# IDA

April 2021 Approved for Public Release. Distribution Unlimited. IDA Document NS D-16393

Log: 2020-000425

INSTITUTE FOR DEFENSE ANALYSES 4850 Mark Center Drive Alexandria, Virginia 22311-1882

#### INSTITUTE FOR DEFENSE ANALYSES

#### Test Science Website Videos - Test Planning and Design

Rebecca Medlin, Project Leader

Kelly Avery Caitlan Fealing Rebecca Medlin Rachel Haga Matthew Avery



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.

#### About This Publication

This work was conducted by the Institute for Defense Analyses (IDA) under contract HQ0034-19-D-0001, Task C9082 "Cross-Divisional Statistics and Data Science Working Group" and Task 229990 "Test Science Applications," for the Office of the Director, Operational Test and Evaluation. The views, opinions, and findings should not be construed as representing the official position of either the Department of Defense or the sponsoring organization.

#### Acknowledgments

The IDA Technical Review Committee was chaired by Mr. Robert R. Soule and consisted of Stephen DeVito, Daniel Hellman, Rachel Haga, Rebecca Medlin, Jason Sheldon, Keyla Pagan-Rivera, and Courtney Au-Yeung from the Operational Evaluation Division.

For more information: Rebecca Medlin, Project Leader rmedlin@ida.org •C9082

Robert R. Soule, Director, Operational Evaluation Division rsoule@ida.org • (703) 845-2482

Copyright Notice © 2021 Institute for Defense Analyses 4850 Mark Center Drive, Alexandria, Virginia 22311-1882 • (703) 845-2000

This material may be reproduced by or for the U.S. Government pursuant to the copyright license under the clause at DFARS 252.227-7013 [Feb. 2014].

Rigorous Analysis | Trusted Expertise | Service to the Nation

#### INSTITUTE FOR DEFENSE ANALYSES

IDA Document NS D-16393

#### Test Science Website Videos - Test Planning and Design

Rebecca Medlin, Project Leader

Kelly Avery Caitlan Fealing Rebecca Medlin Rachel Haga Matthew Avery

#### **Executive Summary**

Understanding the test design process is fundamental to conducting efficient and effective tests of systems. Planning is a critical element of this process. Test planning involves identifying the test objectives, determining the appropriate outcome measure(s), and identifying what factors may influence those outcomes. The first set of website videos introduces the overarching test design process and provides best practices and examples for each of the planning steps.

Test design is the process of actually creating and evaluating a set of experimental test runs. It entails (1) carefully considering the test goals to ensure objectives can be met by the data collected, (2) building the corresponding type of design, and (3) evaluating the design across multiple criteria. The test design should (1) have a large enough sample size to be useful, (2) be executable in the real world, and (3) support the desired analysis. The second set of videos describes and provides recommendations for each element of the test design process, and also introduces considerations for how to manage and store data collected from tests.

#### A. Test Planning

The test design process, which is essentially an example of the scientific method, represents the gold standard for experimentation. It comprises seven steps:

- 1. Identify the test objectives
- 2. Identify the mission-oriented outcome metrics, or response variables
- 3. Identify the factors or variables that may affect those outcomes
- 4. Develop the test design
- 5. Conduct the test and collect the data
- 6. Analyze that data using statistical methods
- 7. Draw conclusions and inform your evaluation

This analytical framework for collecting and analyzing data is applicable to nearly every domain and industry. Successfully performing this process requires iteration and collaboration among subject matter experts. Clearly defined objectives are essential when designing any operational or live fire test, because they will drive the associated test design and analysis. The two most common test objectives used in test and evaluation are (1) screening for influential factors and (2) characterizing a system's performance across a set of conditions. These objectives are best applied sequentially: first a screening experiment narrows down the set of factors that are likely to matter, and then a characterization design details the performance of the system under each of the relevant conditions.

Outcome metrics, also frequently called response variables, are what testers use to assess their test objectives. A good response variable should (1) provide a meaningful measure of performance, (2) support evaluation of requirements, and (3) lend itself well to experimental design, meaning it is measurable, valid, and informative. Oftentimes more than one response variable is necessary to capture the full picture of system performance.

Because defense systems rarely have uniform performance in every situation, it is important for testers to identify factors that may drive changes in performance. These typically relate to the environmental conditions, the threat level, the mission, the user type, or the system configuration. Good factors are significant, controllable, and informative. After brainstorming all potential factors, testers can use a factor management scheme to prioritize which factors should be controlled and systematically varied in a test design, and which will simply be logged, or fixed at a specific level. Modern Design of Experiments (DOE) techniques can investigate a large number of factors very efficiently.

#### **B.** Test Design

The best approach for collecting data from several factors is to conduct a statistically designed experiment. DOE is an experimentation strategy in which testers vary factors in a coordinated, strategic way, instead varying them one at a time, or arbitrarily choosing test cases. A test strategy that employs DOE will always provide the most powerful allocation of test resources for a given number of tests.

The specific choice of design should be tied to the test objectives and the anticipated analysis. For example, designs meant for factor screening have fundamentally different properties than designs meant to characterize or optimize. Sequential test design strategies, where the results of early tests inform the construction of future tests, are highly recommended and can provide additional resource savings.

