**IDA**

INSTITUTE FOR DEFENSE ANALYSES

# Space-Filling Experimental Design and Surrogate Models for U.S. Department of Defense Modeling and Simulation Evaluation

John T. Haman, Project Leader

Curtis G. Miller

**IDA**

The Institute for Defense Analyses is a nonprofit corporation that operates three Federally Funded Research and Development Centers. Its mission is to answer the most challenging U.S. security and science policy questions with objective analysis, leveraging extraordinary scientific, technical, and analytic expertise.

Rigorous Analysis │ Trusted Expertise │ Service to the Nation

INSTITUTE FOR DEFENSE ANALYSES

# Space-Filling Experimental Design and Surrogate Models for U.S. Department of Defense Modeling and Simulation Evaluation

John T. Haman, Project Leader

Curtis G. Miller

# Executive Summary

Operational testing (OT) provides important information on the capability of equipment acquired by the United States Department of Defense (DoD). Operational testing has been a vital component of the acquisition process since Congress reformed OT conduct with the Department of Defense Authorization Act of 1984. Computer modeling and simulation (M&S) can help address the small sample size problem in live, operational testing, but M&S must be validated in order for stakeholders to believe M&S makes meaningful predictions about real-world outcomes.

The Institute for Defense Analyses (IDA) has published documents providing guidance to the OT community on statistical validation of M&S, including a handbook and papers on the use of space-filling designs and statistical surrogates (also known as metamodels). However, important questions remain that need to be answered in order to best validate M&S.

In particular, operational testers need better recommendations on sample size selection, determining the number of replicates in a design (if any), validating with small real-world sample sizes, and incorporating statistical surrogates into hypothesis tests that help determine whether M&S outcomes match real-world outcomes or not.

We hope the larger statistical community can help contribute answers to these questions. This presentation was given at the Joint Statistical Meeting 2023 in Toronto, Canada.

# Space-Filling Experimental Design and Surrogate Models for U.S. Department of Defense Modeling and Simulation Evaluation

Dr. Curtis Miller

August 9, 2023

## Institute for Defense Analyses

730 East Glebe Road ● Alexandria, Virginia 22305

# We want statisticians to appreciate...

… why operational testing matters

… why computer modeling and simulation matters to operational testing

… the challenges in planning computer modeling and simulation studies

… the challenges in statistically validating a computer model

Why does operational testing matter?

# Operational testing studies DOD system effectiveness in warfighting conditions



DOD – U.S. Department of Defense

# OT results matter to both Congress and the warfighter



OPERATIONAL TESTING: ENSURING BETTER WEAPONS FOR OUR TROOPS

4. G 74/9: S. HRG. 103-568

S. Hrg. 103-568

rational Testing: Ensuring Bette...

**HEARING**

BEFORE THE

SUBCOMMITTEE ON FEDERAL SERVICES, POST OFFICE, AND CIVIL SERVICE

OF THE

COMMITTEE ON GOVERNMENTAL AFFAIRS

UNITED STATES SENATE

ONE HUNDRED THIRD CONGRESS

SECOND SESSION

MARCH 22, 1994

Printed for the use of the Committee on Governmental Aff

AUG 4 1994

U.S. GOVERNMENT PRINTING OFFICE

WASHINGTON : 1994

78-061 cc

For sale by the U.S. Government Printing Office
Superintendent of Documents, Congressional Sales Office, Washington, DC 20
ISBN 0-16-044450-0

**Sen. D. Pryor (D-AK)       Sen. W. Roth (R-DE)**

plished.

Senator PRYOR. Now, you know, testing is not a very—I hate to use the word "sexy," but it is not a very high priority item it seems like today. There are not a lot of people interested in this, except the people in the battlefield. They are going to be interested in the outcome of all of this, and they are certainly going to be interested in the outcome of S. 1587. They will also be interested in whether this, as we call it, Mack truck amendment, remains in the bill.

Making certain that these weapons work is a very, very important part of our military preparedness. Whoever our enemies might be today, we had better believe that their intelligence knows if these weapons work or don't work, and I think that we have got to keep that in mind throughout this whole process.

If you would please continue.

Mrs. PRESTON. Let me quickly go through some of the other

OT – Operational Testing

**IDA** | 4

# OT&E should reduce risk and uncertainty regarding the performance of systems in wartime

**United States General Accounting Office**

**GAO**

Report to the Honorable
William V. Roth and the Honorable
Charles E. Grassley, U.S. Senate

October 1997

## TEST AND EVALUATION

### Impact of DOD's Office of the Director of Operational Test and Evaluation

number of unknowns prior to the decision to begin full production, while program and service officials typically sought less testing and were willing to accept greater risk when making production decisions. The additional testing DOT&E advocated, often over the objections of service testers, served to meet the underlying objectives of operational testing—to reduce the uncertainty and risk that systems entering full-rate production would not fulfill their requirements.

GAO/NSIAD-98-22

OT&E – Operational Test and Evaluation

Why does computer modeling and simulation matter to OT?

# Small data problems are alive and well in operational testing
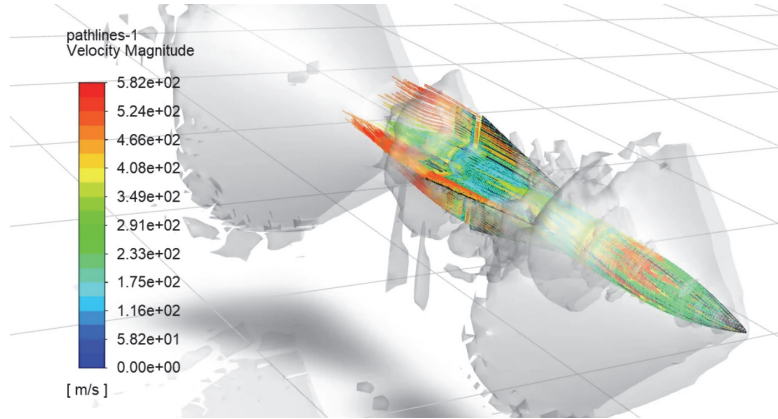
*SM-3 Blk IIA*         *Next Generation Interceptor*         *PAC-3*



**Missile tests can cost $10Ms to $100Ms
per shot**

*https://missiledefenseadvocacy.org/missile-defense-systems-2/missile-defense-systems/missile-interceptors-by-cost/*

# Small data problems are alive and well in operational testing

*SM-3 Blk IIA*          *Next Generation Interceptor*          *PAC-3*



**COST AND ASSET AVAILABILITY GENERATE DEMAND FOR M&S**

*Missile tests can cost $10Ms to $100Ms per shot*

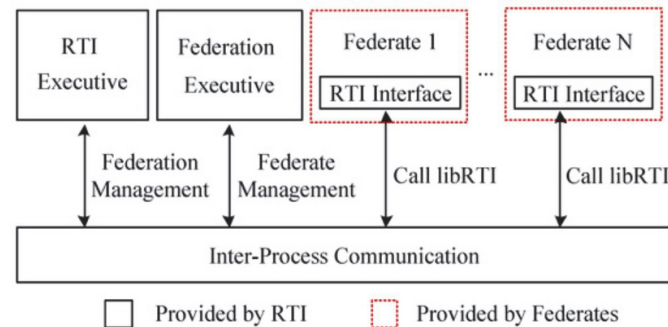# Modeling and simulation comes in a variety of flavors with unique statistical considerations



**Digital Simulation**

https://www.reddit.com/r/dcsworld/comments/mqmvy3/lowquality_steadystate_cfd_simulation_of_an_aim54
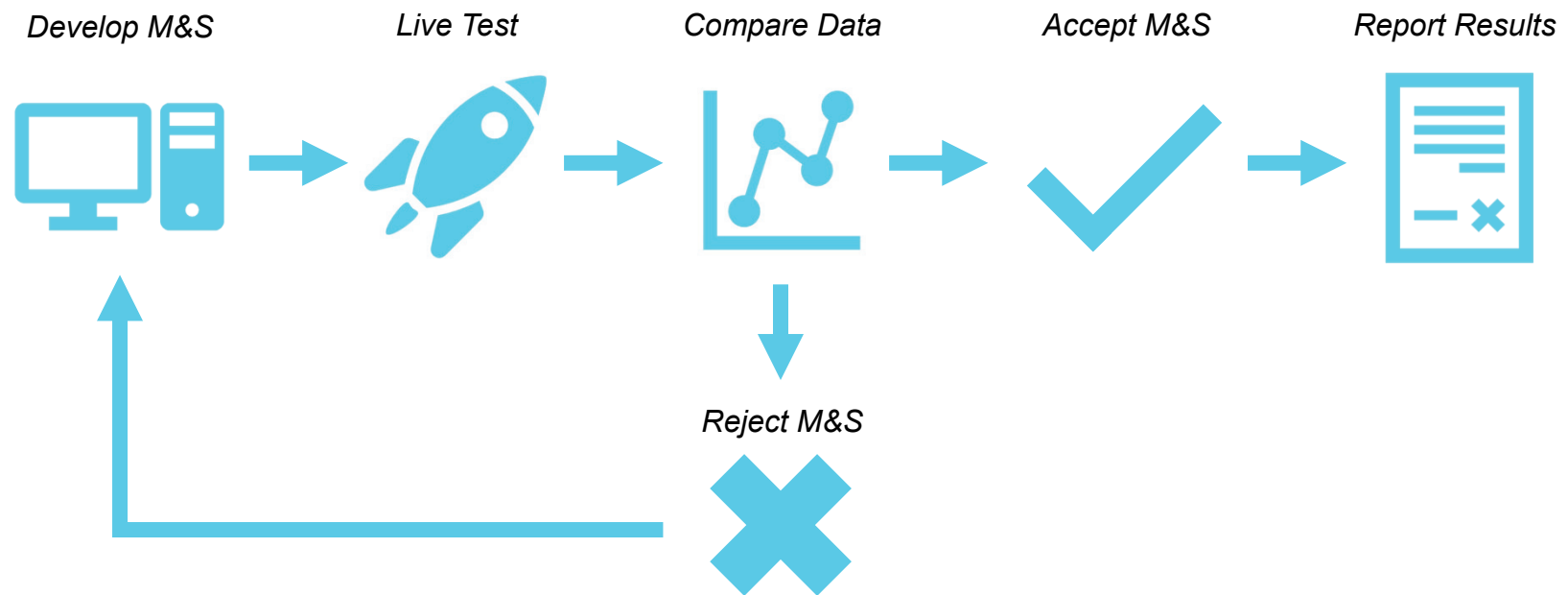


