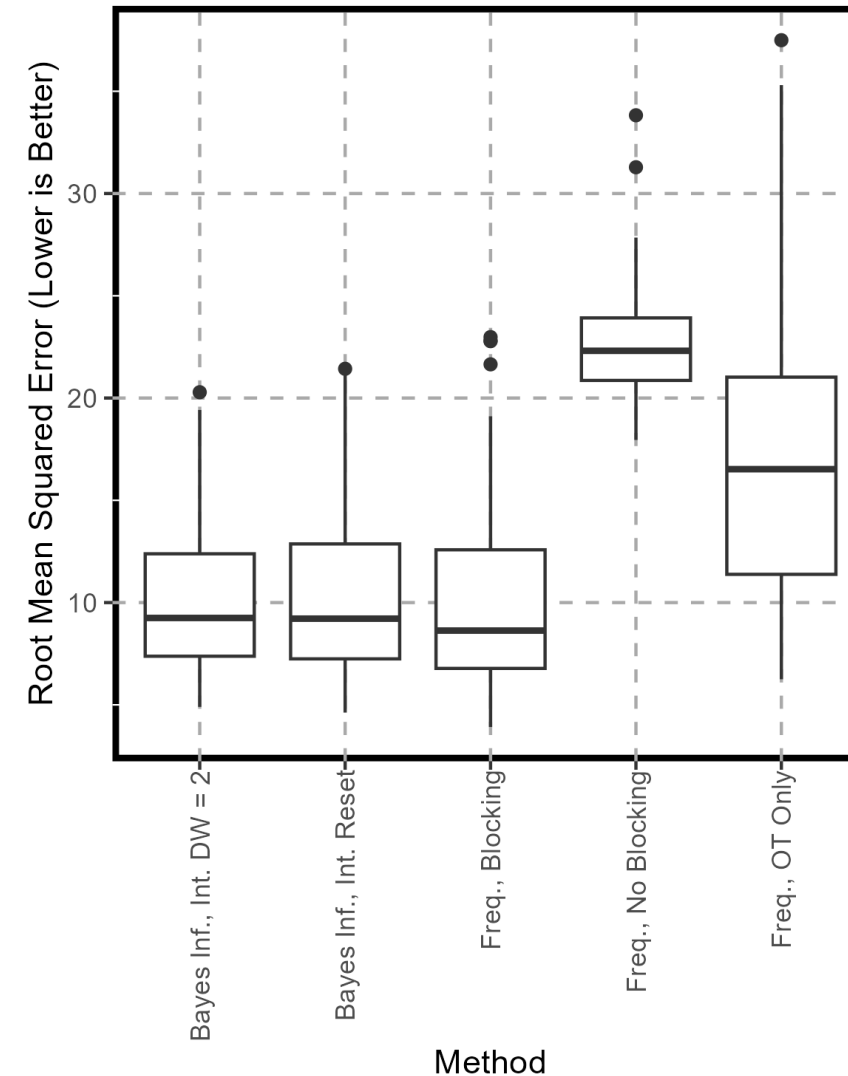
(Additional assumption: No more than 200 DT datapoints.)

Integration without blocking makes estimates

- Worse than using OT only
- Much worse than other integration methods

Takeaways:

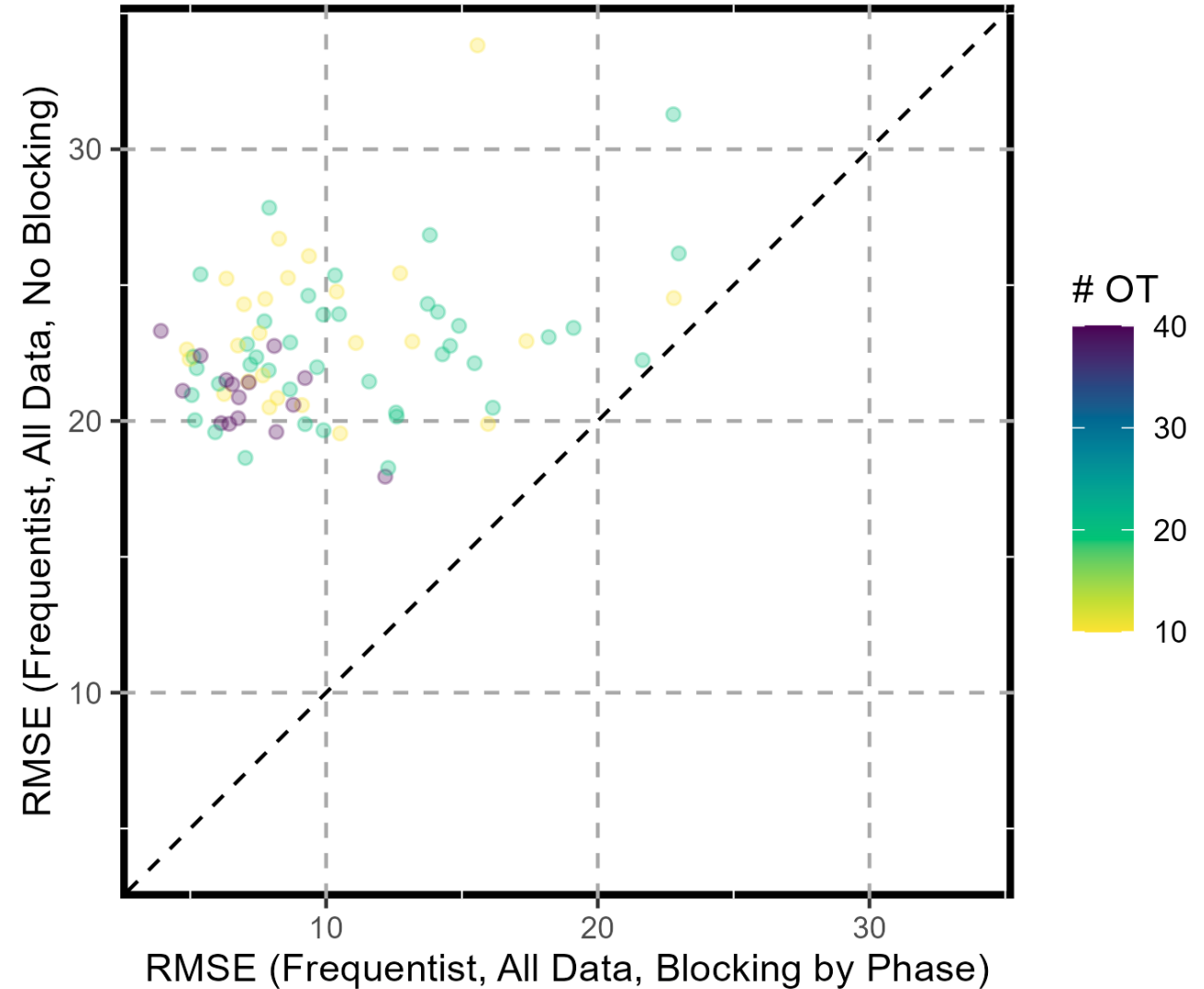
- Integrated testing *mostly* helps provide better models
- **But** some care must be taken in how the data is integrated



Error is dramatically lower for integrated model *with* blocking

Takeaways:

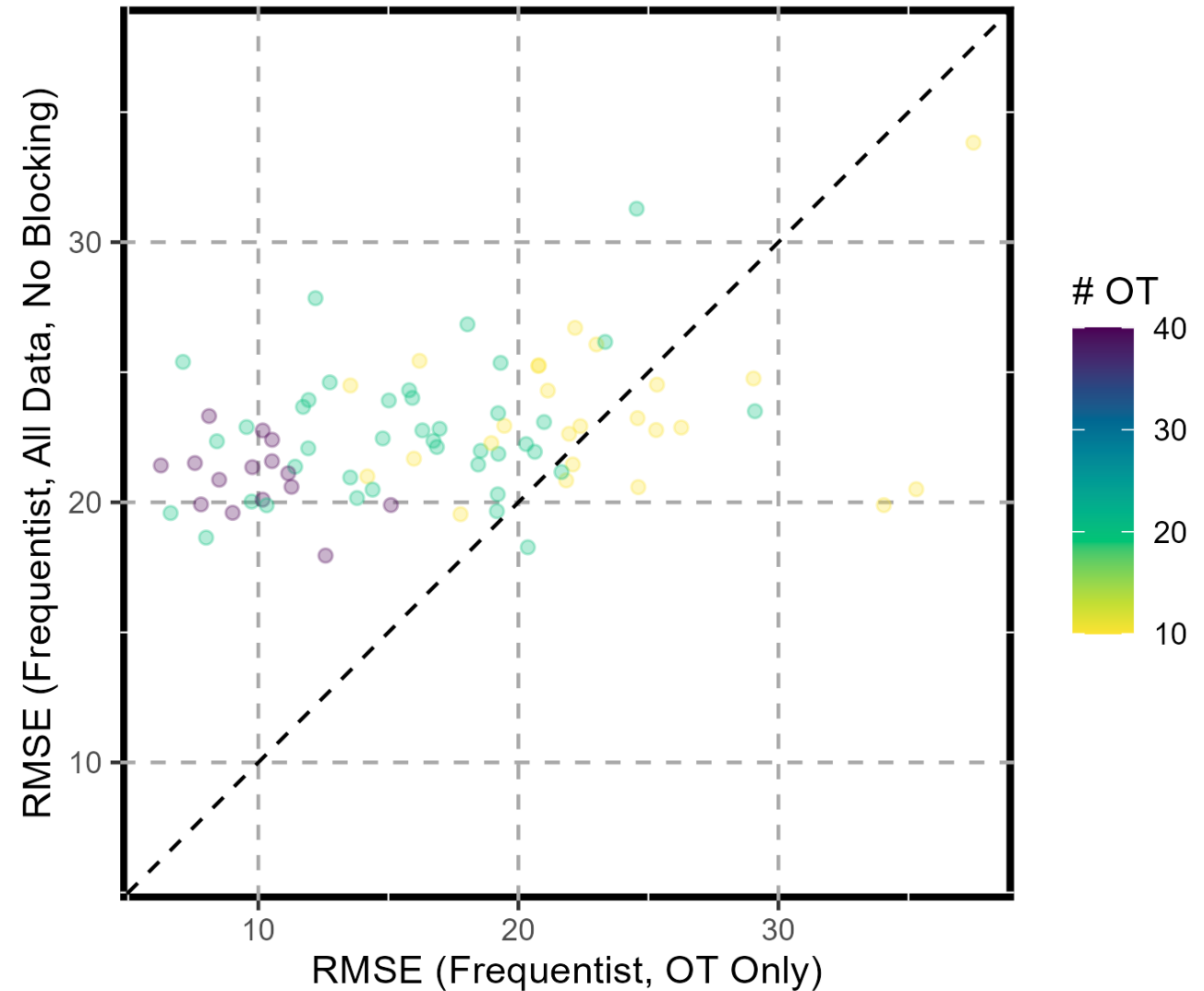
- Blocking allows accounting for differences in data sources
- Blocking is good?