Common test designs include factorial designs, response surface methodology, and optimal designs. The decision on which test design method is most appropriate depends on the question, metrics, types of factors (numeric or categorical), and available test resources. Many software packages are available to build and evaluate each of these designs. The question of "How much testing is enough?" is persistent across Department of Defense test and evaluation endeavors. DOE can help decision makers "right-size" a test, providing a framework for building and evaluating test plans and assessing risk. Leveraging statistical techniques in both test design and test analysis helps decision makers reach correct conclusions about system performance.

Test designs should be realistic and executable. This might mean that factors cannot be perfectly randomized or that certain factor combinations will be disallowed. To avoid reporting artificially inflated metrics, design evaluations should always account for test structure and randomization.

Having a plan for collecting and managing test data is critical to a successful test program. A successful plan will include guidance on (1) how to collect data in the appropriate format, (2) where and how to store acquired data, (3) how to document, organize, and potentially automate data cleaning, reduction, and analysis, and (4) how to handle data access and distribution.





# **Test Planning**

#### **Lesson 1: Introduction to Test Planning**

#### **Institute for Defense Analyses**

4850 Mark Center Drive • Alexandria, Virginia 22311-1882

To provide **warfighters** and **decision makers** with information about military systems that is

RELEVANT

## CREDIBLE

OBJECTIVE



To provide **warfighters** and **decision makers** with information about military systems that is

RELEVANT

## CREDIBLE

OBJECTIVE





To provide **warfighters** and **decision makers** with information about military systems that is

RELEVANT

## CREDIBLE

OBJECTIVE





TestScience Data . Driven . Defense

To provide **warfighters** and **decision makers** with information about military systems that is

RELEVANT

## CREDIBLE

OBJECTIVE







TestScience Data . Driven . Defense

The Design of Experiments process provides a framework and an analytical basis to help us do T&E systematically and efficiently.



- The DOE process, at its core, is the scientific method.
- It is the gold standard for performing experiments and determining cause and effect.
- What we mean by "experiment" goes well beyond the 'chemistry lab' scenario it often evokes.



- An experiment includes any test where variables are purposefully manipulated or inputs are purposefully set.
- So, the types of operational and live fire tests we conduct would certainly count as experiments.

TestScience Data . Driven . Defense





1. Identify capabilities to be tested



2. Identify the mission-oriented outcome metrics



3. Identify factors that may affect outcome



4. Develop the test design



5. Conduct the test



6. Analyze the data



7. Draw conclusions





TestScience Data . Driven . Defense

## **Applying the Testing Framework**

The test design process is applicable to any type of test:

- System Performance
- Human Factors
- Reliability
- Cybersecurity
- Live Fire
- Etc.



## **Applying the Testing Framework**

Keys to success:

- Frequent communication
- Collaboration among technical, operational, and statistical experts
- Openness to iteration



# The Three 'Test Planning' Steps are Crucial for Successful T&E

#### 1. Identify capabilities to be tested

- Objectives of the test
- Questions you can ask about the system

#### 2. Identify mission-oriented outcome metrics

- How you should measure system performance
- Response variables

#### 3. Identify factors that may affect outcome

- Conditions or configurations under which the system may operate
- Operational envelope



# The Three 'Test Planning' Steps are Crucial for Successful T&E

#### 1. Identify capabilities to be tested

- Objectives of the test
- Questions you can ask about the system

#### 2. Identify mission-oriented outcome metrics

- How you should measure system performance
- Response variables

#### 3. Identify factors that may affect outcome

- Conditions or configurations under which the system may operate
- Operational envelope

The next three lessons cover each of these planning steps in more detail...







# **Test Planning**

#### **Lesson 2: Identify Test Objectives**

#### **Institute for Defense Analyses**

4850 Mark Center Drive • Alexandria, Virginia 22311-1882

### **The Test Design Process**



TestScience Data . Driven . Defense

























# It is essential to have a clear goal when designing an experiment

#### 1. Identify capabilities to be tested

- Objectives of the test
- Questions you can ask about the system


# It is essential to have a clear goal when designing an experiment

1. Identify capabilities to be tested

- Objectives of the test
- Questions you can ask about the system

Test objectives will drive all future steps, including the test design and statistical analysis.



### **Common Test Objectives**

- Screen for influential factors driving performance.
- Characterize performance across an operational envelope.
- **Compare** two systems (or more) across a variety of operating conditions.
- Identify problems that degrade system performance.
- **Optimize** system performance with respect to a set of conditions.



# **Common Test Objectives**

- Screen for influential factors driving performance. Characterize performance across an operational envelope.
- Compare two systems (or more) across a variety of operating conditions.
- **Identify problems** that degrade system performance.
- **Optimize** system performance with respect to a set of conditions.



# **Common Test Objectives**

- Screen for influential factors driving performance. Characterize performance across an operational envelope.
- Compare two systems (or more) across a variety of operating conditions.
- **Identify problems** that degrade system performance.
- **Optimize** system performance with respect to a set of conditions.

DoDI 5000.02: The test program should produce "data to characterize combat mission capability across an appropriately selected set of factors and conditions."



### **Screen for Influential Factors**

- Screening is useful early on in testing when we may know little about how the system responds to various factors.
- General approach:
  - Identify all potential factors that are thought to affect the response variable.
  - Choose an initial experimental design that uses minimal test resources.
  - Through analysis, identify the factors that have the largest impact on the response.
  - Update next test design to characterize the response (performance) as a function of only the important factors.