**Hardware in the Loop**

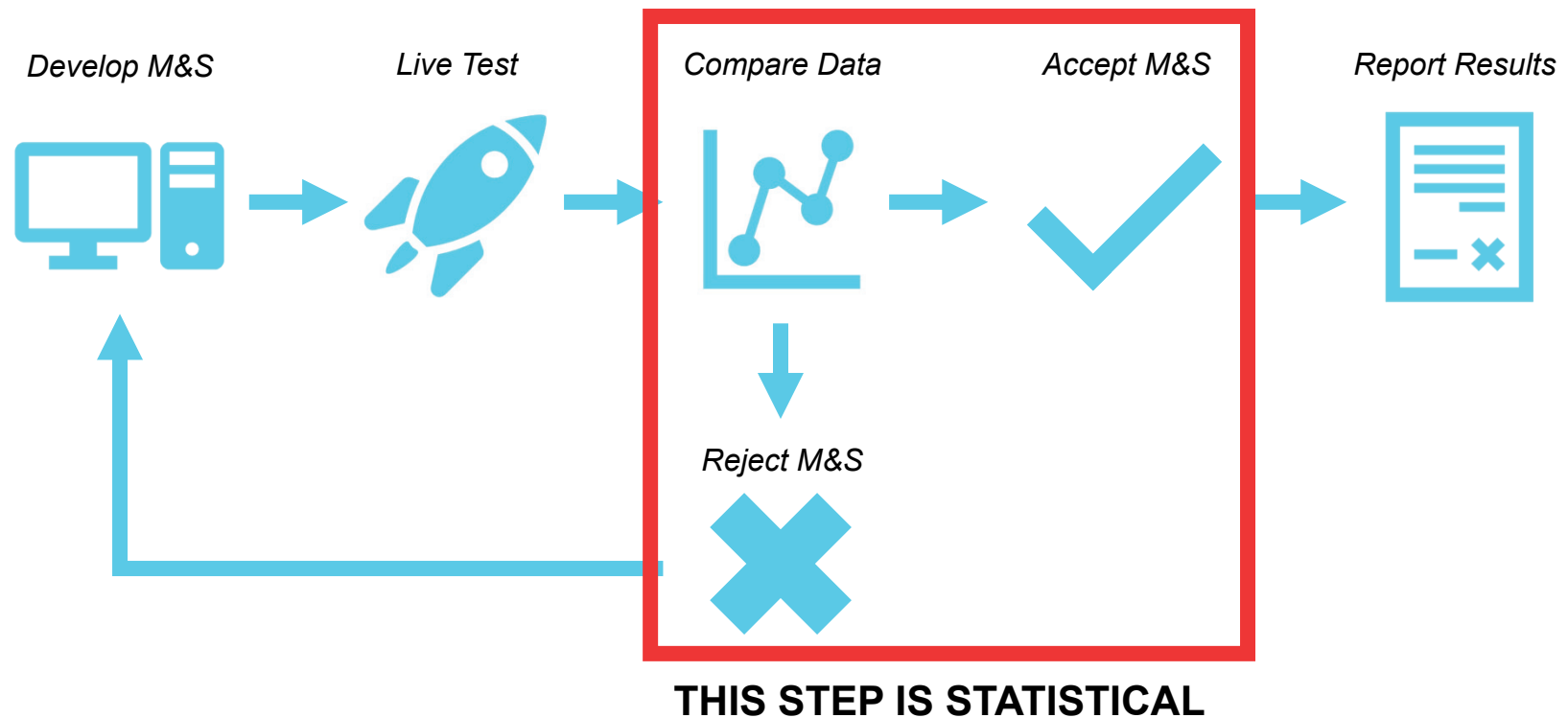https://testscience.org



**Federation**

https://doi.org/10.3390/electronics9030540

M&S – Modeling and Simulation

# If M&S outputs will be used for predicting OT results, we must compare M&S outputs to live test data



M&S – Modeling and Simulation; OT – Operational Testing

# If M&S outputs will be used for predicting OT results, we must compare M&S outputs to live test data



Develop M&S → Live Test → **Compare Data** → **Accept M&S** → Report Results

Reject M&S

**THIS STEP IS STATISTICAL**

M&S – Modeling and Simulation; OT – Operational Testing

# If M&S outputs will be used for predicting OT results, we must compare M&S outputs to live test data

Develop M&S      Live Test      Compare Data      Accept M&S      Report Results

**HOW CAN STATISTICS HELP MAKE THE RIGHT DECISION ON WHETHER TO TRUST M&S PREDICTIONS?**

**THIS STEP IS STATISTICAL**

M&S – Modeling and Simulation; OT – Operational Testing

What are the challenges in planning computer M&S studies?

# IDA publications introduce and recommend M&S DOE best practices

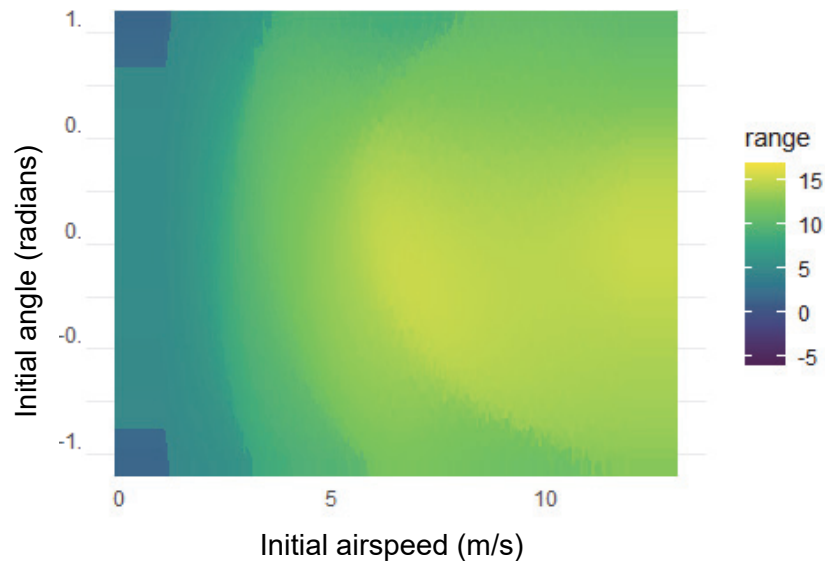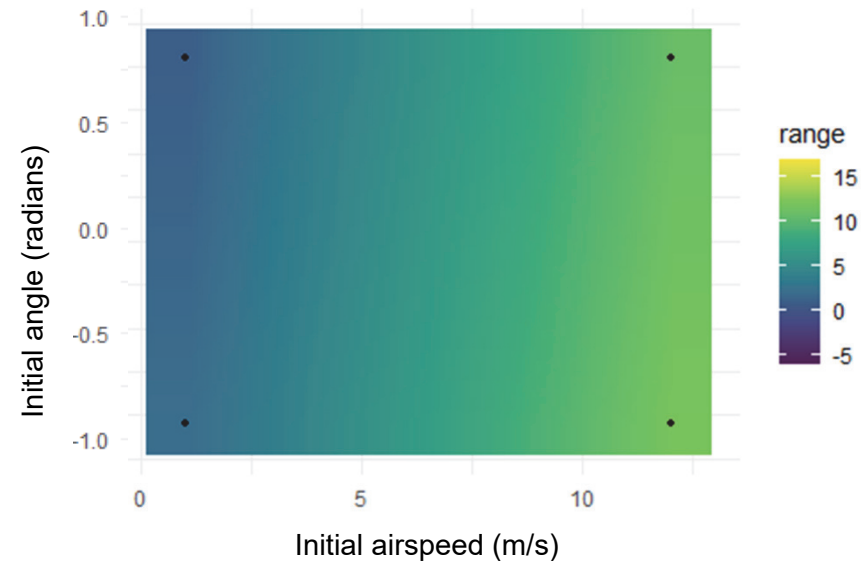# Space-filling design of experiments helps recover more trends in simulation outputs

A factorial design with a simple linear model fitted will not accurately describe the M&S system's behavior.



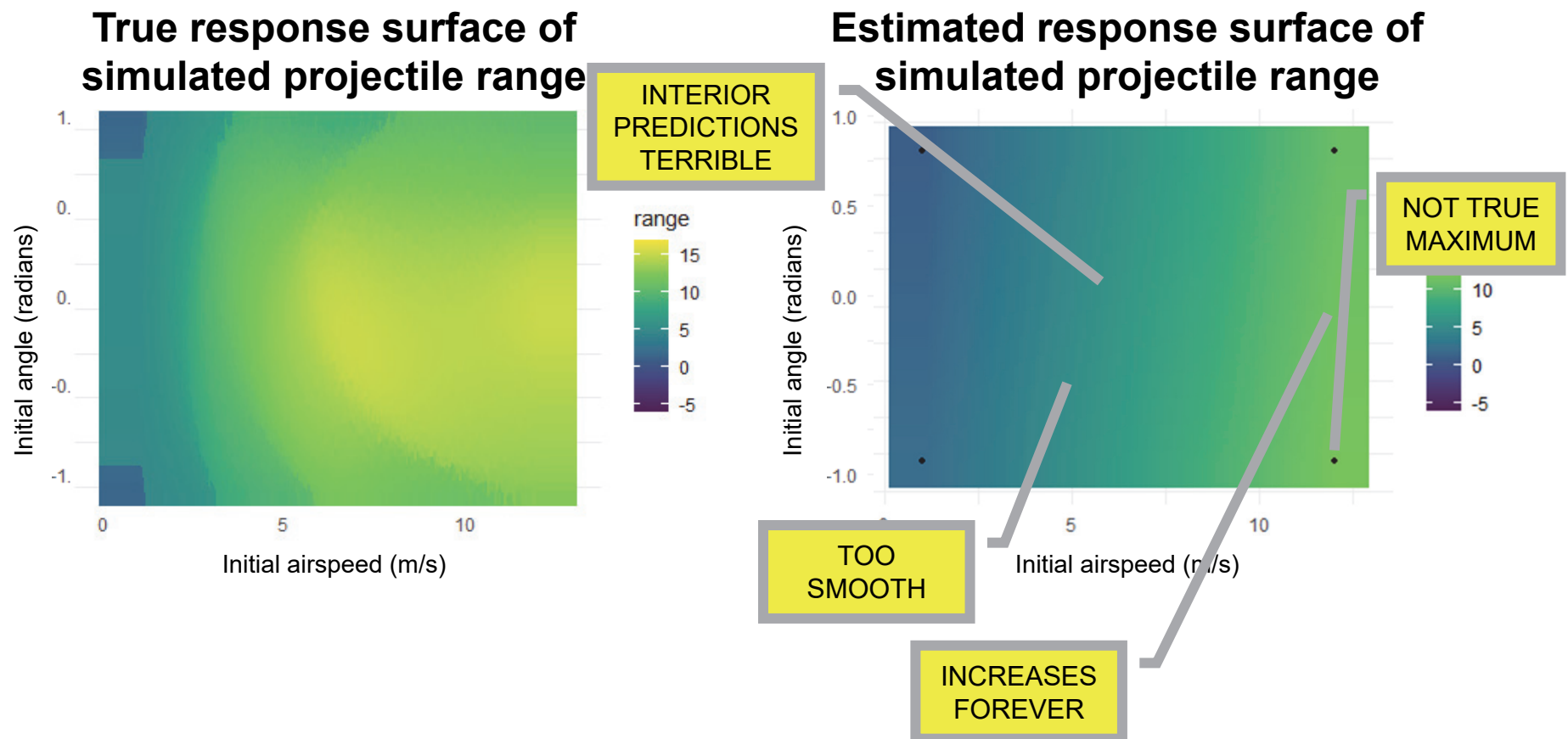**True response surface of simulated projectile range**