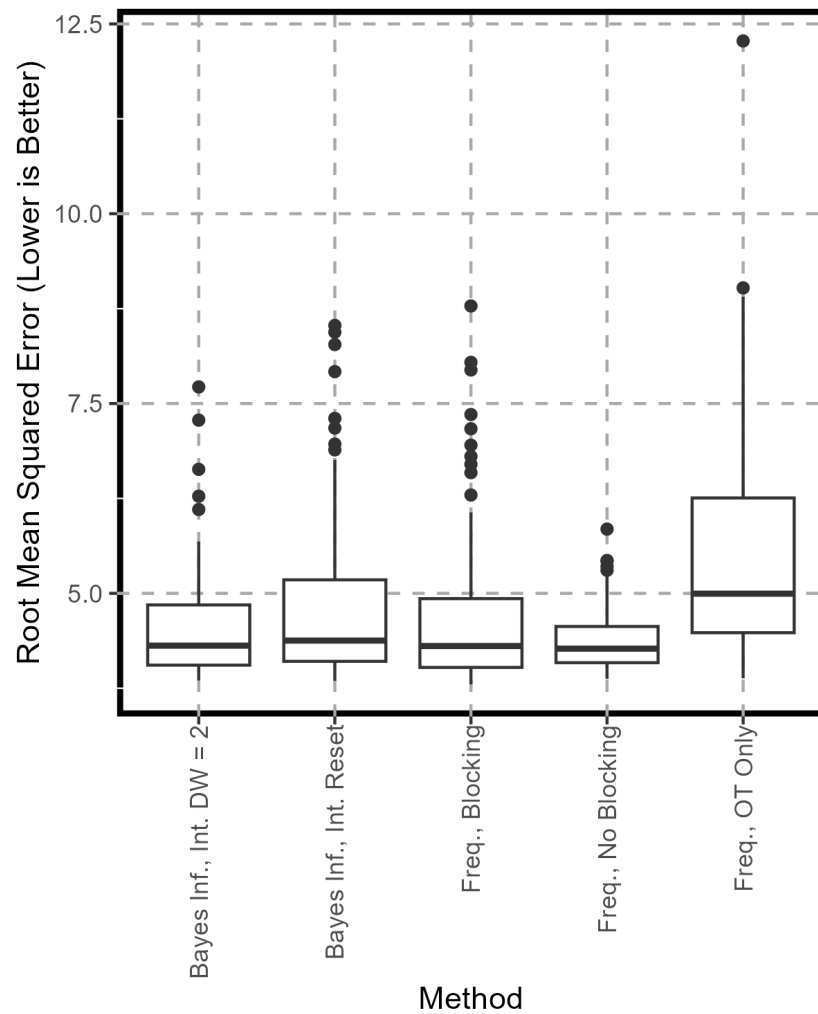
Non-blocking model
actually makes estimates
worse than single phase
model

Takeaways:

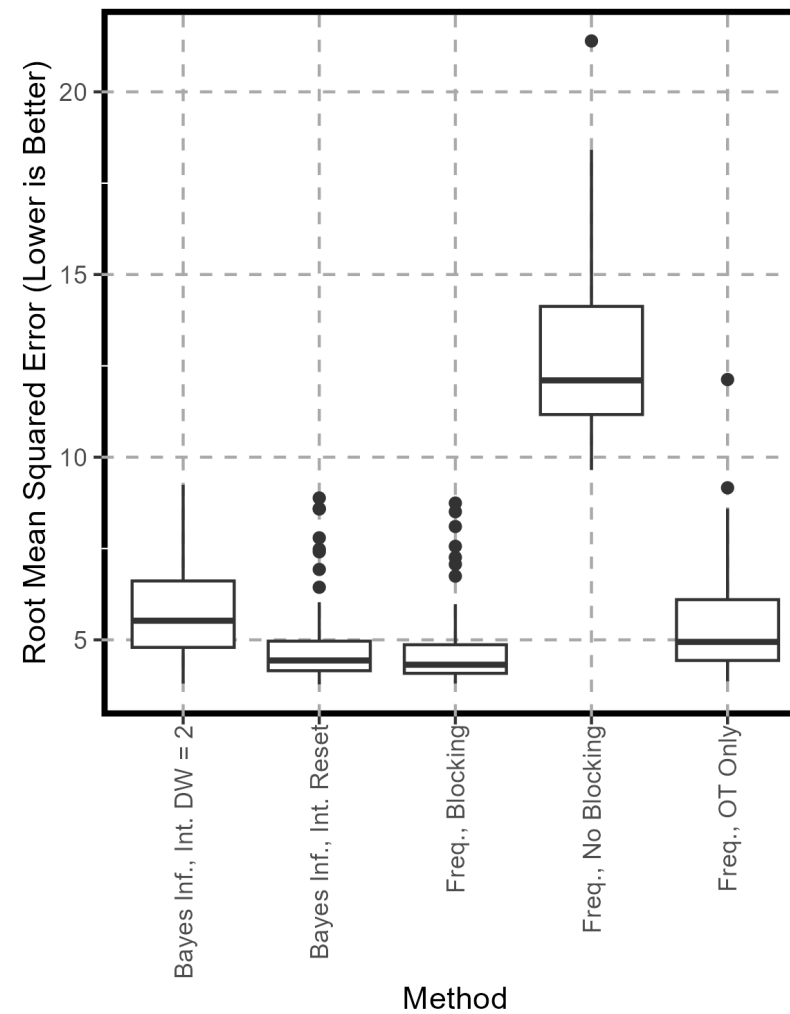
- Integration of information can make analysis worse if not done carefully



Unbiased M&S/DT



Biased M&S/DT



CASE STUDY 3: BINARY DATA



ACQUISITION INNOVATION
RESEARCH CENTER

Given

- prior $Beta(a_0, b_0)$ on probability of failure p (e.g., assembled from downweighted previous test results)
- binary test results a successes and b failures,

the Bayesian posterior distribution is $Beta(a_0 + a, b_0 + b)$.

Also, given a number of tests n , the probability of getting a successes (and $b = n - a$ failures) is given by the binomial distribution

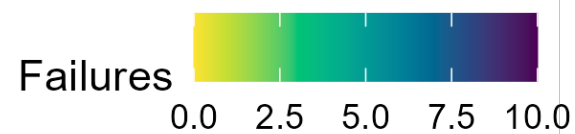
$$\binom{n}{a} p^a (1 - p)^{n-a}$$

where p is the true reliability.

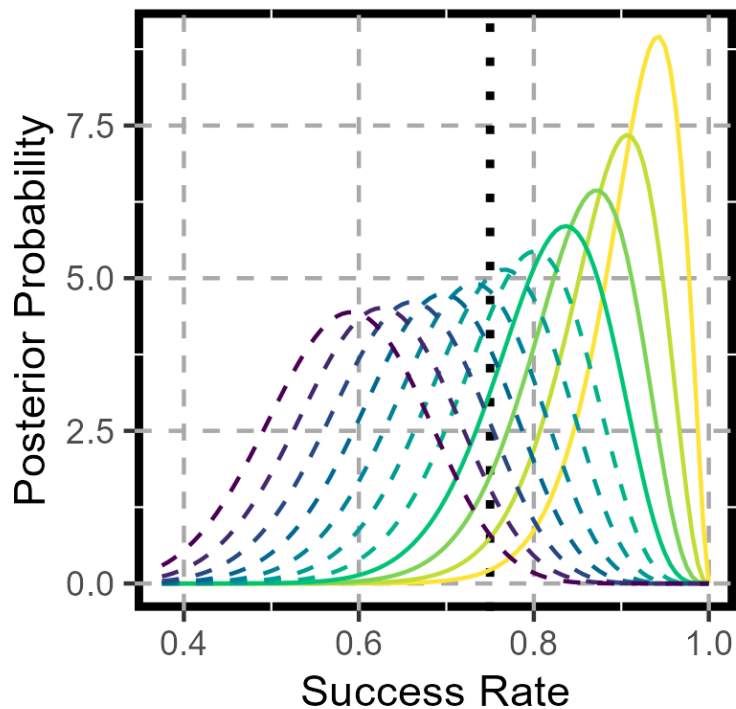
For now, assume the following parameters:

Parameter	Value
Prior successes (a_0)	17.84
Prior failures (b_0)	2.656
Prior mean ($a_0 / (a_0 + b_0)$)	0.87
Prior sample size ($a_0 + b_0$)	20.496
True reliability (p)	0.85
Threshold	0.75
Lower credible interval	0.2