## **Screen for Influential Factors**

- Screening is useful early on in testing when we may know little about how the system responds to various factors.
- General approach:
  - Identify all potential factors that are thought to affect the response variable.
  - Choose an initial experimental design that uses minimal test resources.
  - Through analysis, identify the factors that have the largest impact on the response.
  - Update next test design to characterize the response (performance) as a function of only the important factors.

Screening is essential to integrated and sequential testing.



#### **Characterize performance across conditions**

- Characterization tests describe performance at all relevant conditions and configurations in which the system may operate (i.e., across the operational envelope).
- They are also useful for determining if a system meets requirements across a variety of operational conditions.
- Characterization designs facilitate building precise mathematical models of performance based on test data.

#### **Characterize performance across conditions**

- Characterization tests describe performance at all relevant conditions and configurations in which the system may operate (i.e., across the operational envelope).
- They are also useful for determining if a system meets requirements across a variety of operational conditions.
- Characterization designs facilitate building precise mathematical models of performance based on test data.

Characterization is the most important and common goal for operational testing.





# **Test Planning**

#### **Lesson 3: Select Outcome Measures**

#### **Institute for Defense Analyses**

4850 Mark Center Drive • Alexandria, Virginia 22311-1882

#### **The Test Design Process**



TestScience Data.Driven.Defense

#### **Response variables measure the outcome of your test**

2. Identify mission-oriented outcome metrics

- How you should measure system performance
- Response variables



#### **Response variables measure the outcome of your test**

2. Identify mission-oriented outcome metrics

- How you should measure system performance
- Response variables



TestScience Data . Driven . Defense

- Provide a meaningful measure of system performance or mission capability
  - These measures should encapsulate reasons for procuring the system



- Provide a meaningful measure of system performance or mission capability
  - These measures should encapsulate reasons for procuring the system





- Provide a meaningful measure of system performance or mission capability
- Provide adequate data to evaluate requirements
  - Selection of response variables is influenced by, but not limited to, the requirements.



- Provide a meaningful measure of system performance or mission capability
- Provide adequate data to evaluate requirements
  - Selection of response variables is influenced by, but not limited to, the requirements.





- Provide a meaningful measure of system performance or mission capability
- Provide adequate data to evaluate requirements
  - Selection of response variables is influenced by, but not limited to, the requirements.





- Provide a meaningful measure of system performance or mission capability
- Provide adequate data to evaluate requirements
- Lend themselves well to defensible **experimental design**, which means they are:
  - Measurable. They can be measured at a reasonable cost and without affecting the test outcome.
  - Valid. They directly address the test objective.
  - Informative. Continuous responses always provide more information per test point than pass/fail metrics.

- Provide a meaningful measure of system performance or mission capability
- Provide adequate data to evaluate requirements
- Lend themselves well to defensible experimental design, which means they are:
  - Measurable. They can be measured at a reasonable cost and without affecting the test outcome.
  - Valid. They directly address the test objective.
  - Informative. Continuous responses always provide more information per test point than pass/fail metrics.

Multiple responses are common and almost always necessary!

Data . Driven . Defense

#### Who is the better dart player?

Player 1







#### Who is the better dart player?

 Player 1
 Player 2

We have a much better picture of the players' abilities by looking at the miss distances from bullseye compared to just a bullseye hit/miss!

TestScience Data . Driven . Defense

# We can often convert probability requirements into continuous measures

	Chemical Agent Detector	Submarine Mine Detection	Missile/Bomb	Radar Detection
Requirement	Probability of detection greater than 85% after one minute of exposure	Probability of detection greater than 80% outside 200 meters	Probability of hit at least 90%	Probability of detection greater than 90% at 300 kilometers.
Original Test Measurement	Detect prior to 1 minute? (Yes/No)	Detect/Non-Detect	Hit/Miss	Detect/Non-Detect
Modified Test Measurement	Detection Time	Detection Range	Miss Distance	Detection Range





# **Test Planning**

#### **Lesson 4: Select Factors Affecting Outcomes**

#### **Institute for Defense Analyses**

4850 Mark Center Drive • Alexandria, Virginia 22311-1882

#### **The Test Design Process**



TestScience Data.Driven.Defense

# What might influence the performance of the process or system under test?

#### 3. Identify factors that may affect outcome

- Conditions or configurations under which the system may operate
- Operational envelope



# What might influence the performance of the process or system under test?

#### 3. Identify factors that may affect outcome

- Conditions or configurations under which the system may operate
- Operational envelope



TestScience Data.Driven.Defense

# Terminology

**Factors** are independent variables that are expected to affect the outcome of a test.

**Levels** are the specific values that the factors assume. Factor levels are often referred to as the test conditions.





**Important:** Factors are expected to have an effect on the test outcome.



Important: Factors are expected to have an effect on the test outcome.

**Controllable:** Factors can be controlled (i.e., set to a specific level) at a reasonable cost.



Important: Factors are expected to have an effect on the test outcome.

Controllable: Factors can be controlled (i.e., set to a specific level) at a reasonable cost.

**Informative:** Continuous factors are preferred to categorical factors.



Important: Factors are expected to have an effect on the test outcome.

Controllable: Factors can be controlled (i.e., set to a specific level) at a reasonable cost.