**Estimated response surface of simulated projectile range**

# Space-filling design of experiments helps recover more trends in simulation outputs
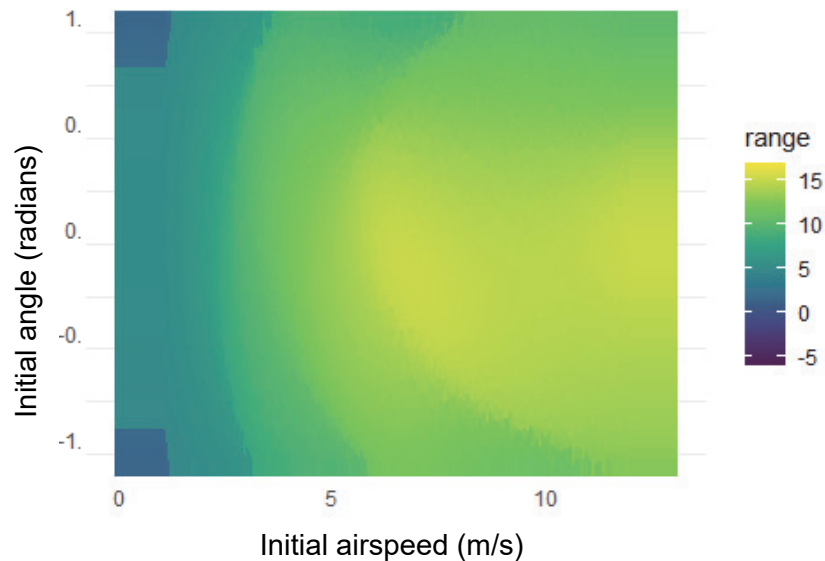
A factorial design with a simple linear model fitted will not accurately describe the M&S system's behavior.



True response surface of simulated projectile range

Estimated response surface of simulated projectile range

INTERIOR PREDICTIONS TERRIBLE
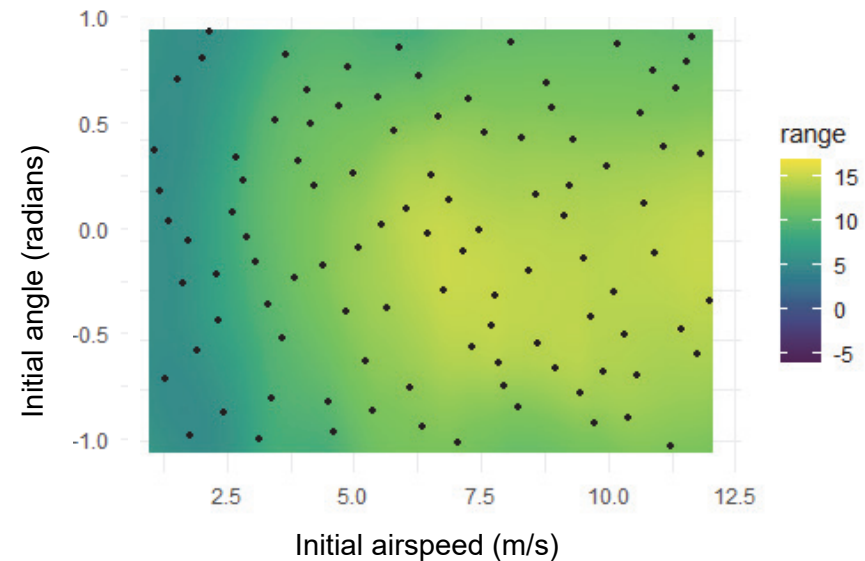
NOT TRUE MAXIMUM

TOO SMOOTH

INCREASES FOREVER

# Space-filling design of experiments helps recover more trends in simulation outputs

Analyzing the flights with a Gaussian Process model via a Space-Filling Design yields a good approximation to simulation output.

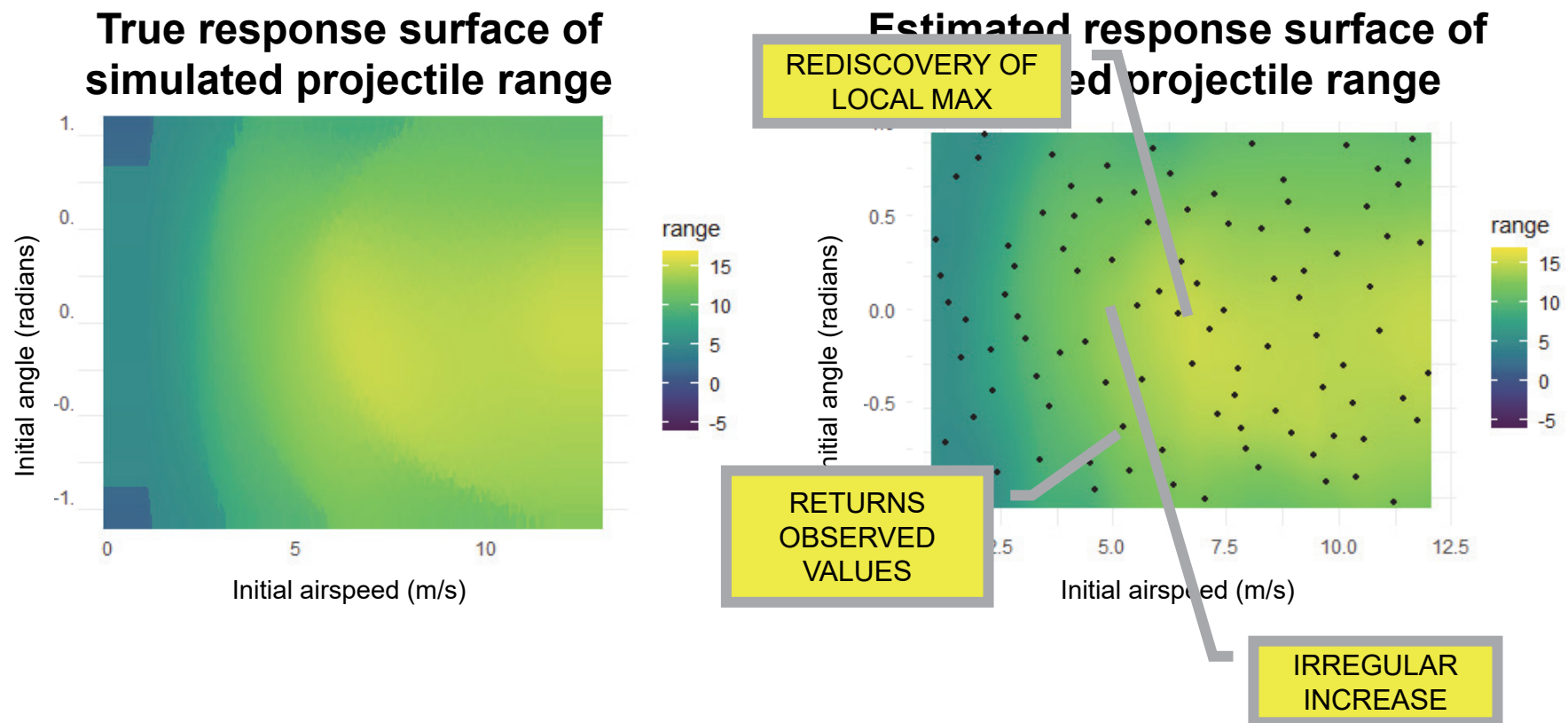**True response surface of simulated projectile range**



**Estimated response surface of simulated projectile range**

# Space-filling design of experiments helps recover more trends in simulation outputs
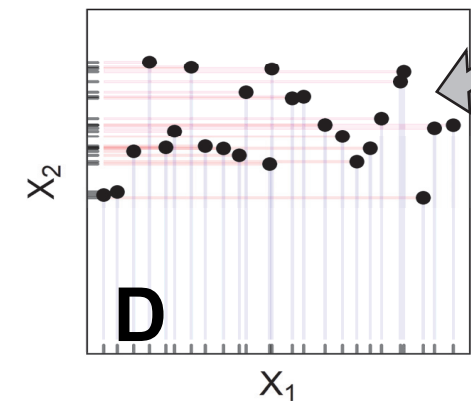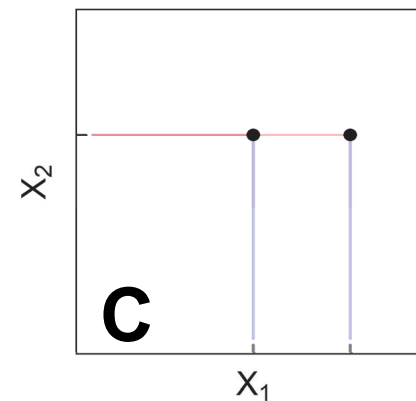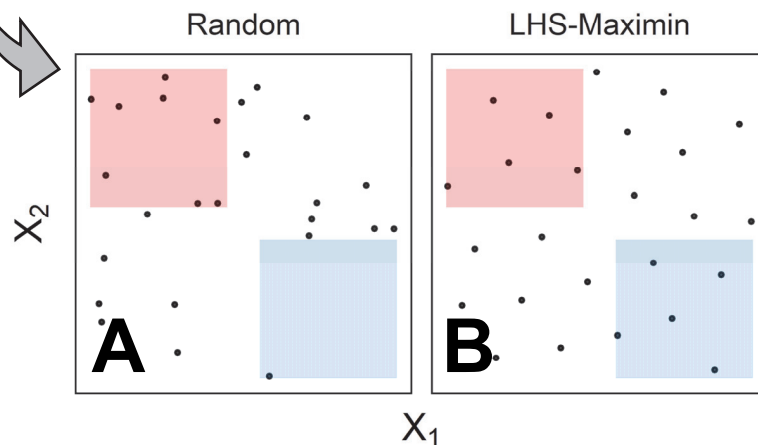
Analyzing the flights with a Gaussian Process model via a Space-Filling Design yields a good approximation to simulation output.