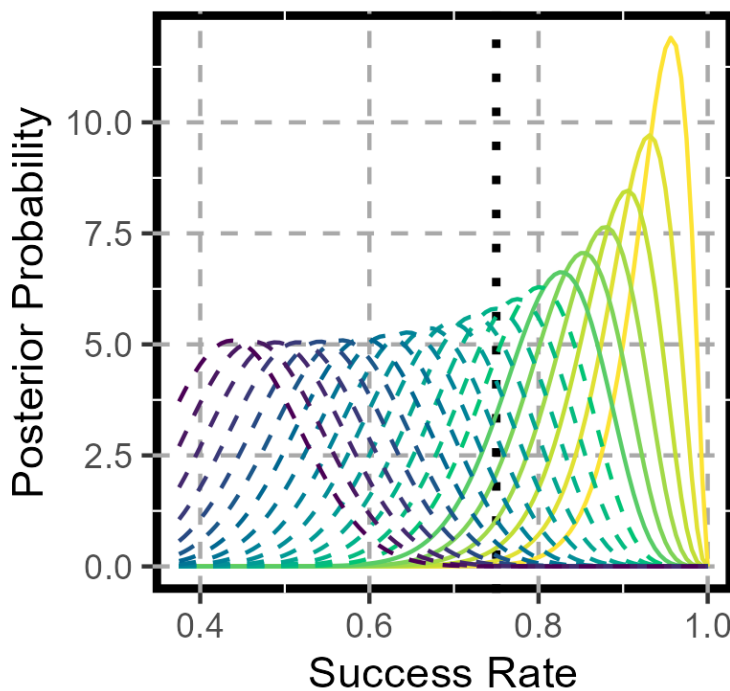
Pass/Fail Pass Fail



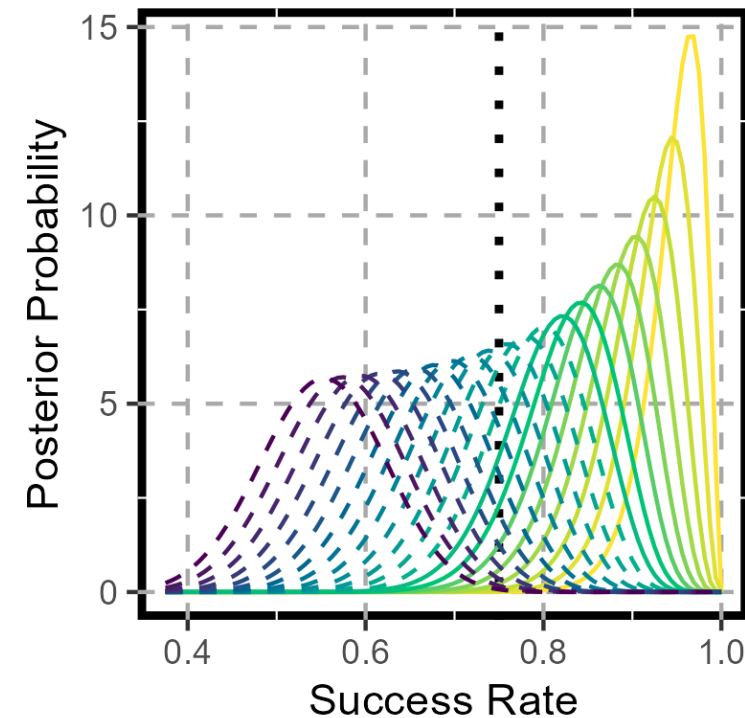
10 shots



20 shots



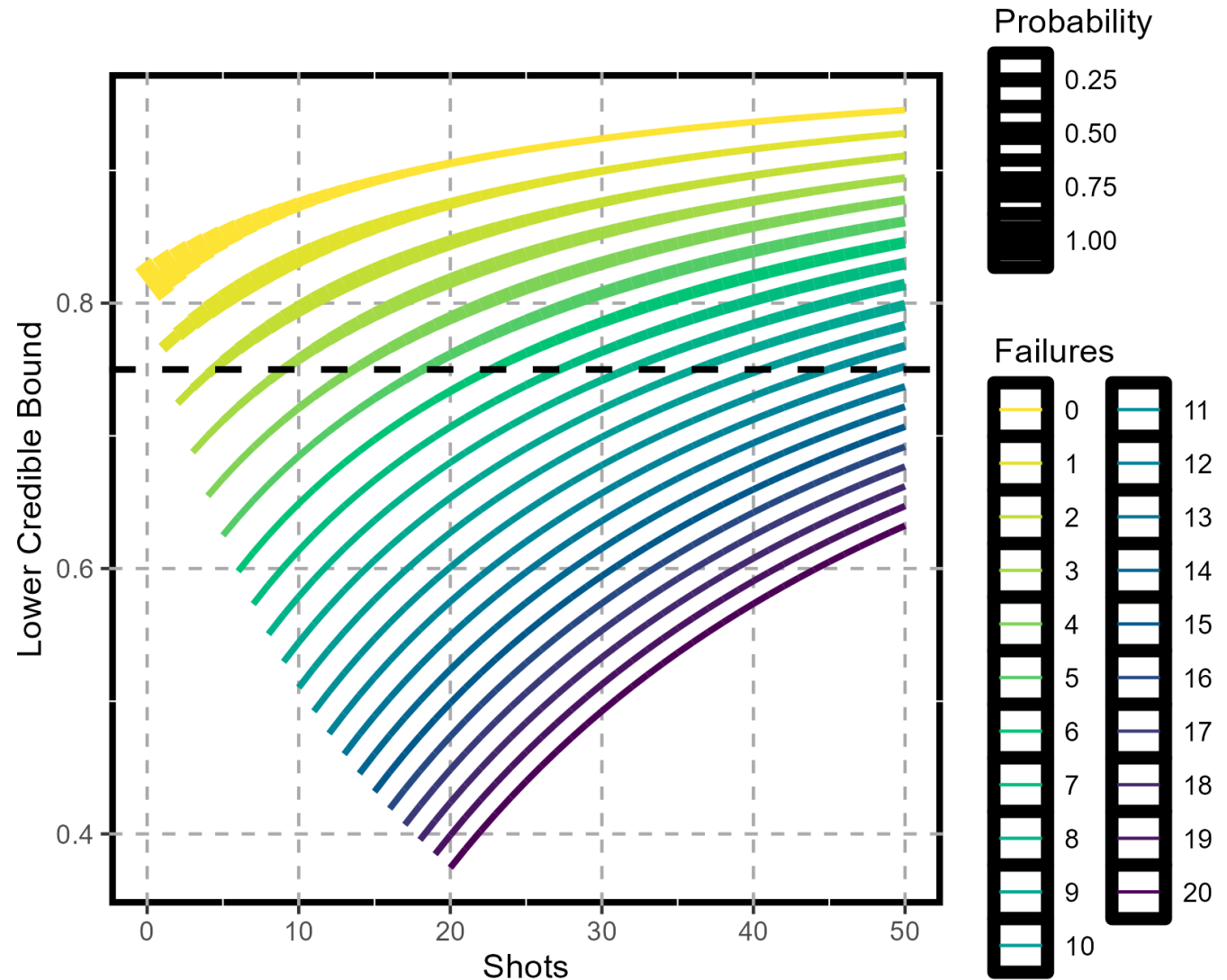
30 shots



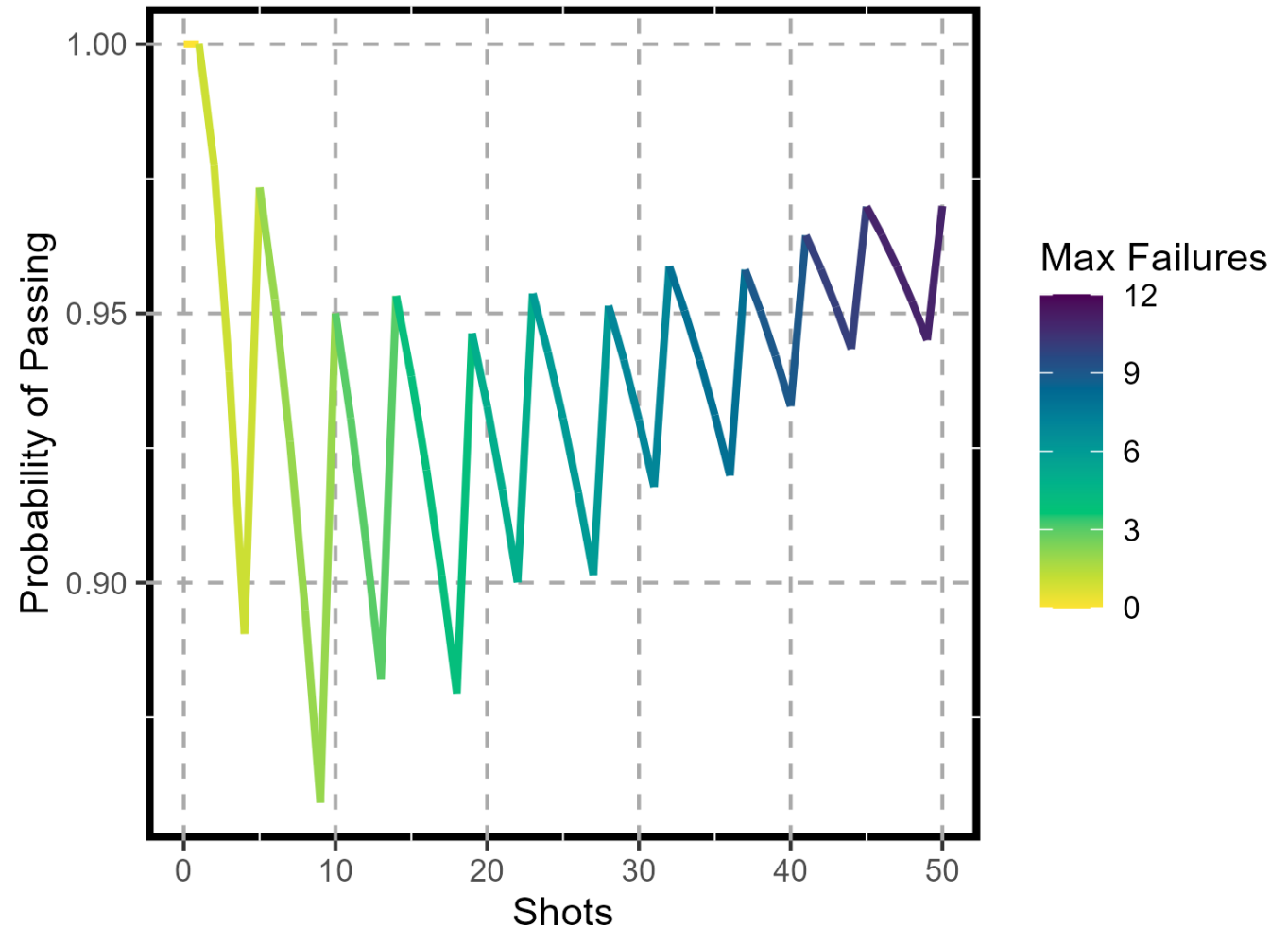
Compare lower credible bounds vs. threshold