**Informative:** Continuous factors are preferred to categorical factors.

- -Continuous factors allow for:
  - oInterpolation
  - o Higher power to detect significance of a factor
  - oStrategic point placement



# Testers should brainstorm ALL potential factors that could affect test outcomes



# Testers should brainstorm ALL potential factors that could affect test outcomes







# Testers should brainstorm ALL potential factors that could affect test outcomes

- Factors of particular interest to the experimenter Factors that are measurable, controllable, and thought to be (very) influential





# Testers should brainstorm ALL potential factors that could affect test outcomes



# Modern experimental designs can investigate a large number of factors efficiently

Testers should err on the side of strategically varying factors

Fixing factors limits conclusions about the system to the one condition tested

For recorded factors there is no guarantee that all levels of interest will be observed

- Potential tradeoff with operational realism

Adding factors does not cause the test size to grow exponentially!

Data . Driven . Defense




# **Test Design**

#### Lesson 1: Introduction to Test Design

#### **Institute for Defense Analyses**

4850 Mark Center Drive • Alexandria, Virginia 22311-1882

#### **Testing Framework: The Gold Standard of Experimentation**



### **Common Challenges to Test**



#### **Test Goals Influence Test Designs**

Screen Characterize Optimize Compare **Identify Problems** Predict Improve Demonstrate









A **one-factor-at-a-time** test strategy is inefficient and fails to detect interactions





A **one-factor-at-a-time** test strategy is inefficient and fails to detect interactions

Changing the variables together prevents the ability to detect cause end effect.



A test strategy of "We'll just do what we did last time" or "We'll test the special cases" might miss important performance shortfalls.



Special/Critical Cases

#### DOE changes "I think" to "I know"

A test strategy that employs **DOE** will provide the most powerful allocation of test resources for a given number of tests.





# **DOE Provides a Structured Approach to Picking Test Points**



#### **DOE Provides a Structured Approach to Picking Test Points**



Data . Driven . Defense

### "How much testing is enough?"

How many test points are enough to get it *right*?

- 3? That's how much money/time we have.
- 8-10? That's what we did last time.
- 30? From my high school stats class, something good happens at 30!

DOE methods provide tools to assess the risk of drawing incorrect conclusions



Executing an operationally realistic test means your factors might not be **randomized** or that certain factor combinations will be **disallowed**.



Executing an operationally realistic test means your factors might not be **randomized** or that certain factor combinations will be **disallowed**.



Factor	Level
Illumination	Light, Dark
Range to Target	Short, Medium, Long
Fuze Type	А, В
Angle of Fire	Low, High



Executing an operationally realistic test means your factors might not be **randomized** or that certain factor combinations will be **disallowed**.



Factor	Level
Illumination	Light, Dark
Range to Target	Short, Medium, Long
Fuze Type	А, В
Angle of Fire	Low, High

Executing an operationally realistic test means your factors might not be **randomized** or that certain factor combinations will be **disallowed**.

### **Response Variable**

Miss Distance



CONOPS: Fire  $\rightarrow$  Move  $\rightarrow$  Fire

		Difficulty of
Factor	Level	Randomization
Illumination	Light, Dark	Hard
Range to Target	Short, Medium, Long	Hard
Fuze Type	А, В	Easy
Angle of Fire	Low, High	Easy



Data . Driven . Defense

#### A Good Test Design Supports Credible and Clear Analyses

A successful analysis should follow the structure of the test design.

Building a statistical model using the factors specified in our test design allows us to quantify differences across relevant conditions and draw conclusions with confidence.

All experiments are designed with an analysis methodology in mind: to reap the benefits, we need to follow through to the analysis!







# **Test Design**

#### **Lesson 2: Determine Design Goals**

#### **Institute for Defense Analyses**

4850 Mark Center Drive • Alexandria, Virginia 22311-1882

Are the system, mission, and operating environment clearly defined? Are the goals of the test clearly defined? What are the question(s) to be answered in testing? Does the overall test strategy support characterization of combat mission capability?



Data . Driven . Defense

Are the system, mission, and operating environment clearly defined?

Are the goals of the test clearly defined?



 $\mathbf{?}$ 

?

What are the question(s) to be answered in testing?



Does the overall test strategy support characterization of combat mission capability?

Are the system, mission, and operating environment clearly defined?

Are the goals of the test clearly defined?



 $\mathbf{\Omega}$ 

**?** 

What are the question(s) to be answered in testing?



Does the overall test strategy support characterization of combat mission capability?

# What are some common design goals?



# What are some common design goals?



What are reasons to screen for factors?



What are reasons to screen for factors?



Screening is useful early in testing when we know little about how the system responds to various factors.



What are reasons to screen for factors?

Screening is useful early in testing when we know little about how the system responds to various factors.



Screening is essential to integrated & sequential testing to reduce the factor space prior to operational testing.

**Consider a novel hypersonic missile** 



How do you pick a screening design?



Many factors

Typically 2 levels per factor



How do you pick a screening design?





How do you pick a screening design?





What are reasons to characterize performance?



What are reasons to characterize performance?