**True response surface of simulated projectile range**

**Estimated response surface of simulated projectile range**

REDISCOVERY OF LOCAL MAX

RETURNS OBSERVED VALUES

IRREGULAR INCREASE

# Just like with classical DOE, there are quantitative ways to evaluate a specific design

Many criteria exist, but it is particularly important that an SFD satisfy the following three criteria in order to be useful:
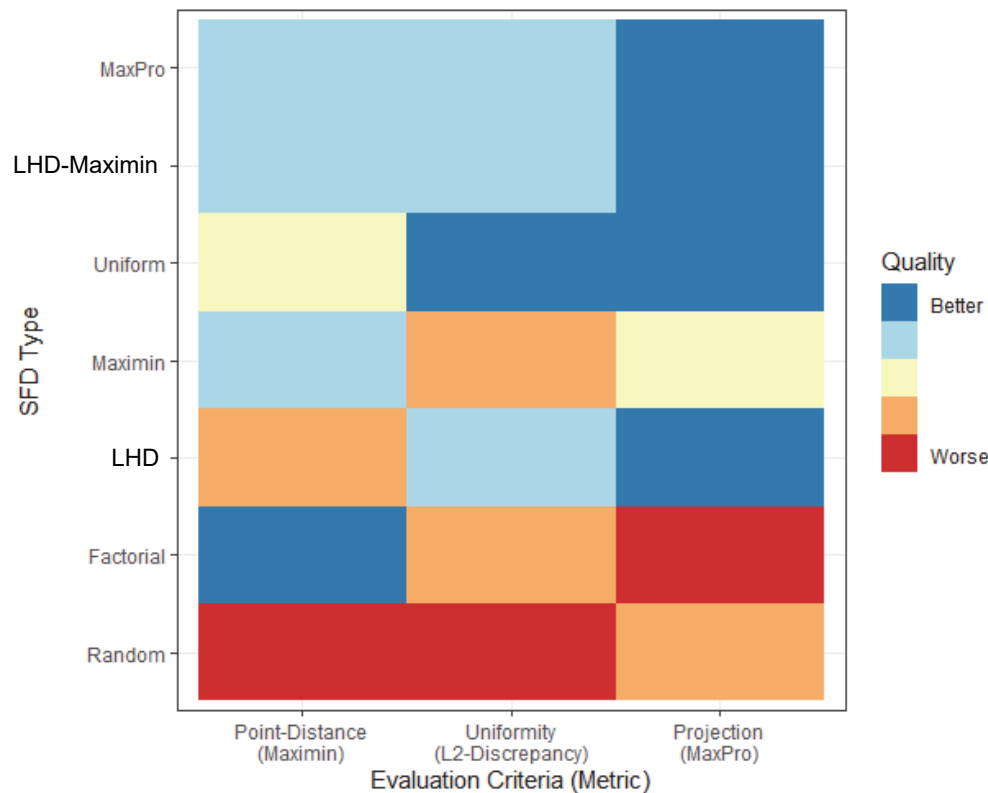
- Point-distance: Samples are placed as far apart from each other as possible. [Maximin]

- Uniformity: All regions of the design space are equally well-represented. [Center $L^2$ Discrepancy]

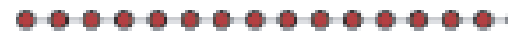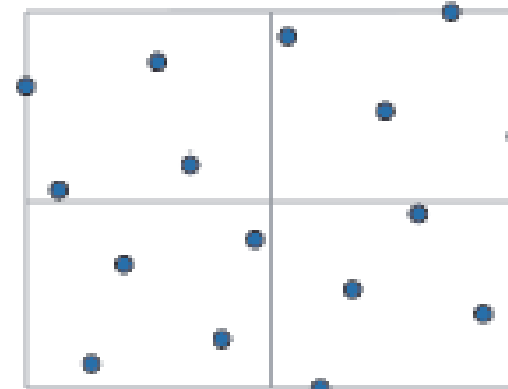- Projection: The design is robust to variables being collapsed. [MaxPro]

# I prefer maximin sliced Latin hypersquare designs (Maximin SLHD) and MaxPro SFDs

General recommendations:
Maximin (Sliced) LHD or MaxPro



Latin Hypersquare



SFD – Space-Filling Design; SLHD – Sliced Latin Hypersquare Design

# I prefer maximin sliced Latin hypersquare designs (Maximin SLHD) and MaxPro SFDs
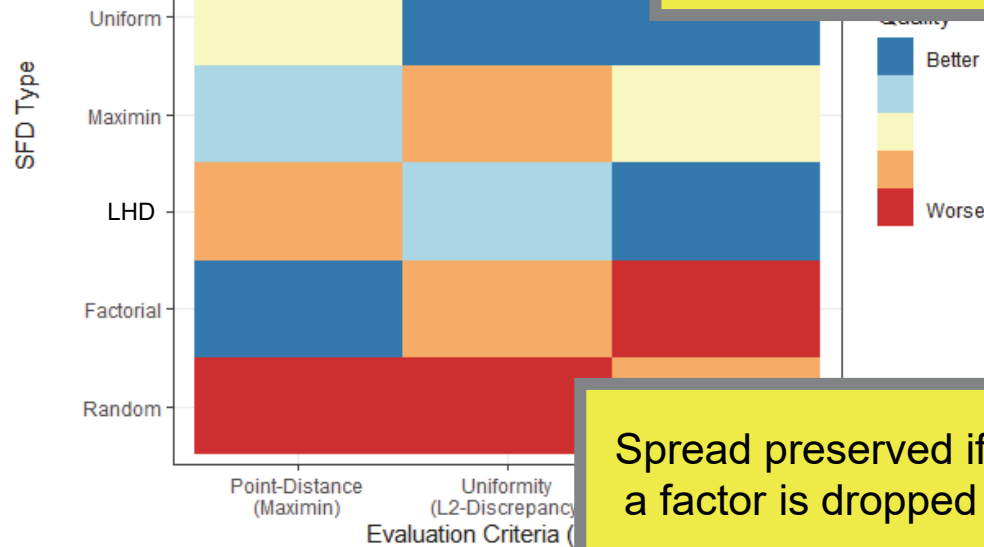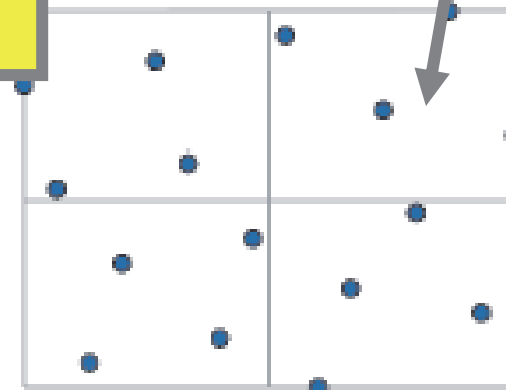
General recommendations:
Maximin (Sliced) LHD or MaxPro

**Handles a small (<5) number of categorical factors**

**Handles more categorical factors better, while also resembling maximin SLHD**

**Good spread in the design space**

Latin Hypersquare

**Quality**
- Better
- Worse

SFD Type (y-axis): Uniform, Maximin, LHD, Factorial, Random

Evaluation Criteria (x-axis): Point-Distance (Maximin), Uniformity (L2-Discrepancy)

**Spread preserved if a factor is dropped**

SFD – Space-Filling Design; SLHD – Sliced Latin Hypersquare Design

IDA

# How should we choose the sample size when generating an SFD?

**SCREE PLOT**



$$n = 10d$$

## Choosing the Sample Size of a Computer Experiment: A Practical Guide
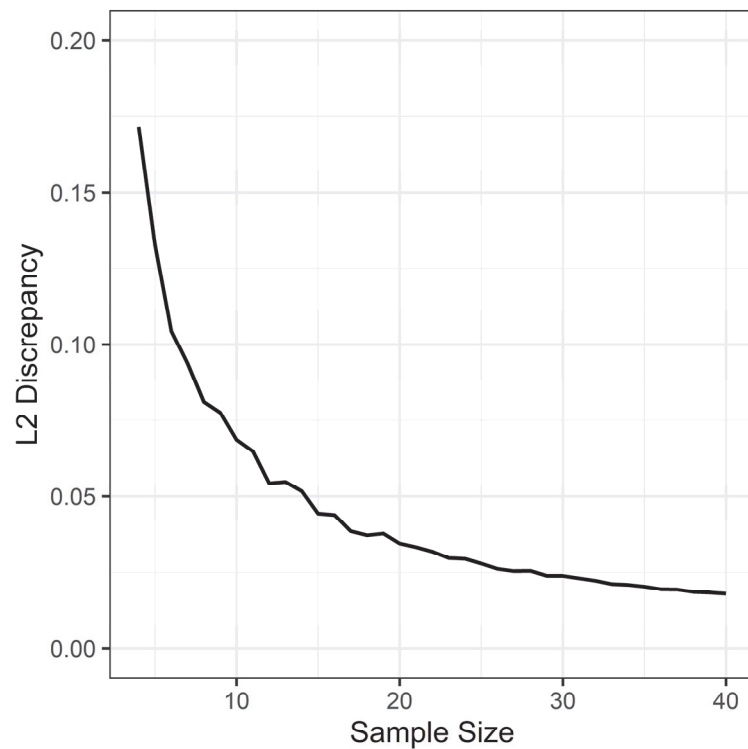
**Jason L. LOEPPKY**

Mathematics, Statistics, and Physics
University of British Columbia, Okanagan
Kelowna, BC V1V 1V7
Canada
(*jason@stat.ubc.ca*)

**Jerome SACKS**

National Institute of Statistical Sciences
Research Triangle Park, NC 27709
(*sacks@niss.org*)

**William J. WELCH**

Department of Statistics
University of British Columbia
Vancouver, BC V6T 1Z2
Canada
(*will@stat.ubc.ca*)

We provide reasons and evidence supporting the informal rule that the number of runs for an effective initial computer experiment should be about 10 times the input dimension. Our arguments quantify two key characteristics of computer codes that affect the sample size required for a desired level of accuracy when approximating the code via a Gaussian process (GP). The first characteristic is the total sensitivity of a code output variable to all input variables; the second corresponds to the way this total sensitivity is distributed across the input variables, specifically the possible presence of a few prominent input factors and many impotent ones (i.e., effect sparsity). Both measures relate directly to the correlation structure in the GP approximation of the code. In this way, the article moves toward a more formal treatment of sample size for a computer experiment. The evidence supporting these arguments stems primarily from a simulation study and via specific codes modeling climate and ligand activation of G-protein.
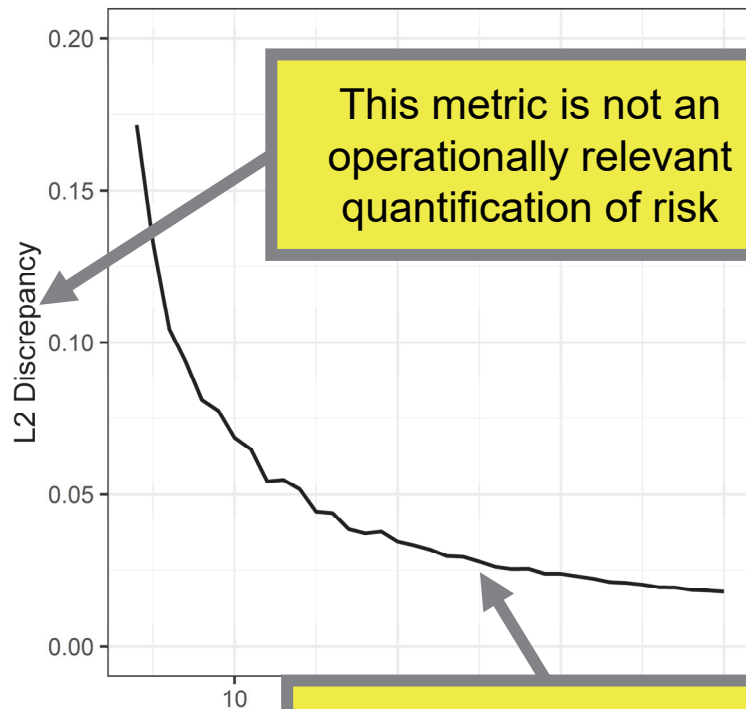
KEY WORDS: Curse of dimensionality; Effect sparsity; Gaussian process; Latin hypercube design; Prediction accuracy; Random function.