Above threshold = pass

Thicker lines = more probable outcome (according to binomial distribution)



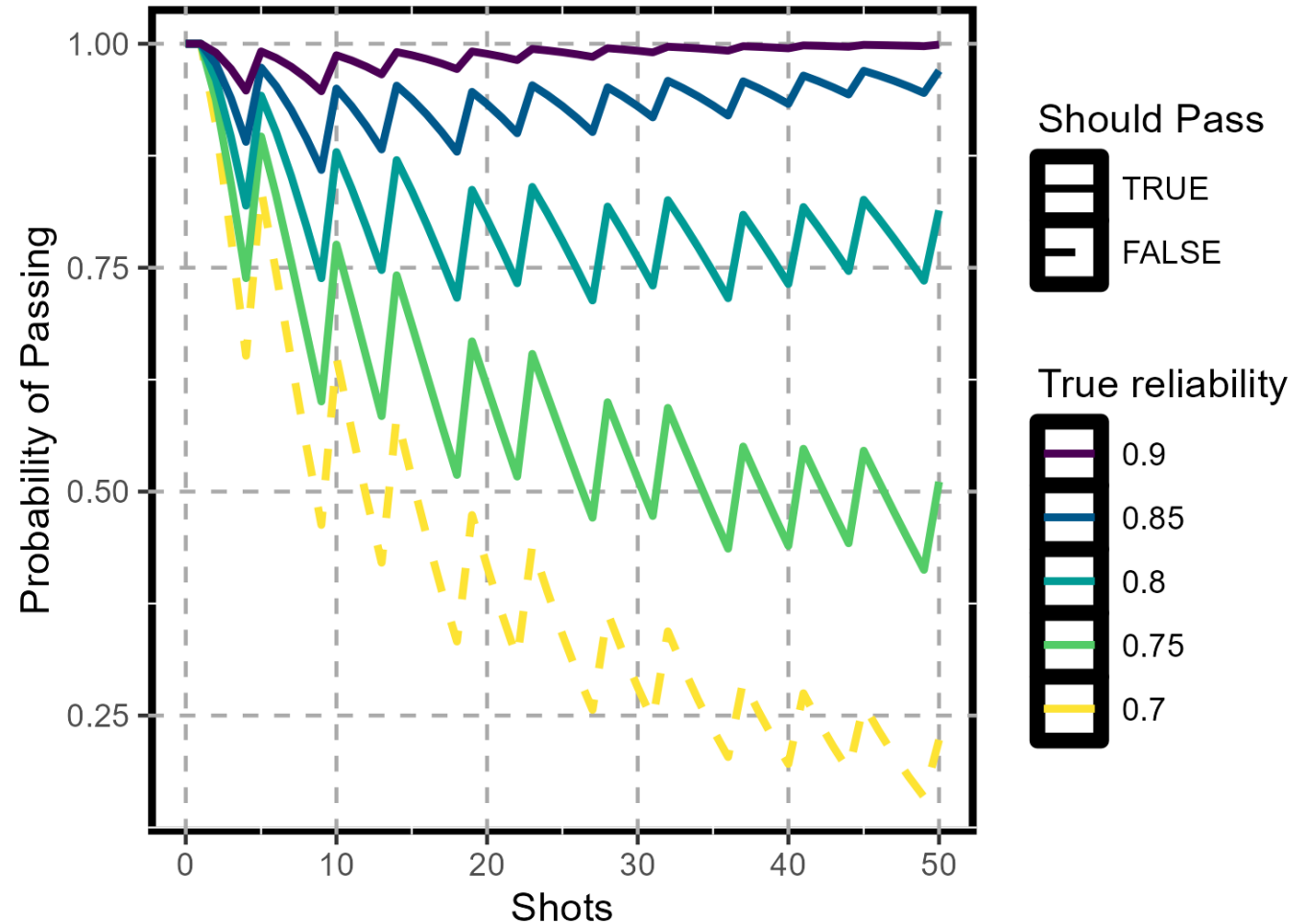
Sum probability (line thickness) above threshold from previous plot to get probability of passing
Increases occur when a new number of failures becomes “acceptable”



Previous results assumed a true reliability to compute binomial probability

Can compare probabilities of passing across true reliability values

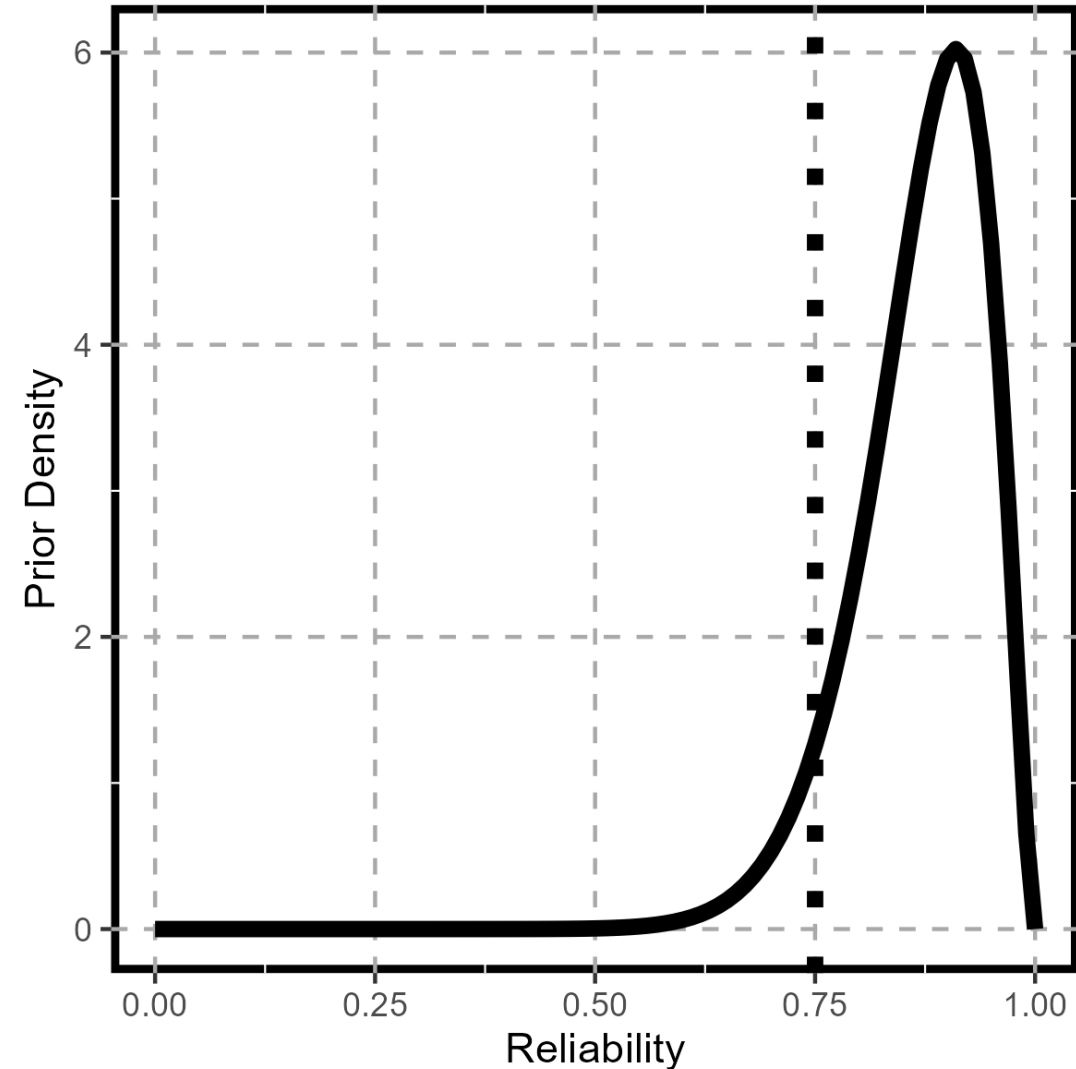
These can be used to size tests to discriminate between candidate reliabilities



In the Bayesian framework, we have beliefs about the value of the true reliability p . Can we use this information?