Characterization tests describe system performance across the operational envelope



#### What are reasons to characterize performance?

Characterization tests describe system performance across the operational envelope



Characterization tests determine if a system meets requirements across operational conditions


## What are reasons to characterize performance?

Characterization tests describe system performance across the operational envelope



Characterization tests determine if a system meets requirements across operational conditions

**DoDI 5000.02**: The test program should produce "data to characterize combat mission capability across an appropriately selected set of factors and conditions."



How will our factors impact its performance?





How will our factors impact its performance?

We identified factors relevant to performance through screening





How will our factors impact its performance?

We identified factors relevant to performance through screening We need to test across a variety of operational conditions ...





How will our factors impact its performance?

We identified factors relevant to performance through screening We need to test across a variety of operational conditions ... ... to adequately understand performance using minimal resources

How do you pick a characterizing design?





How do you pick a characterizing design?





How do you pick a characterizing design?



What are reasons to optimize performance?



Optimization test designs seek the combination of controllable factors that optimizes test outcomes



What are reasons to optimize performance?

Optimization test designs seek the combination of controllable factors that optimizes test outcomes



Optimization tests are useful in design, manufacturing, and the development of tactics, techniques, & procedures





through screening

profile settings ...

distance

How do you pick an optimization design?



Few factors; many levels

**Precision needed** 



How do you pick an optimization design?





How do you pick an optimization design?



An effective test depends upon good, clearly articulated design goals



We have learned about 3 common design goals and how to approach them:

Screening

Characterizing

Optimizing



An effective test depends upon good, clearly articulated design goals



 We have learned about 3 common design goals and how to approach them:

 Screening
 Characterizing
 Optimizing

We have learning that screening is a great tool in sequential and integrated testing to better constrain the design space for characterization and optimization tests.



An effective test depends upon good, clearly articulated design goals



We have learned about 3 common design goals and how to approach them:

 Screening
 Characterizing
 Optimizing



We have learning that screening is a great tool in sequential and integrated testing to better constrain the design space for characterization and optimization tests.



In the next lesson, *Creating the Design*, we will learn more about common experimental designs.







## **Test Design**

#### Lesson 3: Create the Test Design

### **Institute for Defense Analyses**

4850 Mark Center Drive • Alexandria, Virginia 22311-1882

**Design of experiments (DOE)** 

## Which Points?





Which Points? Testing a New Artillery Cannon

### **Response Variable**

#### Accuracy of Fires - Miss Distance





Factor	Level
Time of Day	Day, Night
Range to Target	Short, Medium, Long
Projectile Type	A, B, C
Angle of Fire	Low, High

## Factorial Designs Examine All Possible Combinations of Each Factor Level



### Full Factorial Designs:

- Support the test goals of characterize or compare.
- Examine every possible combination of each level.
- Allow for the estimation of all main effects and all possible interactions.
- Are highly efficient and informative, though potentially prohibitively costly.



# **Factorial Designs** Examine All Possible Combinations of Each Factor Level



Tests all possible combinations. The first fire mission is conducted during the *Day* at *Short* Range for *Low* Angle using *Projectile A* 

### Full Factorial Designs:

- Support the test goals of characterize or compare.
- Examine every possible combination of each level.
- Allow for the estimation of all main effects and all possible interactions.
- Are highly efficient and informative, though potentially prohibitively costly.



## Typically, Screening Or Characterization Experiments Involve a Fractional Factorial Design



#### Fractional Factorial Designs:

- Support the test goals of screen, characterize, or compare.
- Require a subset of the runs required for a full factorial.
- Achieve a large reduction in test points by trading off the ability to estimate high-order interaction effects.



## Typically, Screening Or Characterization Experiments Involve a Fractional Factorial Design



Places test points at the right conditions to support a main effects model

## Fractional Factorial Designs:

- Support the test goals of screen, characterize, or compare.
- Require a subset of the runs required for a full factorial.
- Achieve a large reduction in test points by trading off the ability to estimate high-order interaction effects.



**Response Surface Designs** spread test points to collect data throughout the experimental region to support a detailed model of the response



## **Optimal Designs** Are Most Useful When The Number Of Test Points Is Constrained To Preclude A Factorial Design



Places test points at the right conditions to estimate specified model terms

#### **Optimal Designs**

- Support the test goals of characterize, optimize, predict.
- Useful when the number of test points is constrained to preclude a factorial design.
- Requires researcher-specified model and a fixed sample size (a subset of the runs required for a full factorial).

## **Plan for a Sequential Experimental Campaign**



## **Rarely is the Test Design Process Simple**

Multidimensional Evaluation

Many Categorical Factors with More than Two Levels

Restrictions on Randomization

**Disallowed Factor Combinations** 

Limited Resources for Replication



## Our test of the new cannon includes factors that are Difficult to Completely Randomize and have Disallowed Combinations

			$\sim$
Factor	Level	Difficult to Randomize	Disallowed Combinations
Time of Day	Day, Night	Hard	Night and Short Range
Range to Target	Short, Medium, Long	Hard	
Projectile Type	A, B, C	Easy	A and Long Range
Angle of Fire	Low, High	Easy	