# How should we choose the sample size when generating an SFD?

**Seems too small in stochastic cases**

$$n = 10d$$

**SCREE PLOT**



This metric is not an operationally relevant quantification of risk

Where is an objective cut-off?

## Choosing the Sample Size of a Computer Experiment: A Practical Guide

**Jason L. LOEPPKY**

Mathematics, Statistics, and Physics
University of British Columbia, Okanagan
Kelowna, BC V1V 1V7
Canada
(jason@stat.ubc.ca)

**Jerome SACKS**

National Institute of Statistical Sciences
Research Triangle Park, NC 27709
(sacks@niss.org)

**William J. WELCH**

Department of Statistics
University of British Columbia
Vancouver, BC V6T 1Z2
Canada
(will@stat.ubc.ca)

We provide reasons and evidence suppor[...]
initial computer experiment should be ab[...]
key characteristics of computer codes tha[...]
when approximating the code via a Gauss[...]
of a code output variable to all input varia[...]
distributed across the input variables, spe[...]
and many impotent ones (i.e., effect spar[...]
in the GP approximation of the code. In [...]
sample size for a computer experiment. T[...]
a simulation study and via specific codes [...]

KEY WORDS: Curse of dimensionalit[...]
Prediction accuracy; Ra[...]

**ADJACENT QUESTION: WHEN SHOULD REPLICATES BE USED (AND HOW MANY)?**

The minimum may be all you get

SFD – Space-Filling Design

What are the challenges in statistically validating a computer model?

# IDA publications and presentations introduce and recommend surrogate modeling best practices

INSTITUTE FOR DEFENSE ANALYSES

**IDA**

**Metamodeling Techniques for Verification and Validation of Modeling and Simulation Data**

John T. Haman, Project Leader

Curtis G. Miller

September 2022
This publication has not been approved by the sponsor for distribution and release. Reproduction or use of this material is not authorized without prior permission from the responsible IDA Division Director.

IDA Paper P-33230
Log: H 2022-000374

INSTITUTE FOR DEFENSE ANALYSES
730 East Glebe Road
Alexandria, Virginia 22305

**IDA**

**Space Filling Designs and Metamodeling for Understanding Modeling & Simulation Behavior**
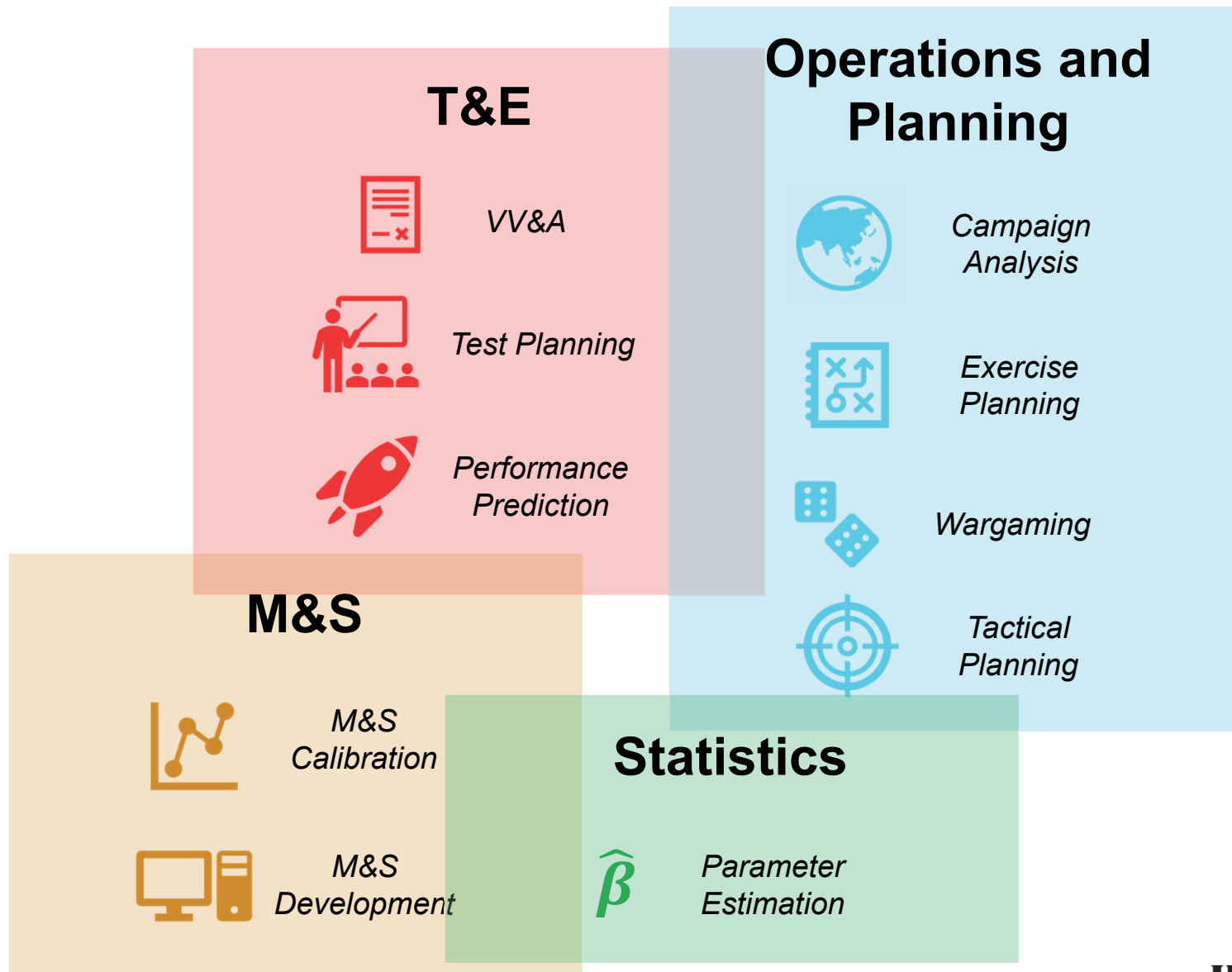
Curtis Miller

April 26, 2022

**Institute for Defense Analyses**
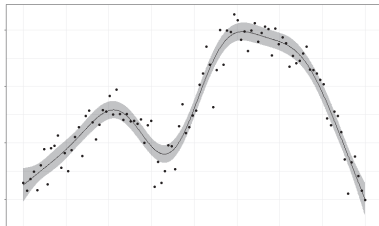730 East Glebe Road ● Alexandria, Virginia 22305

https://dataworks.testscience.org/

https://testscience.org

**IDA** | 25

# A statistical surrogate is a useful product in and of itself



**T&E**
- VV&A
- Test Planning
- Performance Prediction

**Operations and Planning**
- Campaign Analysis
- Exercise Planning
- Wargaming
- Tactical Planning

**M&S**
- M&S Calibration
- M&S Development

**Statistics**
- $\widehat{\beta}$ Parameter Estimation

# Statistical surrogates must allow for relevant assessments of M&S prediction quality
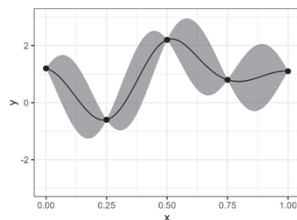
*Statistical surrogates should…*



*… discover unanticipated trends between factors and response variables*
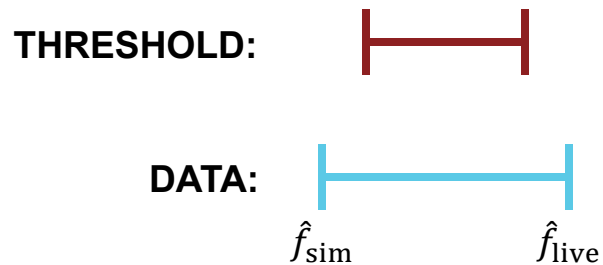


*… predict observed responses well*



*… quantify uncertainty in M&S outputs*

# Statistical surrogates must allow for relevant assessments of M&S prediction quality

*Statistical surrogates should…*

**THRESHOLD:**

**DATA:**

$\hat{f}_{\text{sim}}$         $\hat{f}_{\text{live}}$

*… allow comparing real-world outcomes to M&S predictions*

*… allow expert judgement on prediction quality*

# We recommend different statistical surrogate procedures based on the M&S output

| | CONTINUOUS OUTPUT | DISCRETE OUTPUT |
|---|---|---|
| **DETERMINISTIC** | GAUSSIAN PROCESS (GP) | NEAREST NEIGHBOR (NN) DECISION TREE |
| **STOCHASTIC** | GENERALIZED ADDITIVE MODEL (GAM) | |

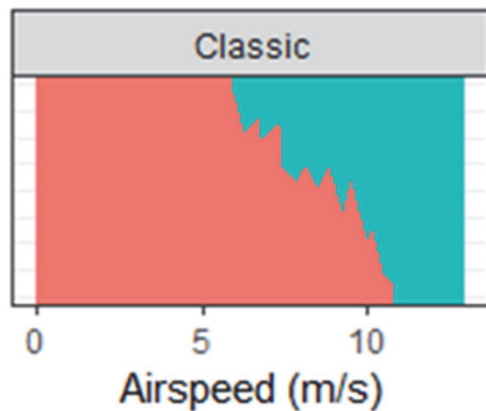# We demonstrate methods using a simple numerical ODE solver

$$\dot{V} = -C_D \left(\rho V^2/2\right) S/m - g\sin(\gamma)$$
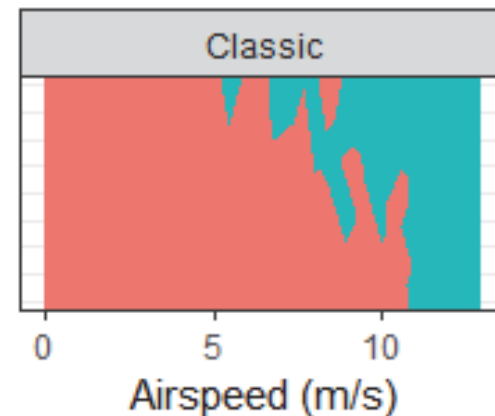$$\dot{\gamma} = \left(C_L \left(\rho V^2/2\right) S/m - g\cos(\gamma)\right)/V$$
$$\dot{h} = V\sin(\gamma)$$
$$\dot{r} = V\cos(\gamma)$$