Integrate potential test results across true reliability p , weighted by the prior (right)

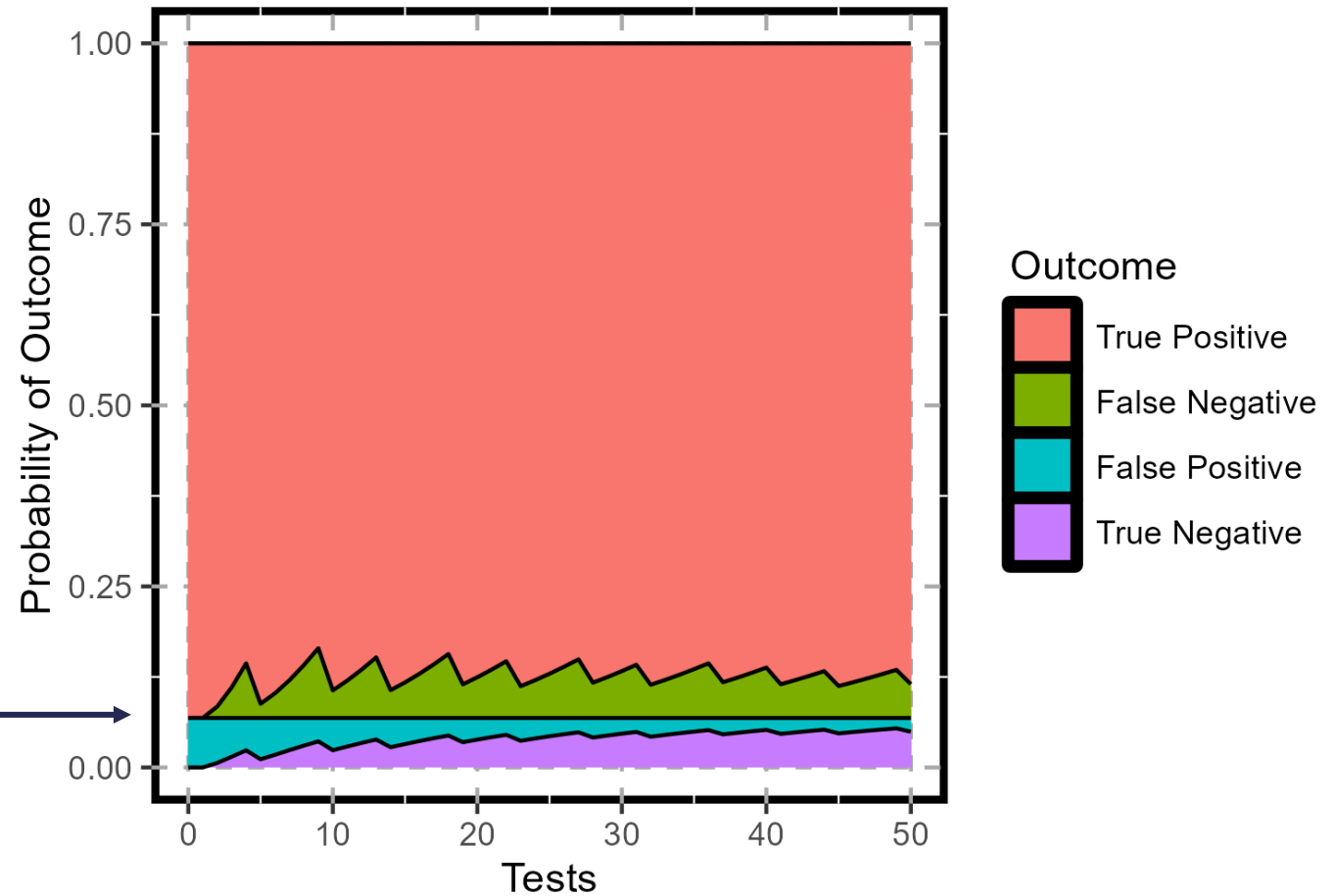
Gives the probability of test outcomes given our *current* knowledge



1. Make a fine grid of p values and evaluate the prior probability at each
2. Similarly, make a grid of number of tests (n) and number of successes (a) values
3. Evaluate:
 1. The probability of getting a successes for each values of n and p
 2. Whether the posterior associated with n tests and a successes meets the acceptance threshold
4. Use quadrature to approximately integrate the acceptance/rejections across p , weighted by the prior, considering the cases where p meets or fails to meet the threshold

Confusion matrix: Bin results by system (above/below threshold) and outcome (pass/fail)

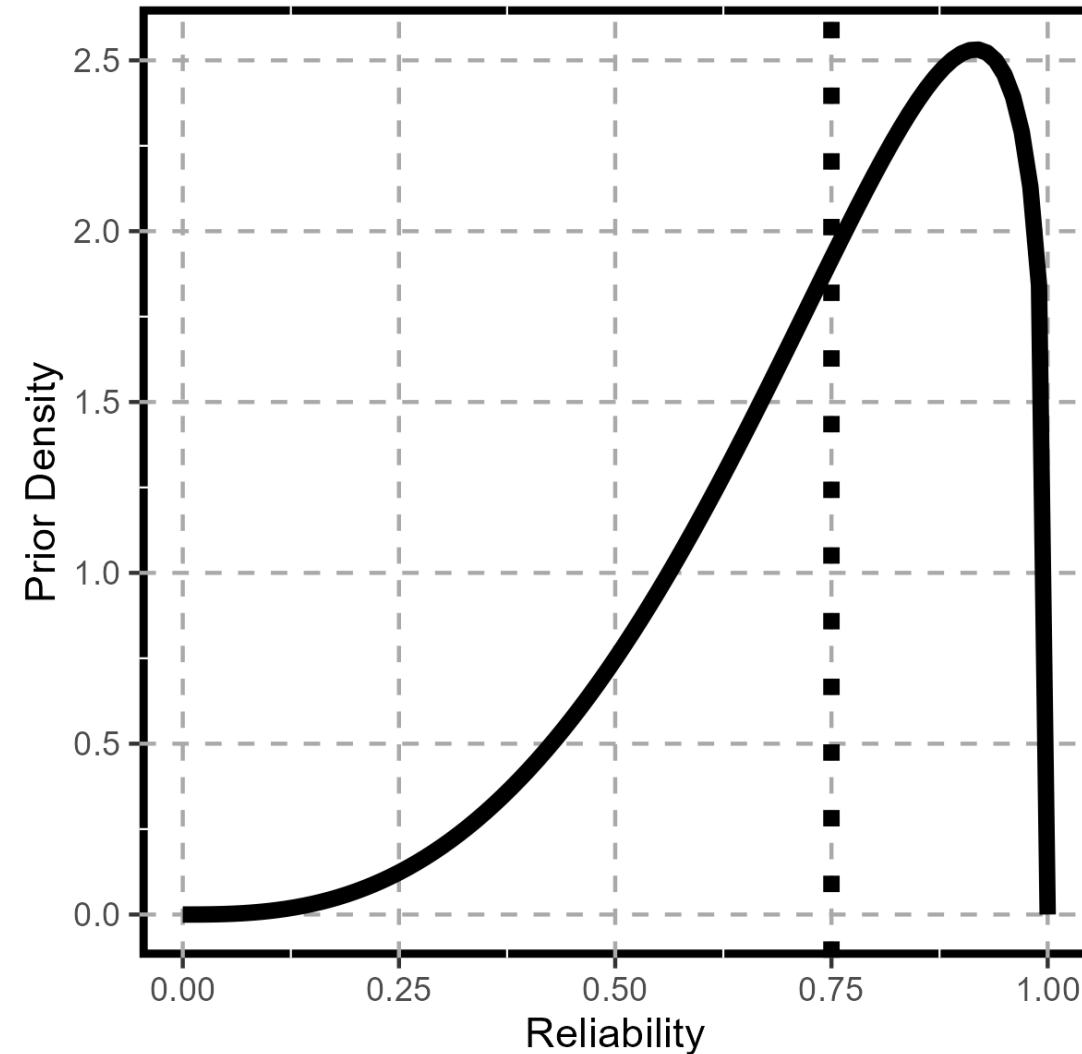
Horizontal line: Probability of whether system is good or bad is constant (given by prior) because we do not yet have additional information

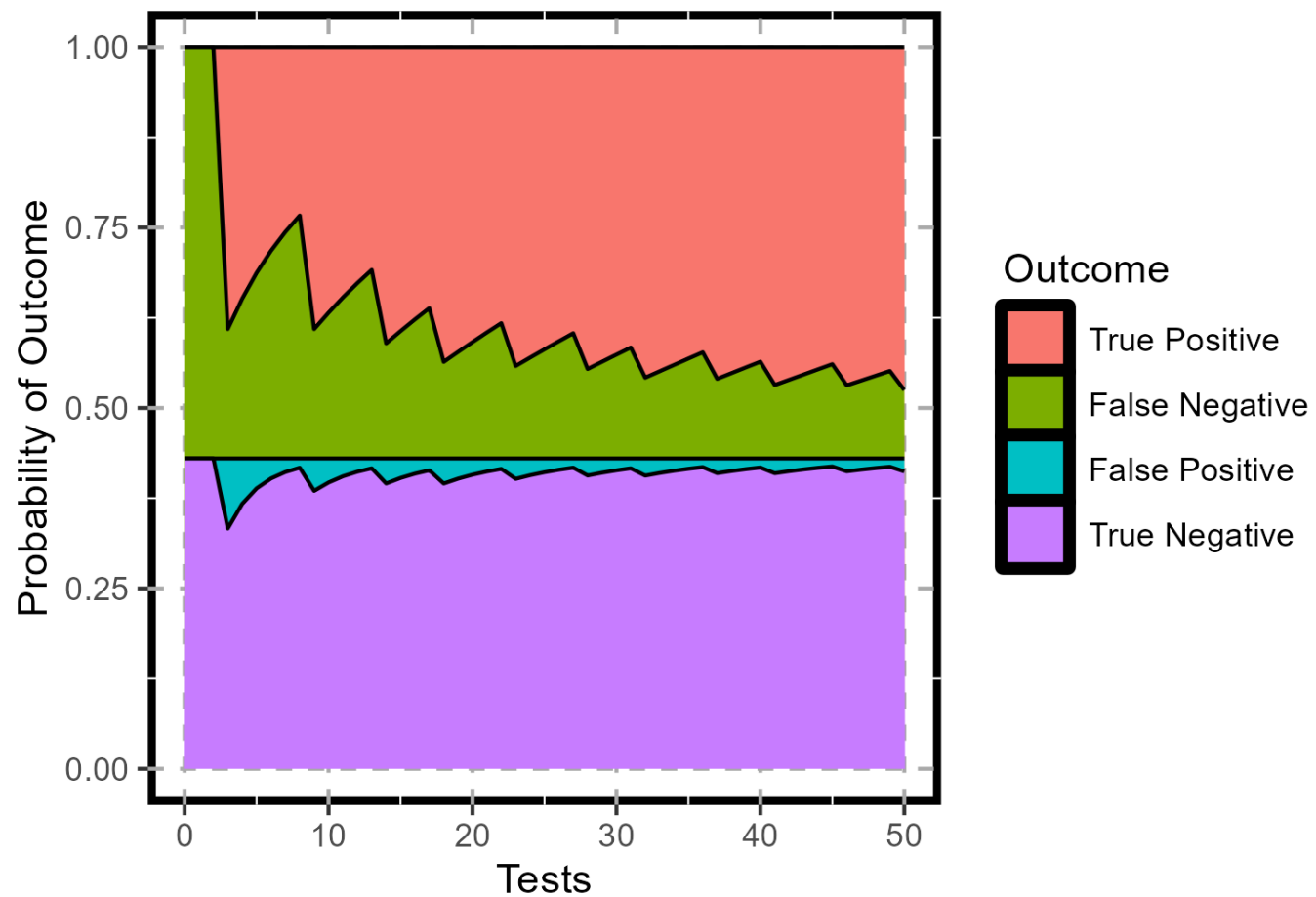


(Note: The prior for this example already meets the acceptance threshold, so negative values start at 0.)