	Time of Day	Range	Angle of Fire	Projectile Typ	e				
Γ			Low	В	Snlit-nlot designs				
	Dav	Long	Low	C					
	Day		High	В					
			High	С					
_		<u>т</u>	actical Move						
			Low	A	- to reflect your				
	Day	Short	Low	В					
	,		High		- test execution				
-		<u> </u> т	High	В					
		1		L C	_				
			High	B					
	Night	Medium	Low	B	Failing to correctly design				
			Low	A					
		T.	actical Move		and account for execution				
			High	A	order could lead to wrong				
	Day	Medium	High	В					
			Low	C	conclusions.				
			Low	A					
		т.							
Γ	Night	Medium	High	A	- Sub-Plots				
Vhole-Plot -			Low	В					
	-		Low	A					
L—			High						
		1		B					
Lighting and			low	B	Angle of Fire and				
Lighting and	Night	Long	High	C	Aligie of File allu				
Range are			Low	C	Projectile Type are				
hard to change		, T	actical Move		easy to vary from				
			Low	В					
petween	Night		Low	C	one fire mission to				
Tactical Moves			High	С	the next				
			High	C C					
_		T:	actical Move						
			High		_				
	Day	Medium	High	В					
			High	A A	TestScience				
			LOW	Г В	Data . Driven . Defense				

# Which Points? Testing for Problem Cases in a Software System Application

System required to pull and integrate information from multiple sources and display via applications.

There are 22 applications total; personnel typically would not need to open more than 8 at a time.

Are there any combinations of applications that do not work together?





# **Combinatorial Designs** are useful for finding bugs and conditions where the system breaks.

C#	W 1	W 2	W 3	W 4	W 5	W 6	W 7	W 8	W 9	W 10	W 11	W 12	W 13	W 14	W 15	W 16	W 17	W 18	W 19	W 20	W 21	W 22
1	0	0	0	1	1	1	1	0	1	0	1	0	0	1	1	0	0	0	0	0	0	1
2	0	0	1	0	0	0	0	1	1	0	1	1	1	1	0	1	0	0	0	0	0	1
3	0	1	0	0	0	0	0	0	1	0	1	1	0	0	1	1	1	1	0	0	0	1
4	0	1	1	0	0	0	1	1	0	0	1	0	1	0	1	0	1	0	0	0	0	1
5	1	0	0	0	0	0	0	0	0	0	1	1	1	0	1	0	0	0	1	1	1	1
					•													•				
70	1	1	0	0	0	0	1	1	0	1	1	1	0	1	0	0	0	0	0	0	0	1
71	1	1	1	0	0	0	1	0	1	1	0	1	1	0	0	0	0	0	0	0	0	1
72	0	0	0	1	0	1	0	1	0	1	1	0	0	0	0	1	1	0	1	0	0	1
73	0	0	0	0	1	0	0	1	1	1	0	0	1	0	1	1	1	0	0	0	0	1
74	0	0	0	1	0	0	1	1	1	0	0	1	1	0	1	1	0	0	0	0	0	1
75	0	0	0	1	1	0	0	0	0	1	1	1	1	1	0	0	1	0	0	0	0	1

BUT because of their inability to determine cause and effect, Combinatorial Designs are not the best choice for Operational Tests.

## Software programs assist testers in selecting and constructing the test designs

skprGUI								×
skpr <b>GUI</b>	Results					𝔗 Save State	<b>?</b> Tutorial	R
Generate Evaluate Design Design	Design Design Evaluation Generating Code				Generating Code			
Basic Advanced Power Trials	De • o	ESIGN rder Desig	n					- 1
36		Range	Angle	TimeofDay	Projectile			
Model	1	Medium	High	Day	В			
~ Range + Angle +TimeofDay + Projectile + I	2	Long	High	Day	А			
Number of Factors	3	Medium	Low	Day	В			
4	4	Short	High	Night	A			
-	5	Medium	Low	Day	С			
	6	Long	High	Night	С			
Factor 1	7	Long	High	Night	А			
Changes Type	8	Short	High	Night	С			
Easy   Categorical	9	Short	Low	Day	В			
Name	10	Medium	High	Night	А			
	11	Short	Low	Night	С			•

skpr: Design of Experiments Suite: Generate and Evaluate Optimal Design

## **Summary of Test Designs**

Type of DOE	Size	Size Scales With	Uses	Robustness
Factorial Designs (Full factorial, fractional factorial, factorial with center point, replicated factorial)	Large to medium	Number of factors	Estimate performance at extremes/"corners" of the design space	Very robust to lost data points
Optimal Designs (D- optimal, I-optimal)	Small	Nothing	Design minimal tests when there is minimal risk to lost data, and factors affecting performance are well understood	Fragile design; lost data can cause major problems
Response Surface Design	Large	Number of factors, flexibility of response surface	Estimate performance with non-linear factor effects and interactions	Somewhat robust to lost data points
Hierarchical Designs (blocking, split-plot designs)	Medium to small	Number of groups/blocks, group/block size	Properly randomize, while accounting for natural grouping in your data	Somewhat robust to lost data points; less robust to losing full groups of data points
Combinatorial Designs	Large	Number of factors, size of combinations	Find problems through systematically testing all possible cases	Not robust; must check all combinations





## **Test Design**

#### Lesson 4: Evaluate the Design

### **Institute for Defense Analyses**

4850 Mark Center Drive • Alexandria, Virginia 22311-1882

## When evaluating a test design ask:

- 1. Does this design ensure that **adequate data** are collected?
- 2. If the system doesn't perform well, will I be able to **determine why**?
- 3. Does the DOE reflect **the way the test will be conducted**?
- 4. Will the data collected let me do the **analysis** I want to do?