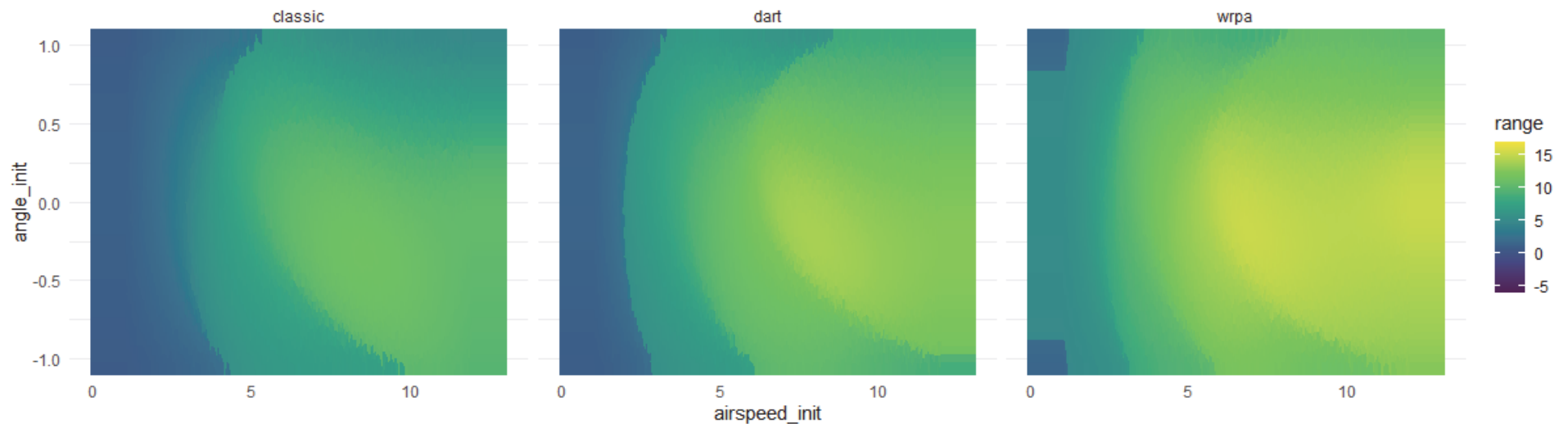Randomize initial conditions for random output (i.e., variance in toss)



**DETERMINISTIC**



**STOCHASTIC**
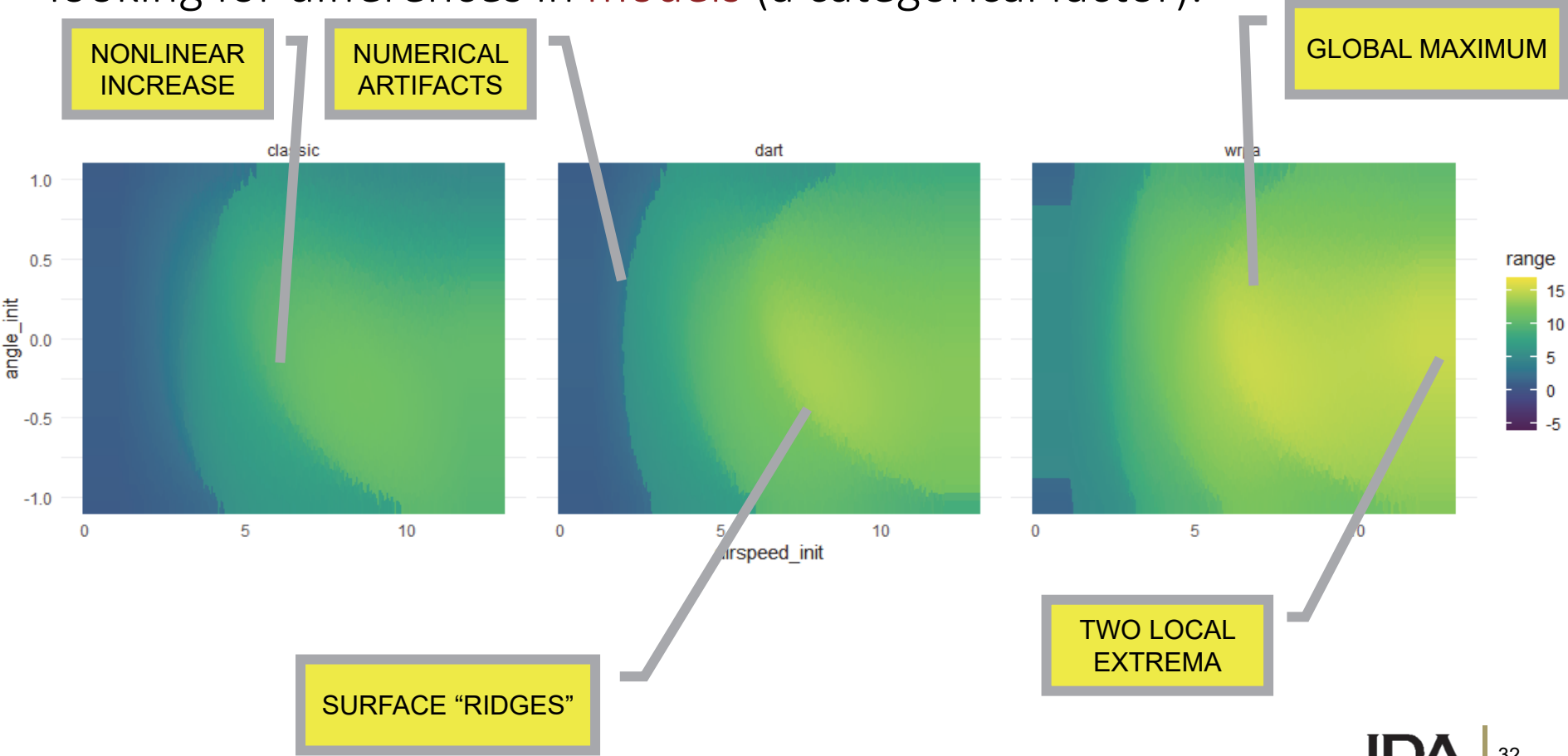
# Model, Airspeed, and Angle Study

Consider a study of the effects of initial airspeed and angle of flight (continuous factors) on flight terminal range (continuous), while also looking for differences in models (a categorical factor).
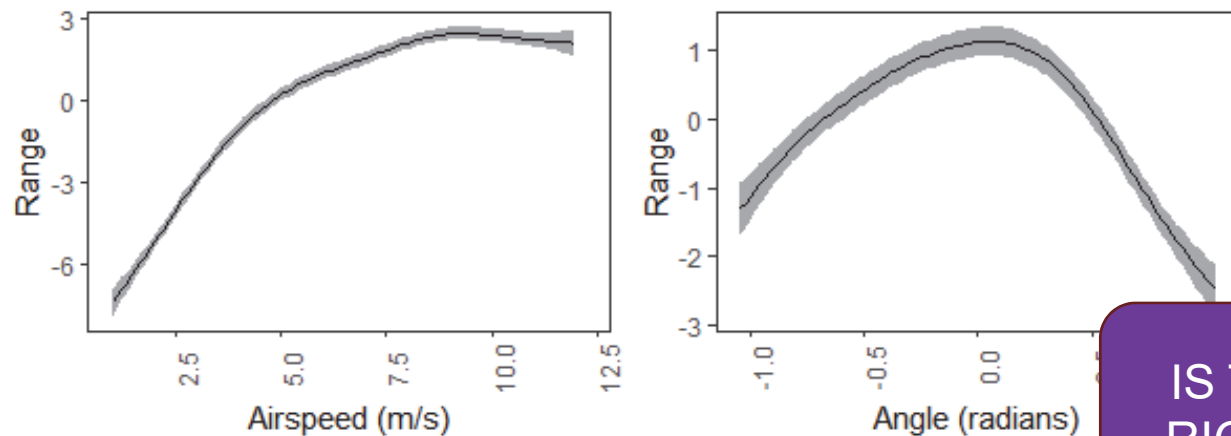
# Model, Airspeed, and Angle Study

Consider a study of the effects of initial airspeed and angle of flight (continuous factors) on flight terminal range (continuous), while also looking for differences in models (a categorical factor).
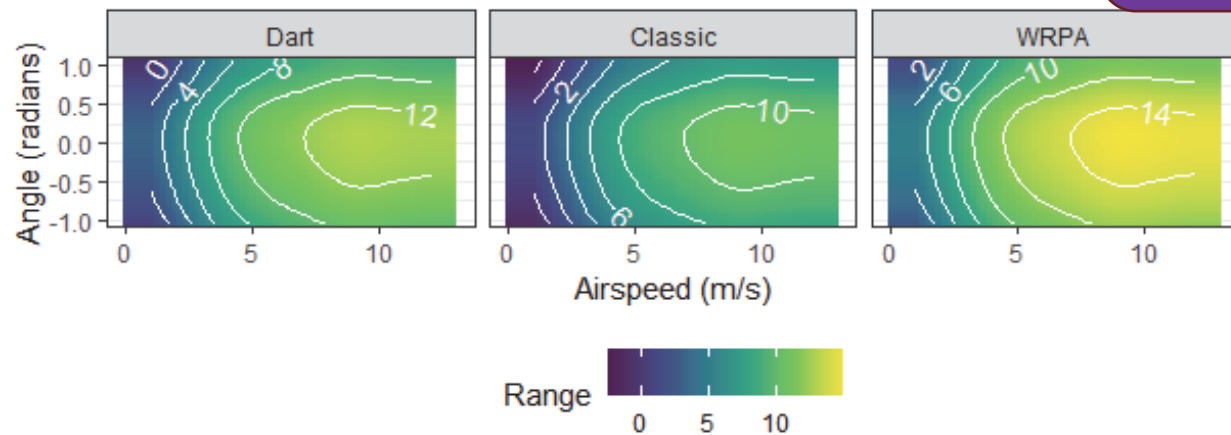
# GAMs allow for qualitative assessment of simulation system performance

$$\text{range}_i = \beta_0 + \beta_{\text{cl}}\text{classic}_i + \beta_{\text{wr}}\text{wrpa}_i + f_{\text{as}}(\text{airspeed}_i) + f_{\text{an}}(\text{angle}_i) + \epsilon_i$$