Assume the following parameters, representing a more ambivalent prior:

Parameter	Value
Prior successes (a_0)	3.75
Prior failures (b_0)	1.25
Prior mean ($a_0 / (a_0 + b_0)$)	0.75
Prior sample size ($a_0 + b_0$)	5
Threshold	0.75
Lower credible interval	0.2







ACQUISITION INNOVATION
RESEARCH CENTER

ITERATE AS TEST RESULTS COME IN



ITERATE AS TEST RESULTS COME IN (WORSE LUCK)



The Bayesian approach provides a way to incorporate prior knowledge and test results

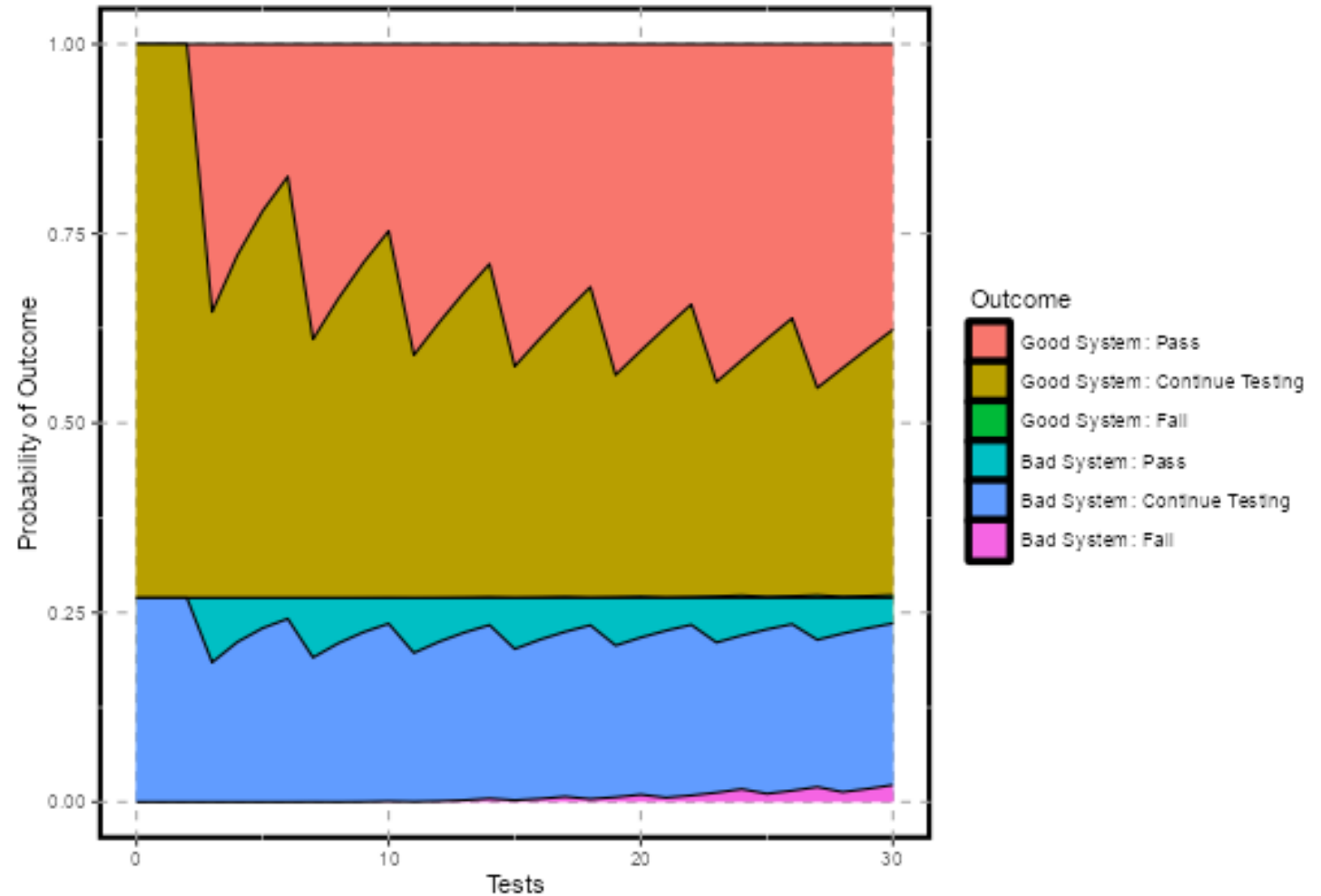
The prior predictive distribution provides full characterization of probable test outcomes given current knowledge

Future work: Add a “keep testing” category for when the posterior is not confident in either direction

The Bayesian approach provides a way to incorporate prior knowledge and test

The prior predictive distribution characterizes our prior knowledge

Future work: Add a “key” posterior is not confident



TOOL DEVELOPMENT



ACQUISITION INNOVATION
RESEARCH CENTER

Advanced statistical methods can be hard to implement

Tools allow users to:

- Conduct analysis using pre-written and validated code
 - Externally hosted: Without installation of special software or libraries
- Allow experimentation & exploration of assumptions
- Illustrate and disseminate methods
- Provide examples for training materials and tutorials

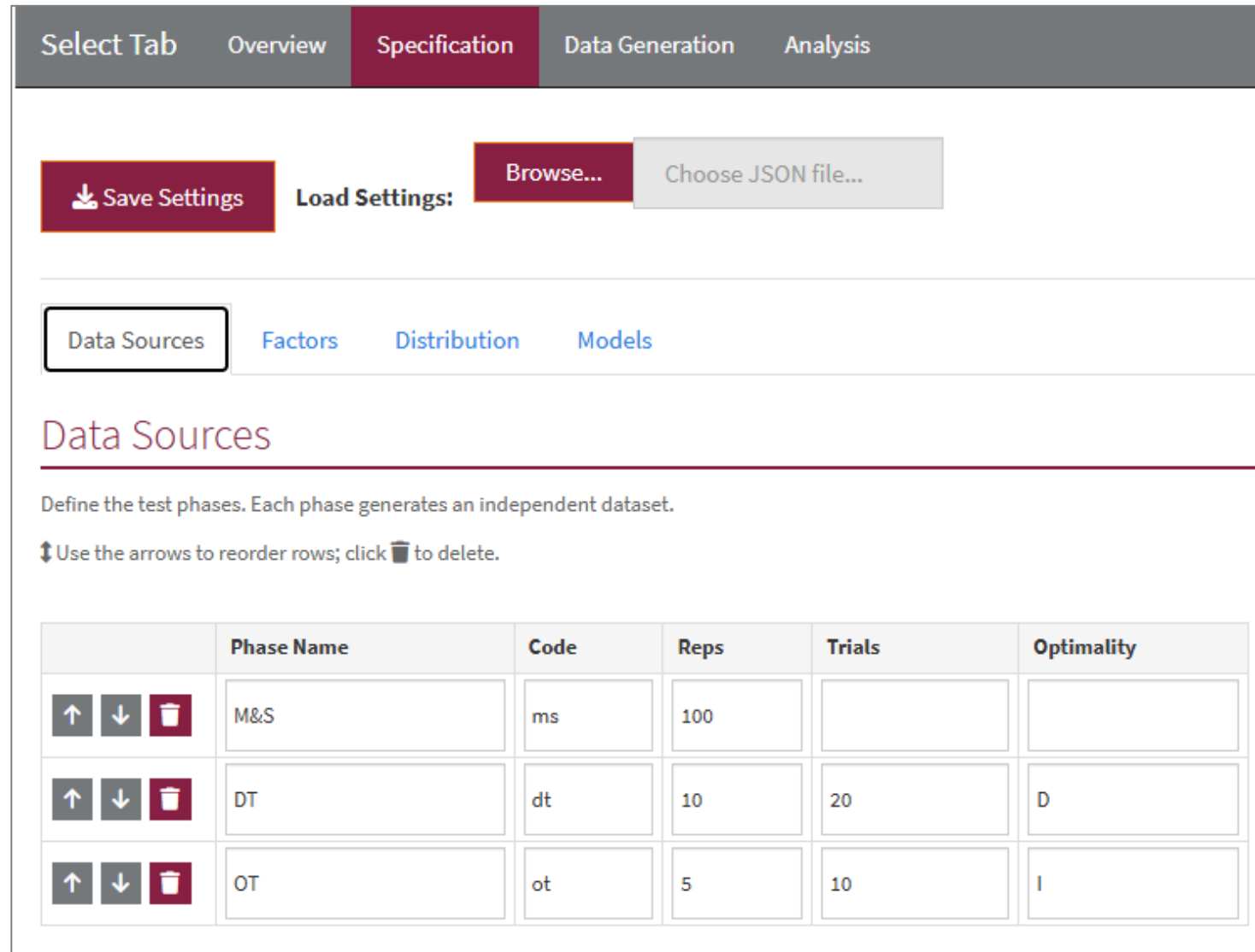
Key technology: Shiny

- Facilitates creation of graphical interfaces from R or Python
- Server version allows browser-based access for multiple users