### When evaluating a test design ask:

1. Does this design ensure that **adequate data** are collected?

**Power and Confidence** 







True Combat Performance













# Power curves are good ways to visualize and compare power for different test designs



Data . Driven . Defense

### When evaluating a test design ask:

2. If the system doesn't perform well, will I be able to **determine why**?

Correlation and Identifiability





# Tests should be designed so that we can identify the situations or missions where performance drops off

Poor performance during the day?

Poor performance against wheeled targets?

Poor performance against wheeled targets during the day?

Time of Day	Target Speed	Target Type	Update Rate	Clutter
Day	Fast	Wheeled	30	Y
Day	Fast	Wheeled	12	Ν
Day	Fixed	Wheeled	12	Ν
Day	Fixed	Wheeled	30	Ν
Day	Slow	Wheeled	12	Y
Day	Slow	Wheeled	30	Y
Night	Fast	Tracked	30	Ν
Night	Fast	Tracked	12	Y
Night	Fixed	Tracked	30	Y
Night	Fixed	Tracked	12	Y
Night	Slow	Tracked	12	Ν
Night	Slow	Tracked	30	Ν

#### Narrow Radius Explosive\* Test Design

\*Notional

# Use correlation maps to quickly visualize how hard it will be to distinguish between factors



- Blue means perfectly uncorrelated (No problems!)
- Red is perfectly (100%) correlated (Cannot differentiate)
- Colors in between indicate some degree of correlation (Can differentiate, but estimates may be biased/incorrect)



### Ideally, your test design will have factors completely independent to ensure identifiability

Target Speed 1 Target Speed 1 Target Type Update Rate Clutter	
Time of Day Target Speed 1 Target Type Update Rate Update Rate	
	<b> r </b> 0
	1

	Time of Day	Target Speed	Target Type	Update Rate	Clutter
	Day	Fast	Tracked	30	Y
	Day	Fast	Wheeled	12	Ν
	Day	Fixed	Tracked	12	Ν
	Day	Fixed	Wheeled	30	Ν
	Day	Slow	Tracked	12	Y
	Day	Slow	Wheeled	30	Y
	Night	Fast	Tracked	30	Ν
	Night	Fast	Wheeled	12	Y
	Night	Fixed	Tracked	30	Y
	Night	Fixed	Wheeled	12	Y
	Night	Slow	Tracked	12	Ν
7	Night	Slow	Wheeled	30	N

### When evaluating a test design ask:

3. Does the DOE reflect **the way the test will be conducted**?

## Restricted Randomization and Disallowed Combinations





## **RECALL:** Our test of the new cannon includes factors that are **Difficult to Completely Randomize** and have **Disallowed Combinations**

		$\boldsymbol{\mathcal{C}}$	
Factor	Level	Difficult to Randomize	Disallowed Combinations
Time of Day	Day, Night	Hard	Night and Short Range
Range to Target	Short, Medium, Long	Hard	
Projectile Type	А, В, С	Easy	A and Long Range
Angle of Fire	Low, High	Easy	
Traverse Angle	0-15, 16-30, 31-45	Easy	
000			
			TestSc

Data . Driven . Defense

Γ	Time of Day	Range	Angle of Fire	Projectile Typ	e	
Γ			Low	В		
	Dav	Long	Low	C	Split-plot designs	
	Day		High	В		
L			High	C	allow your DOE to	
_		<u>т</u>	actical Move	1	roflact your tast	
			Low	A		
	Day	Short	Low	В	execution	
	,		High			
		<u> </u> т.	High	В	_	
		1		L C	_	
			High	B		
	Night	Medium	Low	B	Eailing to correctly design	
			Low	A		
		r Ta	actical Move	<u> </u>	and account for execution	
		Medium	High	A	order could lead to wrong	
	Davi		High	В	conclusions.	
	Day		Low	C		
			Low	A		
		T	actical Move			
[]	Night	Medium	High	A		
Nhole-Plot -			Low	В	- Sub-Plots	
			Low	A		
L		<u> </u> т.	High		J	
		1		B		
Lighting and			Low	B	Angle of Fire and	
	Night	Long	High	C	Aligie of File allu	
Range are			Low	C	Projectile Type are	
hard to change		, Ta	actical Move	1	easy to vary from	
			Low	В		
petween	Night	Long	Low	С	one fire mission to	
Tactical Moves			High	С	the next	
			High			
		T	actical Move	-		
	Day	Medium	High	C C		
			High	В		
			High	A	TestScience	
L			LOW	ГВ	Data . Driven . Defense	

### When evaluating a test design ask:

4. Will the data collected let me do the **analysis** I want to do?

Statistical Model Supported





Determine what you want to learn from the test and build the design accordingly

## Factors?

When resources are constrained:

- 1. Carefully describe the operational conditions in which the system will be employed, and
- 2. Make sure to get the data points required to fully describe that space.

## Interactions?

Curvature?



# DOE supports getting the information you need in the most efficient way possible









# **Test Design**

#### Lesson 5: Manage and Store the Test Data

#### **Institute for Defense Analyses**

4850 Mark Center Drive • Alexandria, Virginia 22311-1882

## What is Data Management?

*Data Management* is the development and execution of policies and practices that organize, protect, and control data and information throughout their lifecycle\*

- Access restrictions
- Naming conventions
- Software and hardware systems for storage, search, and discovery
- Roles and responsibilities for data management
- Etc.