Average Range:
Dart: 9.3 m
Classic: 7.4 m
WRPA: 11.3 m

IS THIS RIGHT?
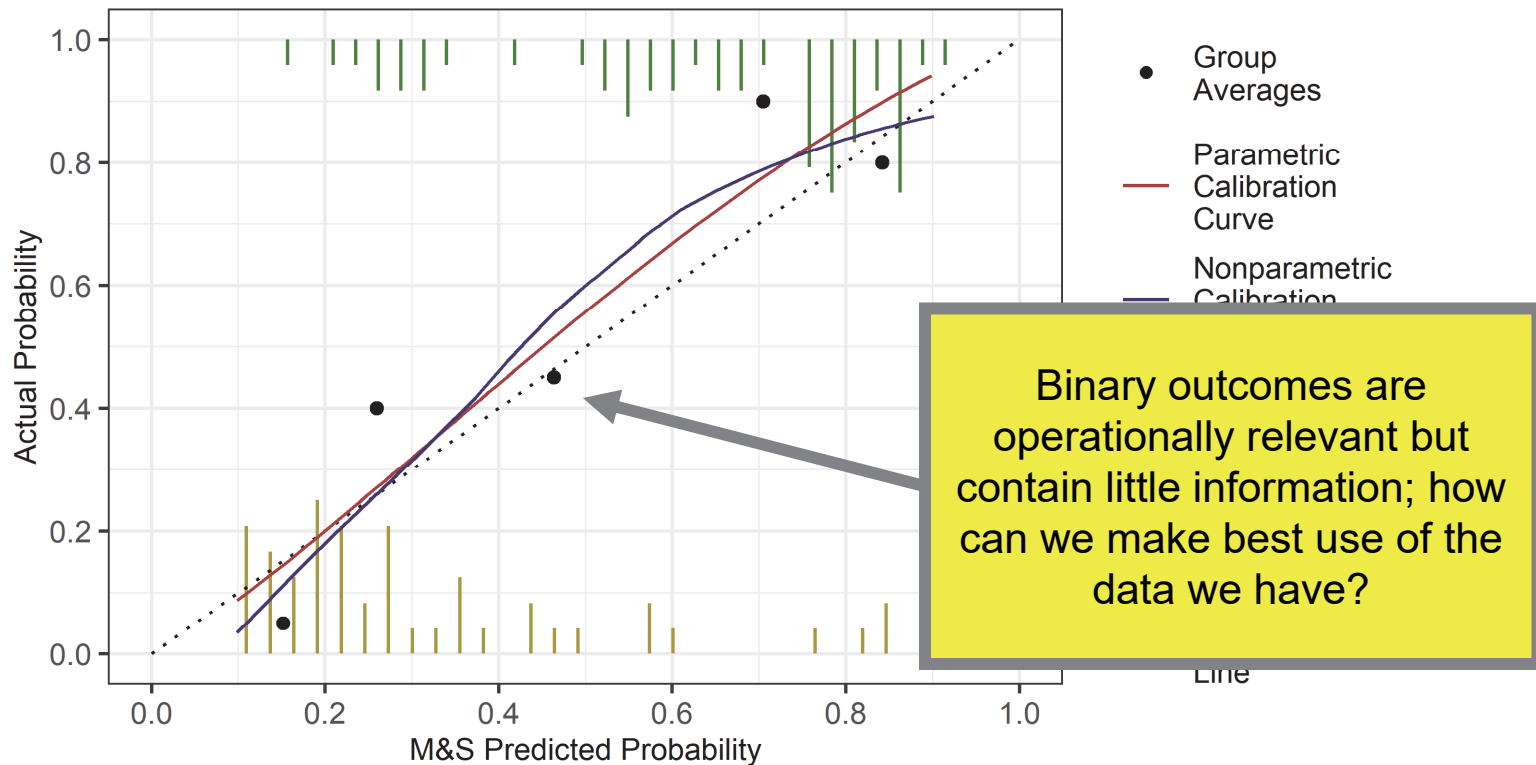
GAM – Generalized Additive Model; WRPA – World Record Paper Airplane

# How can statistical surrogates allow for M&S validation with small sample sizes?

# How can statistical surrogates allow for M&S validation with small sample sizes?



Binary outcomes are operationally relevant but contain little information; how can we make best use of the data we have?

**Legend:**
- Group Averages
- Parametric Calibration Curve
- Nonparametric Calibration
- Line

*Axes: M&S Predicted Probability (x), Actual Probability (y)*

# How should statistical testing incorporate statistical surrogates?

$H_0$: Sim=Live

$H_A$: $H_0$ is false

$$W = \frac{(\bar{Y}_{\text{sim}} - \bar{Y}_{\text{live}})^2}{\text{Var}(\bar{Y}_{\text{sim}} - \bar{Y}_{\text{live}})}$$

**THRESHOLD:**

**DATA:**

$\bar{Y}_{\text{sim}}$  $\bar{Y}_{\text{live}}$

Reject if $W > C_\alpha$.

M&S – Modeling and Simulation

IDA | 36

# How should statistical testing incorporate statistical surrogates?
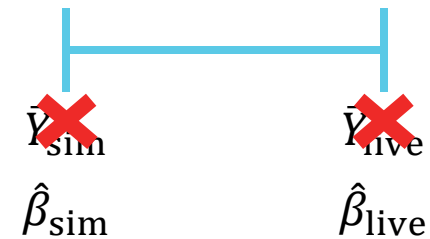
$H_0$: Sim=Live

$H_A$: $H_0$ is false

$$W = \frac{(\bar{Y}_{\text{sim}} - \bar{Y}_{\text{live}})^2}{\text{Var}(\bar{Y}_{\text{sim}} - \bar{Y}_{\text{live}})}$$

$$W = (\hat{\beta}_{\text{sim}} - \hat{\beta}_{\text{live}})^{\text{T}} V^{-1} (\hat{\beta}_{\text{sim}} - \hat{\beta}_{\text{live}})$$

**THRESHOLD:**

**DATA:**

$\bar{Y}_{\text{sim}}$

$\hat{\beta}_{\text{sim}}$

$\bar{Y}_{\text{live}}$

$\hat{\beta}_{\text{live}}$

Reject if $W > C_\alpha$.

M&S – Modeling and Simulation

# How should statistical testing incorporate statistical surrogates?

$H_0$: Sim=Live

$H_A$: $H_0$ is false

$$W = \frac{(\bar{Y}_{\text{sim}} - \bar{Y}_{\text{live}})^2}{\text{Var}(\bar{Y}_{\text{sim}} - \bar{Y}_{\text{live}})}$$

$$W = (\hat{\beta}_{\text{sim}} - \hat{\beta}_{\text{live}})^{\text{T}} V^{-1} (\hat{\beta}_{\text{sim}} - \hat{\beta}_{\text{live}})$$
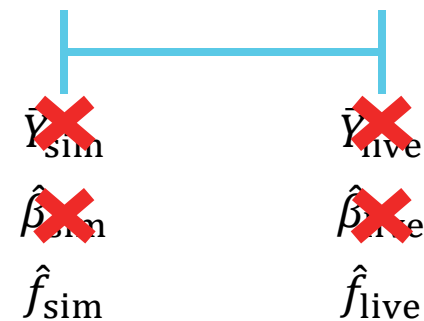
$$W = \left\| \hat{f}_{\text{sim}} - \hat{f}_{\text{live}} \right\|_{V^{-1}}^2$$

Reject if $W > C_\alpha$.

**THRESHOLD:**

**DATA:**

$\bar{Y}_{\text{sim}}$     $\bar{Y}_{\text{live}}$

$\hat{\beta}_{\text{sim}}$     $\hat{\beta}_{\text{live}}$

$\hat{f}_{\text{sim}}$     $\hat{f}_{\text{live}}$

# How should statistical testing incorporate statistical surrogates?

$H_0$: Sim=Live

$H_A$: $H_0$ is false

$$W = \frac{(\bar{Y}_{sim} - \bar{Y}_{live})^2}{\mathrm{Var}(\bar{Y}_{sim} - \bar{Y}_{live})}$$

$$W = (\hat{\beta}_{sim} - \hat{\beta}_{live})^{\mathrm{T}} V^{-1} (\hat{\beta}_{sim} - \cdots)$$

$$W = \left\| \hat{f}_{sim} - \hat{f}_{live} \right\|_{V^{-1}}^2$$

Reject if $W > C_\alpha$.

**THRESHOLD:**

**DATA:**

$\bar{Y}_{sim}$          $\bar{Y}_{live}$

Should these be
- Regression model coefficients?
- Response functions (e.g., GAM smooth)?
- Average response at select factor combinations?

How will we overcome the challenges of modeling and simulation validation?

# M&S validation faces many challenges that must be overcome



T&E workforce statistical literacy
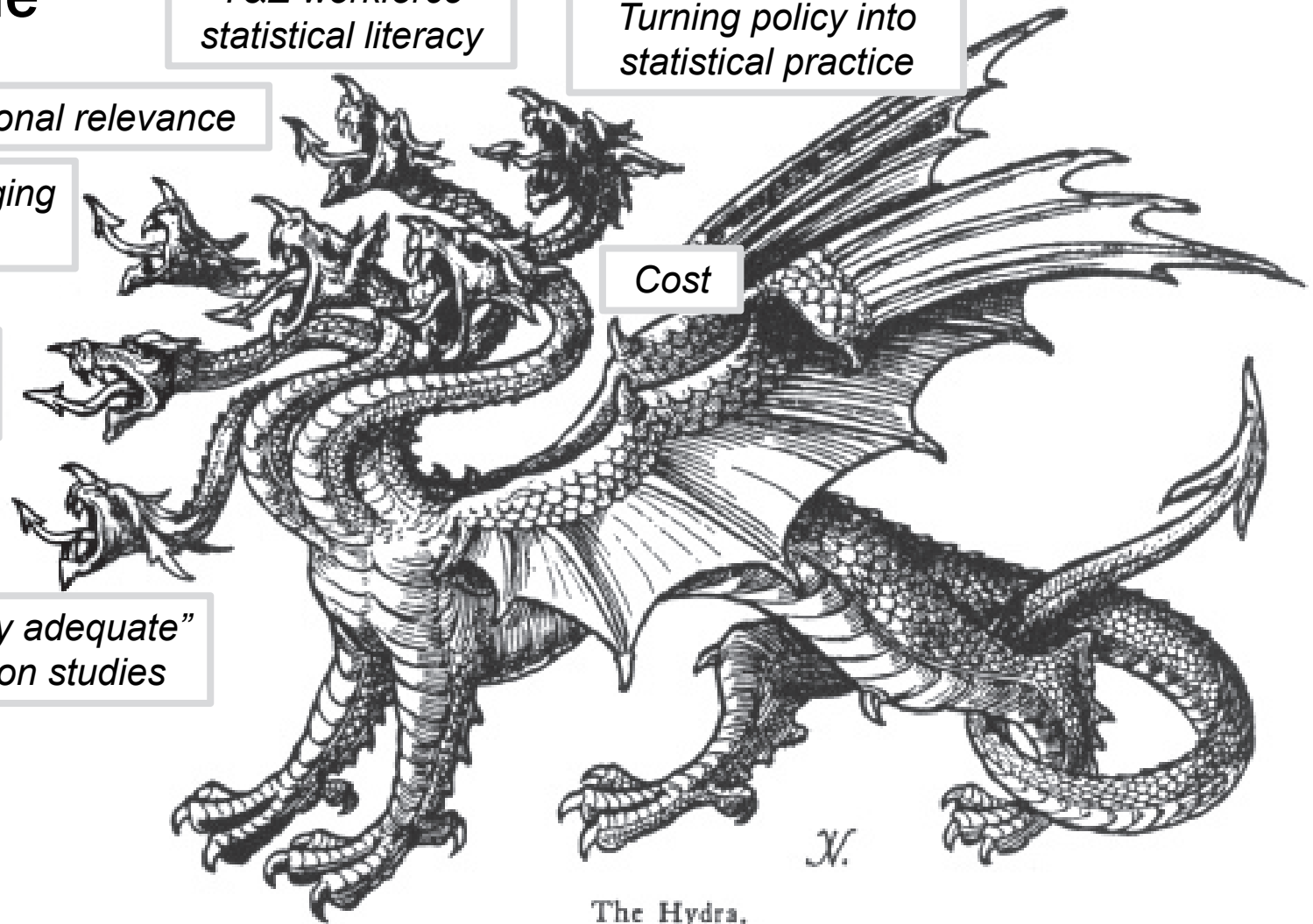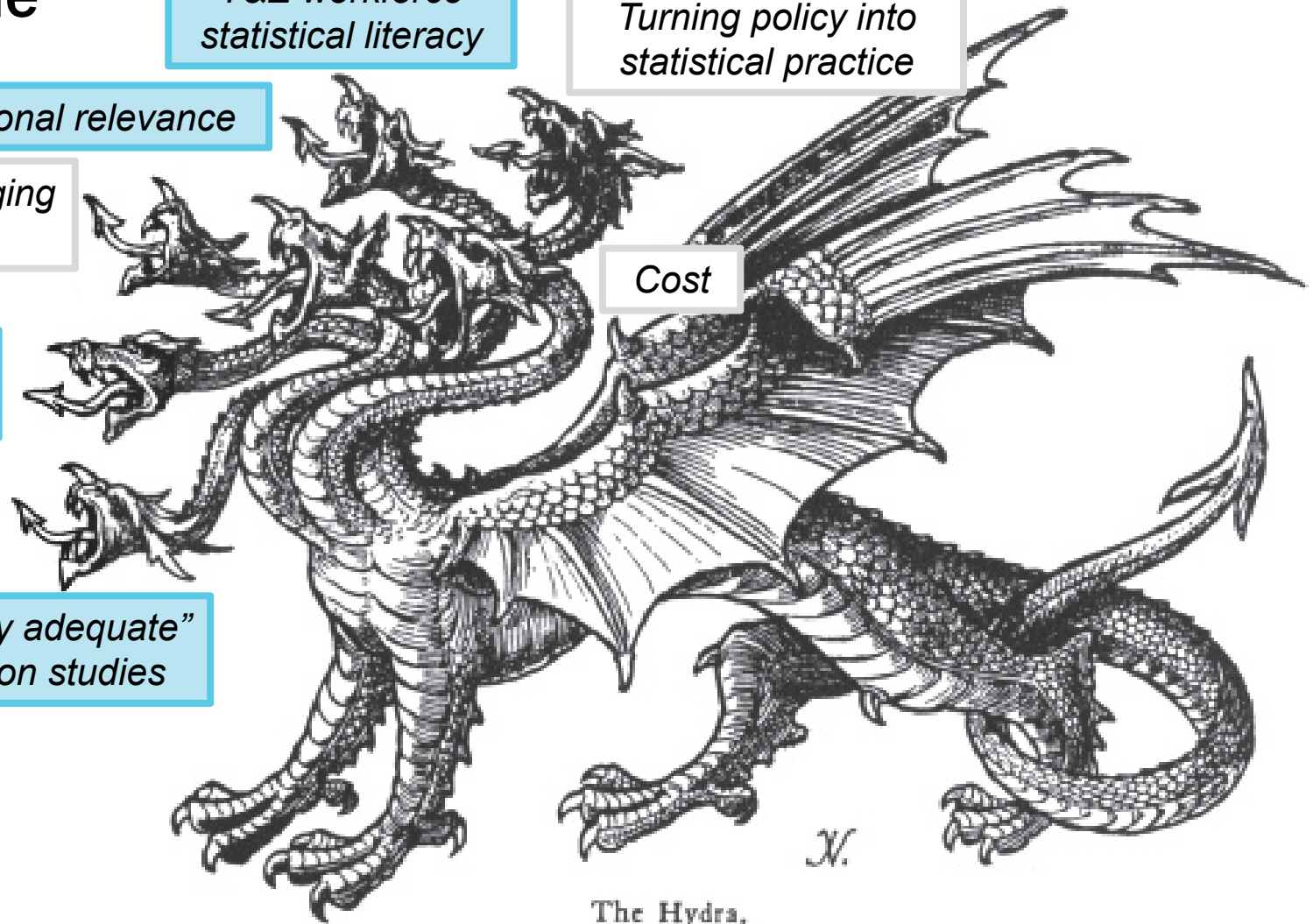
Turning policy into statistical practice