Mimic DoW/D challenges

- Different numbers of test phases/data sources
- Varying data sizes, e.g., trials and reps by phase
- Evolving test factors
- Shifts/biases in test data (e.g., in M&S data)
- Different error/noise in measurements

Basic analysis (right now)



The screenshot shows the 'Specification' tab of the software interface. At the top, there are navigation tabs: 'Select Tab', 'Overview', 'Specification' (active), 'Data Generation', and 'Analysis'. Below the tabs, there are buttons for 'Save Settings', 'Load Settings', 'Browse...', and 'Choose JSON file...'. A sub-menu is open under 'Load Settings', showing 'Data Sources' (selected), 'Factors', 'Distribution', and 'Models'. The 'Data Sources' section is titled 'Data Sources' and includes the instruction: 'Define the test phases. Each phase generates an independent dataset. Use the arrows to reorder rows; click [trash icon] to delete.' Below this is a table with the following data:

	Phase Name	Code	Reps	Trials	Optimality
↑ ↓ [trash icon]	M&S	ms	100		
↑ ↓ [trash icon]	DT	dt	10	20	D
↑ ↓ [trash icon]	OT	ot	5	10	I

Mimic DoW/D challenges

- Different numbers of test phases/data sources
- Varying data sizes, e.g., trials and reps by phase
- Evolving test factors
- Shifts/biases in test data (e.g., in M&S data)
- Different error/noise in measurements

Basic analysis (right now)

Select Tab Overview Specification Data Generation Analysis

Data Generation

Generates synthetic data using the specification defined in the previous tab. Each phase is generated independently and combined into a single dataset.

Options

Random Seed

42

▶ Generate Data

⬇ Download CSV

✔ Done: 850 rows.

Data Preview

Show 10 entries

const	B	D	mean	
1	-1	-1	81	27.273475442
1	1	-1	69	19.848719048
1	-1	0	85	107.28783729
1	1	0	73	35.505004489
1	-1	1	89	66.884107827
1	1	1	77	70.353757111
1	-1	-1	81	50.536640020
1	1	-1	69	110.24433204
1	-1	0	85	62.164369887

Capabilities:

- Upload data
- Select priors
- Conduct DT inference
- Downweight
- Conduct OT Inference

Challenge:

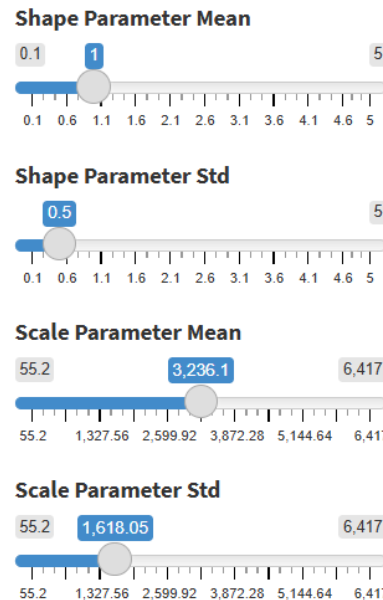
- Inference (MCMC) breaks on free versions of Shiny servers

Priors and Implied Mean Reliability

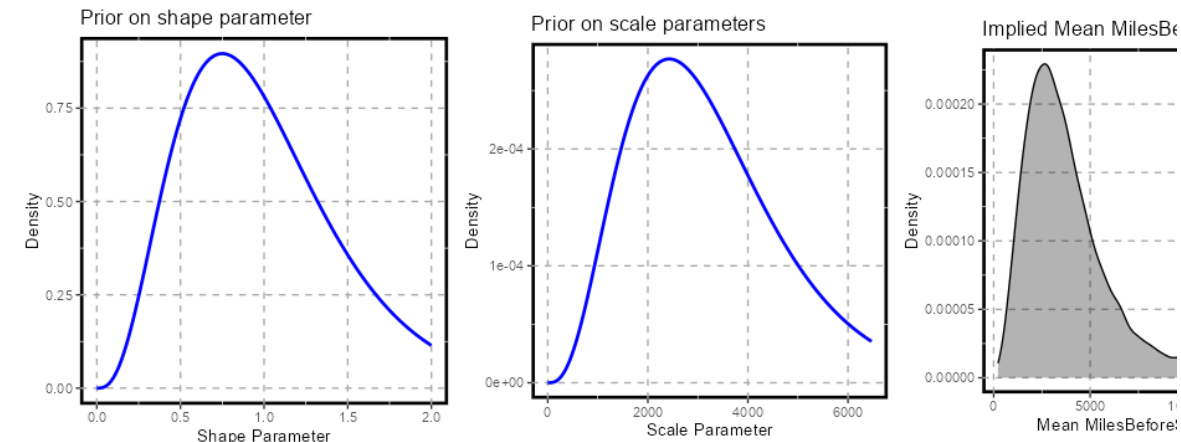
In this tab, you can explore the impact of adjusting the mean and standard deviation of the gamma priors on the scale and shape parameters.

Use the text fields below to specify the minimum value and maximum value for the Scale Parameter Mean and Scale Parameter Standard Deviation values that you would like to explore. If you choose your own slider ranges, press the 'Create Scale Parameter Mean and Std Sliders' button to update the sliders.

Once you have adjusted the prior parameters as you wish, click the 'Select these Priors' button and proceed to the next tab for DT inference.



Select these Priors



Implied Mean MilesBeforeSystemAbort Summary Statistics:

Min: 224

1st Quantile: 2314

Functionality

- Upload data
 - Arbitrary number of test phases
- Compute posteriors
- Evaluate test plans
- Compute prior predictive

Select Tab Overview Dataset Inference **Planning** Prior Predictive

Binary Test Planning

This app provides a Bayesian approach to planning for binary tests. Given a set of true success rates (the probability of a single trial yielding a success), this computes the probability of the system passing the overall test, where a pass is defined as the Bayesian lower credible interval meeting a specified threshold. For example, if the lower credible bound is 0.8 and the threshold is 0.85, then the system will pass the test if the Bayesian posterior distribution gives 80% probability the true success rate is at least 0.85. If the user specifies a single true success rate, the app also produces a series of lower-level plots that describe how the results are calculated, providing additional details about why the results look as they do.

Pass Threshold

Credibility Level

Maximum Number of Tests

Enter prior values:

Prior 1

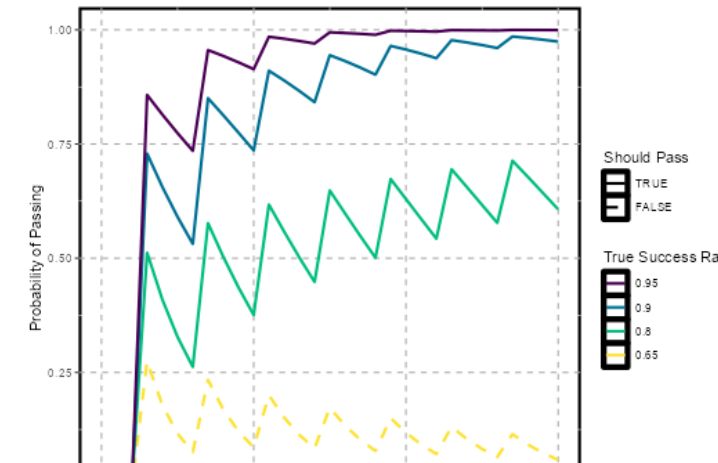
Prior 2

Reset Priors

Enter True Success Rate Values (Multiple rates is optional)

Probability of Passing by Number of Tests and Success Rate

This plot shows the probability of the system passing (meeting the minimum credible interval) by tests (x-axis) and true success rate (color). Success rate values that exceed the threshold (i.e., that the test) are shown as solid lines, while success rate values that fail to meet the requirements are dashed lines.



Download: <https://github.com/krometis/dataworks2026>

Hosted:

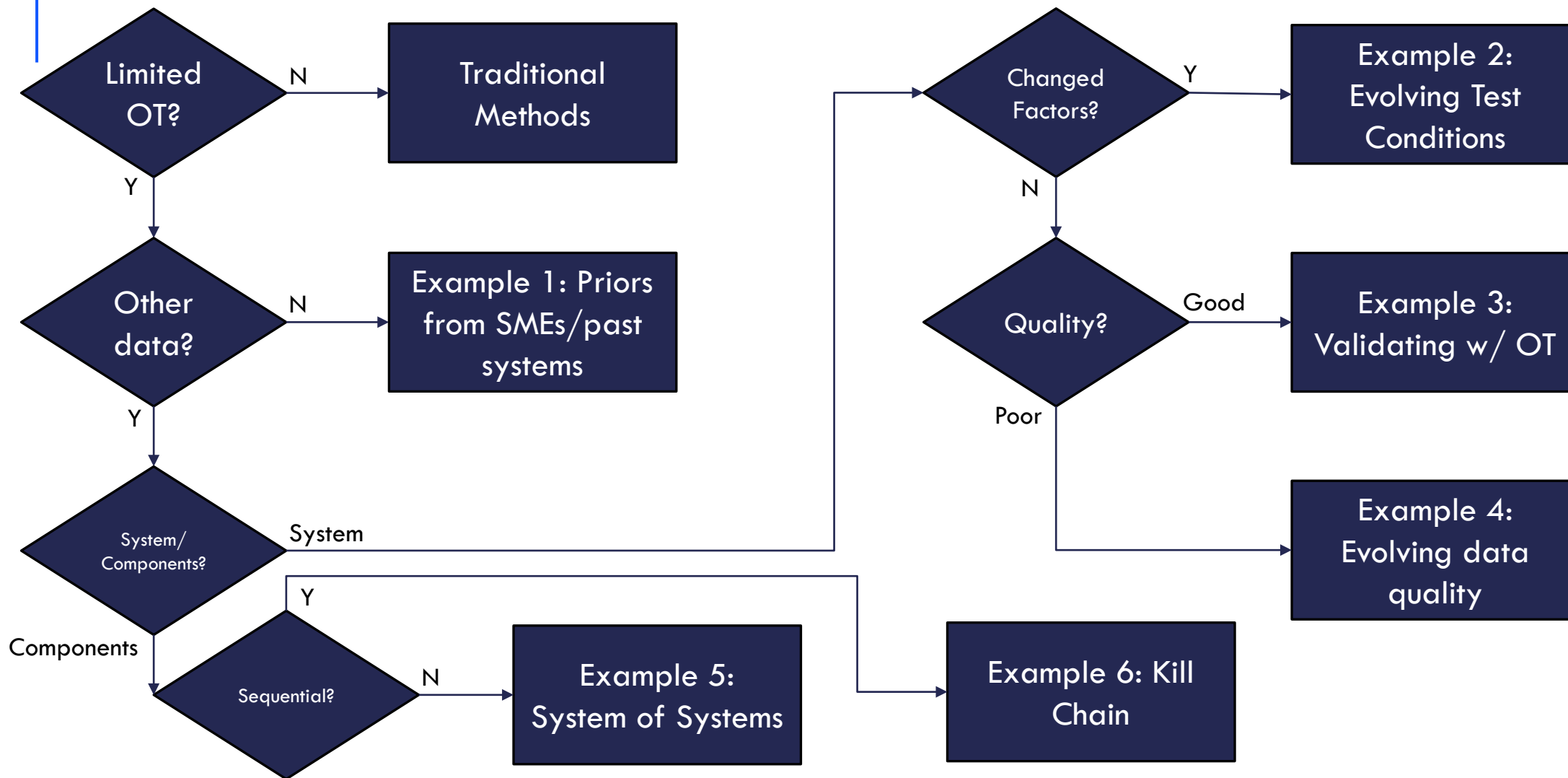
- Data Generation: <https://krometis.shinyapps.io/datagen/>
- Reliability: <https://krometis.shinyapps.io/reliability/>
 - Warning: Analysis will crash!
- Binary: https://krometis.shinyapps.io/binary_multi/