## Data Management plays a key role throughout the T&E Lifecycle



#### Each organization must define its own goals for data management

FAIR – Nature article identifies four properties (Findable, Accessible, Interoperable, Reusable) as goals for research data





**Reproducibility** – Manage data to facilitate Reproducible Research\*

Security – Ensure that classified, proprietary, and PII data are handled appropriately

\* See, e.g., Flack, Kirshenbaum, and Haman tutorial from 2019 DATAWorks, available at www.testcience.org/archive







Data . Driven . Defense

### **Data Acquisition**

- Storage volume
- Backup options
- Preserving pedigree
- Metadata
- Security





## **Data Cleaning**

- Version control
- Documentation
- Intermediate data products
- Naming conventions
- Automation and repetition





### **Data Analysis**

- Version control
- Documentation
- Organization
- Reproducibility





## **Data Archiving**

- Future users
- Findability
- Searchability
- Metadata
- Stewardship





### Just the beginning...

# DMBOK identifies 10 key data management functions:



Data are managed – whether through conscious, deliberate process or not.

By building a data management strategy and executing on it, organizations can ensure that their people, practices, and systems are helping them to get the most out of limited data resources and do the best analysis possible.

REPORT DOCUM	ENTATION PAGE	OMB No. 0704-0188		
Public reporting burden for this collection of information is e maintaining the data needed, and completing and reviewin suggestions for reducing this burden to Department of Defe Suite 1204, Arlington, VA 22202-4302. Respondents sho of information if it does not display a currently valid OMB co	estimated to average 1 hour per response, including the g this collection of information. Send comments regard ense, Washington Headquarters Services, Directorate Id be aware that notwithstanding any other provision of ontrol number. <b>PLEASE DO NOT RETURN YOUR FC</b>	e time for reviewing instructions, searching existing data sources, gathering and ding this burden estimate or any other aspect of this collection of information, including for Information Operations and Reports (0704-0188), 1215 Jefferson Davis Highway, of law, no person shall be subject to any penalty for failing to comply with a collection <b>DRM TO THE ABOVE ADDRESS</b> .		
1. REPORT DATE (DD-MM-YYYY)2. REPORT TYPE04-2021IDA Publication		3. DATES COVERED (From - To)		
4. TITLE ANDSUBTITLE		5a. CONTRACT NUMBER Separate Contract		
Test Science Website Videos - Test Pla	anning and Test Design	5b. GRANT NUMBER		
		5c. PROGRAM ELEMENT NUMBER		
6. AUTHOR(S)		5d. PROJECT NUMBER C9082		
Kelly M. Avery (OED); Caitlan A. Fealing (OED); Rebecca M. Medlin (OED); Rachel A. Haga (OED); Matthew R. Avery (OED);		5e. TASK NUMBER		
		5f. WORK UNIT NUMBER		
7. PERFORMING ORGANIZATION NAME( Institute for Defense Analyses 4850 Mark Center Drive Alexandria, Virginia 22311-1882	S) AND ADDRESS(ES)	8. PERFORMING ORGANIZATION REPORT NUMBER NS-D-16393 H 2020-000425		
<b>9. SPONSORING / MONITORING AGENCY</b> Institute for Defense Analyses 4850 Mark Center Drive Alexandria, Virginia 22311-1882	Y NAME(S) AND ADDRESS(ES)	10. SPONSOR/MONITOR'S ACRONYM(S) IDA		
		11. SPONSOR/MONITOR'S REPORT NUMBER		
<b>12. DISTRIBUTION / AVAILABILITY STAT</b> Approved for Public Release. Distribut	EMENT ion Unlimited.	I		
13. SUPPLEMENTARY NOTES				
<b>14. ABSTRACT</b> Understanding the test design process i element of this process and requires co	s fundamental to conducting efficien llaboration and iteration among subj	nt and effective tests of systems. Planning is a critical ject matter experts. Test planning involves identifying		

**-** -

Δ.,

- -1

element of this process and requires collaboration and iteration among subject matter experts. Test planning involves identifying the test objectives, determining the appropriate outcome measure(s), and identifying what factors may influence those outcomes. The first set of website videos introduces the overarching test design process and provides best practices and examples for each of the three planning steps.

Test design is the process of actually creating and evaluating a set of experimental runs. It entails careful consideration of the test goals to ensure objectives can be met by the data collected, building the corresponding type of design, and then evaluating it across multiple criteria. The test design should have a large enough sample size to be useful, be executable in the real world, and support the desired analysis. The second set of videos describes and provides recommendations for each element of the test design process, and also introduces considerations for how to manage and store data collected from tests.

**15. SUBJECT TERMS** 

Design of Experiments (DOE); OT&E; Statistics; Test Planning

16. SECURITY CLASSIFICATION OF:			17. LIMITATION	18. NUMBER	<b>19a. NAME OF RESPONSIBLE PERSON</b>
			OF ABSTRACT	OF PAGES	Rebecca Medlin (OED)
a. REPORT Unclassified	<b>b. ABSTRACT</b> Unclassified	<b>c. THIS PAGE</b> Unclassified	Unlimited	171	<b>19b. TELEPHONE NUMBER</b> (include area code) (703) 845-6731