Operational relevance

Continually changing M&S systems

Cost

Small real world samples

"Minimally adequate" simulation studies

The Hydra.

M&S – Modeling and Simulation; T&E – Test and Evaluation

# M&S validation faces many challenges that must be overcome

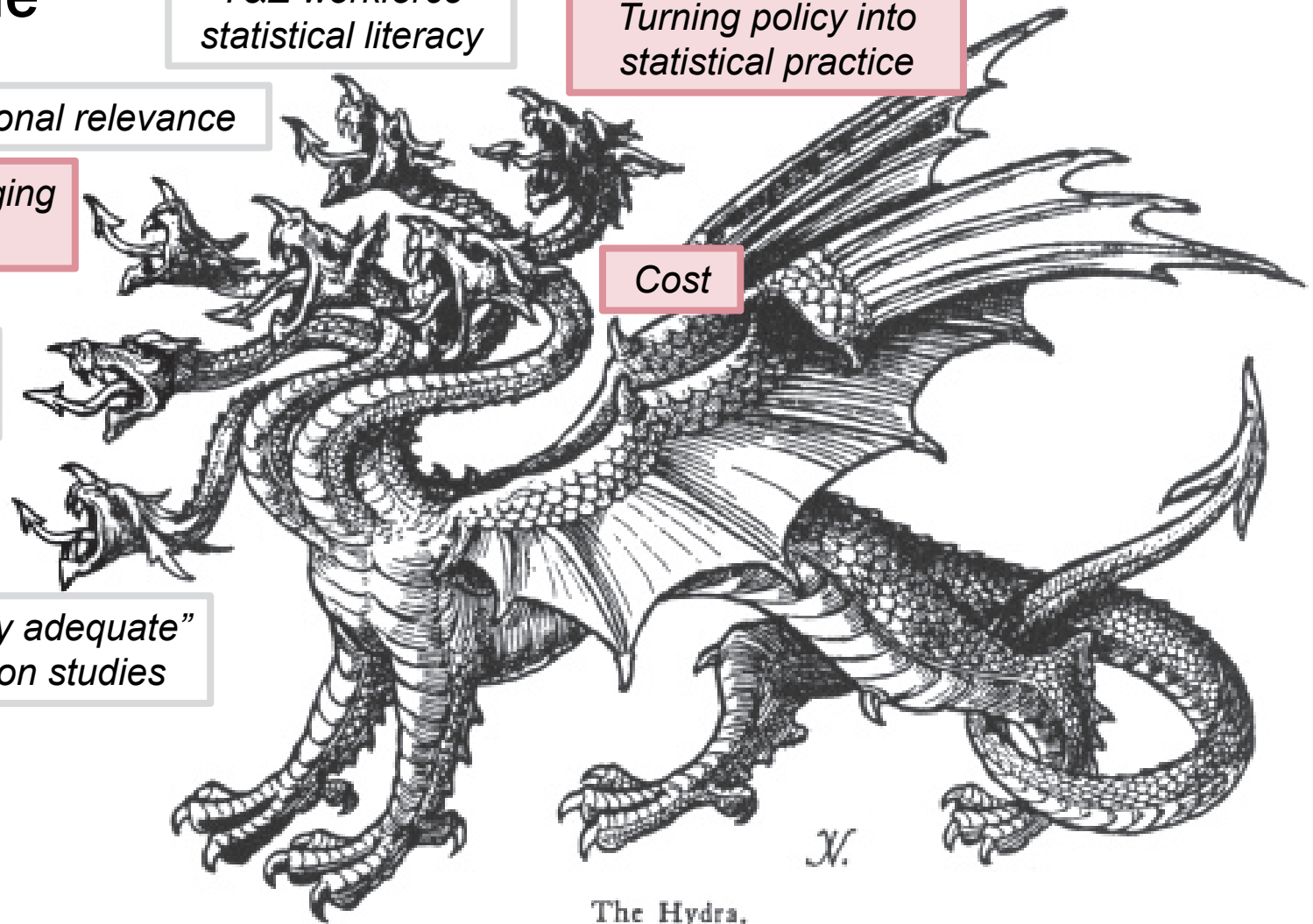T&E workforce statistical literacy

Turning policy into statistical practice

Operational relevance

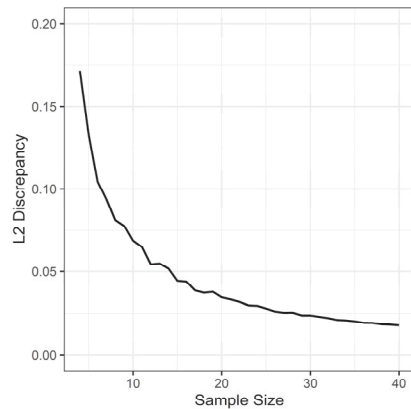Continually changing M&S systems

Cost

Small real world samples

"Minimally adequate" simulation studies

The Hydra.

M&S – Modeling and Simulation; T&E – Test and Evaluation

# M&S validation faces many challenges that must be overcome

T&E workforce statistical literacy

Turning policy into statistical practice

Operational relevance

Continually changing M&S systems

Cost

Small real world samples

"Minimally adequate" simulation studies
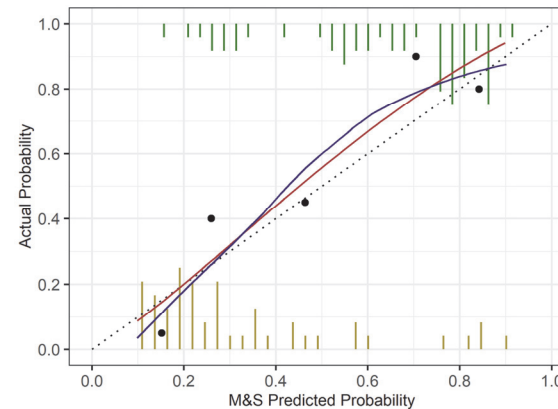
The Hydra.

M&S – Modeling and Simulation; T&E – Test and Evaluation

# We want statisticians to address...



*… how many observations to collect from M&S*



*… how to best handle binary responses*

$$W = \left\| \hat{f}_{\text{sim}} - \hat{f}_{\text{live}} \right\|_{V^{-1}}^{2}$$

*… how to make statistical decisions with surrogate models*

M&S – Modeling and Simulation

# REPORT DOCUMENTATION PAGE

**1. REPORT DATE** *(DD-MM-YYYY)*

**2. REPORT TYPE**

**3. DATES COVERED** *(From - To)*

**4. TITLE AND SUBTITLE**

**5a. CONTRACT NUMBER**

**5b. GRANT NUMBER**

**5c. PROGRAM ELEMENT NUMBER**

**6. AUTHOR(S)**

**5d. PROJECT NUMBER**

**5e. TASK NUMBER**

**5f. WORK UNIT NUMBER**

**7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)**

**8. PERFORMING ORGANIZATION REPORT NUMBER**

**9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)**

**10. SPONSOR/MONITOR'S ACRONYM(S)**

**11. SPONSOR/MONITOR'S REPORT NUMBER(S)**

**12. DISTRIBUTION/AVAILABILITY STATEMENT**

**13. SUPPLEMENTARY NOTES**

**14. ABSTRACT**

**15. SUBJECT TERMS**

**16. SECURITY CLASSIFICATION OF:**

**a. REPORT**

**b. ABSTRACT**

**c. THIS PAGE**

**17. LIMITATION OF ABSTRACT**

**18. NUMBER OF PAGES**

**19a. NAME OF RESPONSIBLE PERSON**

**19b. TELEPHONE NUMBER** *(Include area code)*