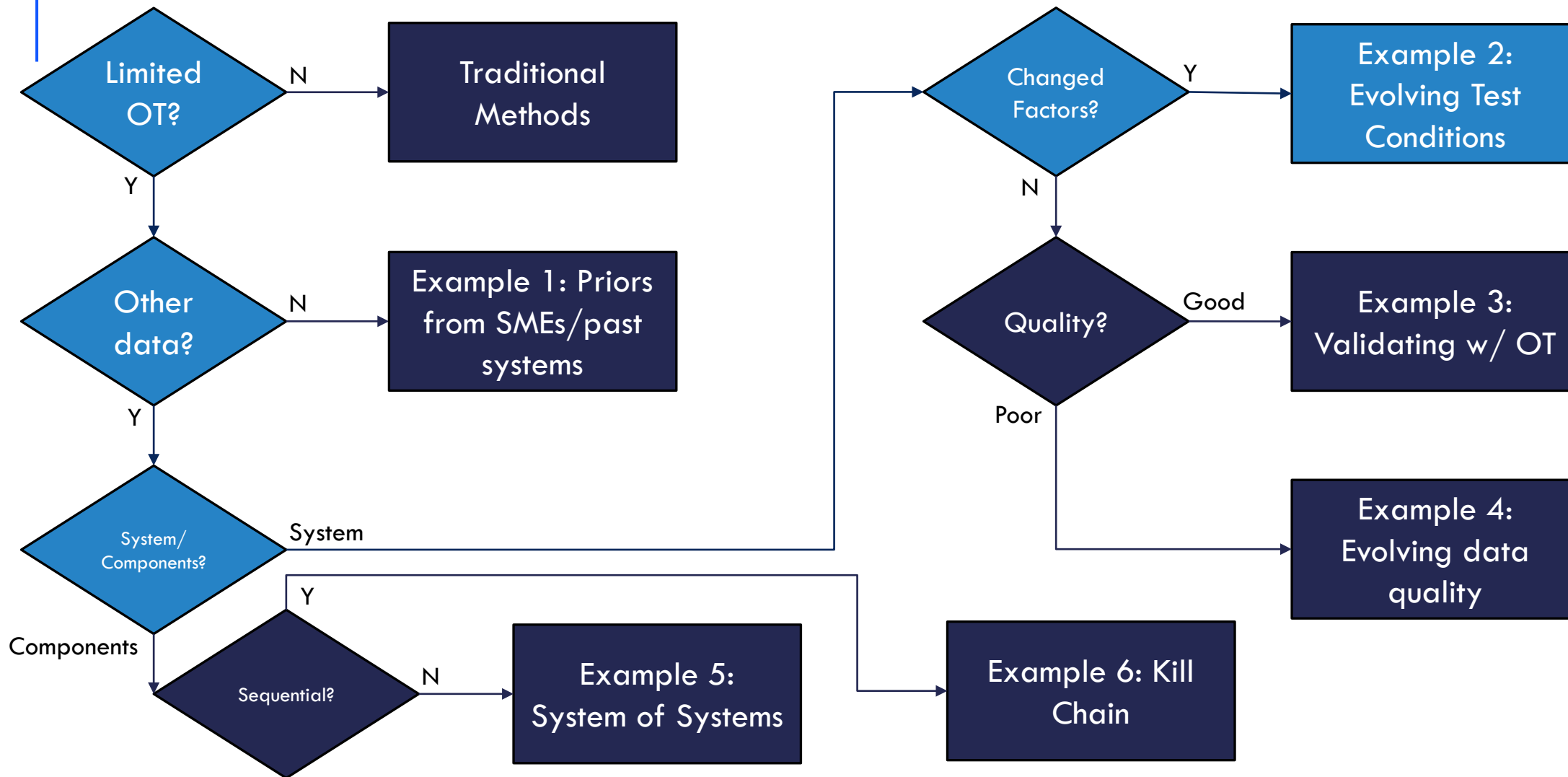


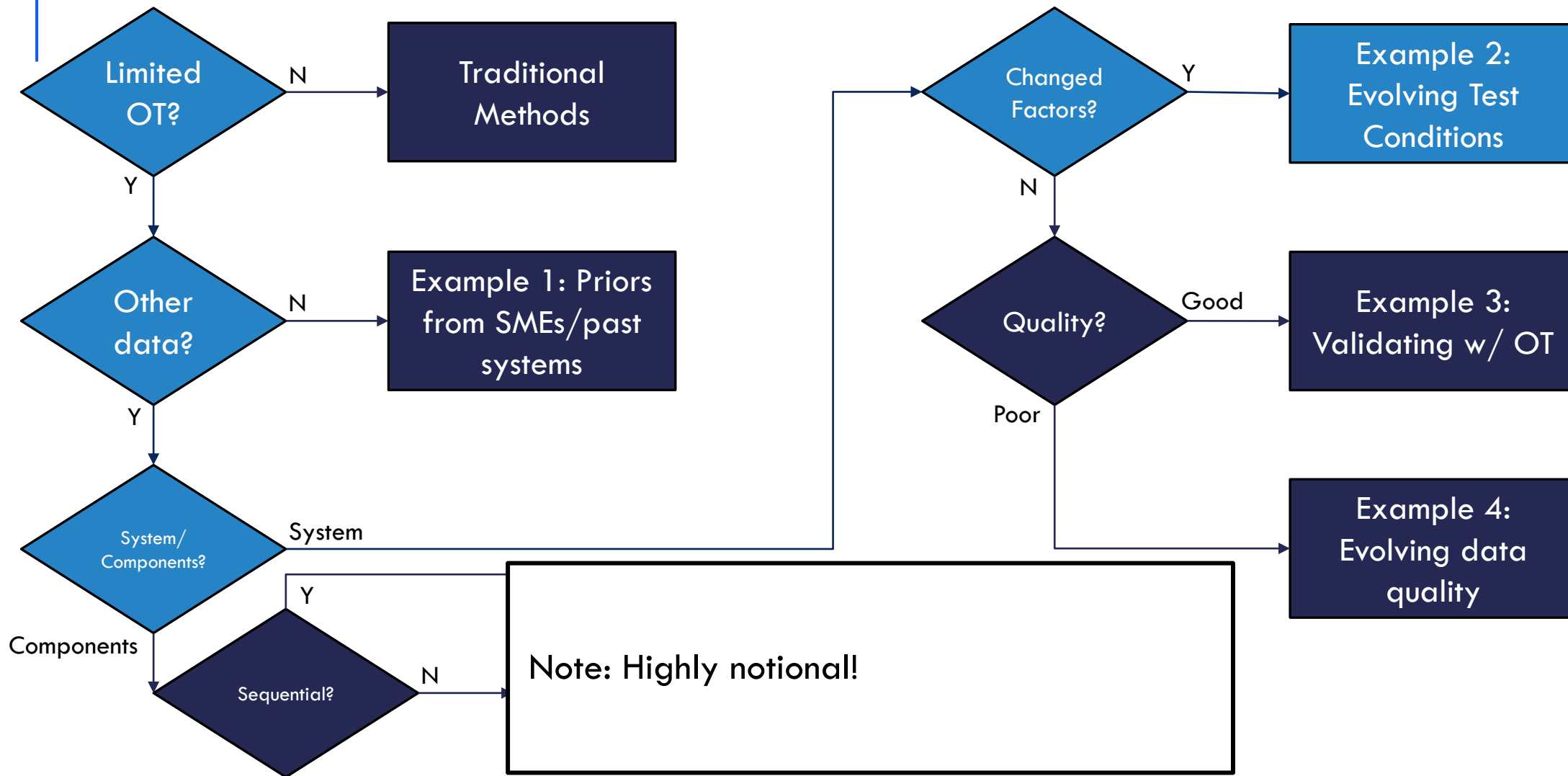
NEXT STEPS



ACQUISITION INNOVATION
RESEARCH CENTER







Concerns/questions about methods

Barriers to implementation & adoption