

Evaluating and Improving
Test Effectiveness
Using Statistical Test Optimization

Kedar M. Phadke
Phadke Associates
kedar@phadkeassociates.com

“Approximately 50% of programs entering Initial Operational Test and Evaluation (IOT&E) have not been evaluated as Operationally Effective and Operationally Suitable.”

Hon. Kenneth Krieg
Under Secretary of Defense
Acquisition, Technology, and Logistics
April 30, 2007

Defense Science Board (DSB) study from 2009 had similar findings.

Industry Perspective of the T&E Challenge

Raytheon

The Test Optimization Challenge

“We are being challenged by our customers and by the marketplace to develop and deliver increasingly complex systems with smaller performance margins that meet the user’s requirements in the shortest time, with high reliability, open and adaptable, and at the lowest cost.”

Given this challenge, there is more pressure than ever on Integration, Verification & Validation activities to deliver performance results on time and within budget.

Industry studies have estimated test and rework to represent between 30 and 50% of product development costs. Given this investment, test represents fertile ground for optimization techniques. Typical benefits of statistically-based test optimization include:

- Increased Mission Assurance
- Optimized performance
- Improved cycle time
- Increased Productivity
- Reduced cost

Raytheon IDS Quality & Process
Performance Objectives established for:

- On-time Deliverables
- Cost and Schedule Performance
- Engineering Productivity
- Reducing Rework

Weapons Fire and Detection
System Program (WFDS)

**Note: Information has been changed for this
presentation**

Program Summary

- Program was upgrading a system that has capability to detect type and direction of enemy fire under many operating scenarios. Operating scenarios include different platforms, arms types, times of day, weather conditions, single/multiple threats, mounting methods, etc.
- The program plans to operationally evaluate prototypes to demonstrate performance for the next phase of acquisition.
- T&E Questions
 - Does the test plan provide high confidence that the delivered prototype will provide required performance (Is the test plan truly validating performance of the system under all intended operating conditions)? What is the risk?
 - What corrective actions can technical staff take to improve efficacy of the test plan to reduce risk?
 - Is there a more optimal way to test the weapons targeting system to reduce test cost and schedule?

Program Status

- **Already developed a suite of tests (71 test scenarios)**
 - **Used best practices, expert guidance, and legacy test scenarios**
- **Systematic review of the tests uncovered multiple operational test gaps**
 - **750+ test gaps identified – all deemed important**

Optimization

- **Identify the key test parameters and their values**
- **Generate a factor table and test design**
- **Ensure that test design covers conditions of the original 71 test and eliminates test gaps**
- **Do all this while optimizing risk, cost, and schedule!**
 - **It should be easy for program staff to implement**

Test Factors and Values

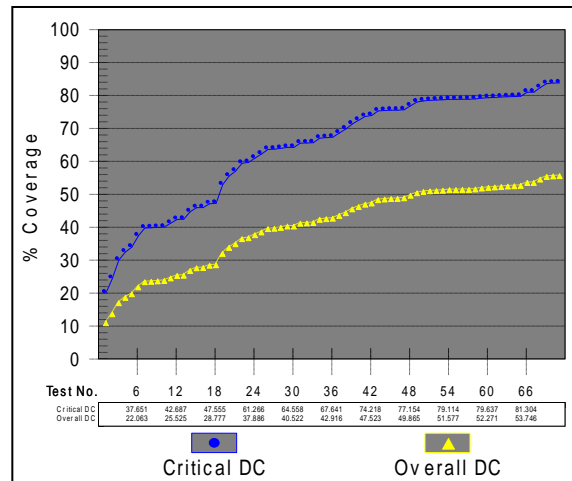
Factor Name	No. Lvl's	Dep. On	Level 1	Level 2	Level 3	Level 4	Level 5	Level 6	Level 7
Ambient Temperature, degrees F	3		-51	54.5	160				
Time of Day	2		Daylight	Night					
Weather Conditions	4		Clear Day	Cloudy	Rain	Sand/Dust			
Sensor position	2		on the move	stationary					
Platforms/Mounting Type	7		Ground - Tripod	Ground - UGV	Air - Fixed wing	Ground - Yoke	Ground - Gimble	Air - UAV	Air - Rotary Wing
Weapons System Mount	2		Not attached to Friendly Weapons System	Attached to Friendly Weapons System					
Configuration	2		stand alone sensor	multiple sensor with cross cueing and validation					
Power Source	2		Battery Power	Standard Military Power					
5.56mm threats	4		1 present	6 present	2 present	10 present			
5.56mm threat location	4		outside of range with natural barriers	within range with natural barriers	outside of range with manmade barriers	within range with manmade barriers			
7.62mm threats	4		1 present	2 present	6 present	10 present			
7.62mm threat location	4		outside of range with natural barriers	within range with manmade barriers	within range with natural barriers	outside of range with manmade barriers			
50 Cal threats	3		1 present	4 present	2 present				
50 Cal threat location	4		within range with natural barriers	outside of range with natural barriers	outside of range with manmade barriers	within range with manmade barriers			
RPG threats	2		1 present	2 present					
RPG threat location	4		within range with natural barriers	within range with manmade barriers	outside of range with manmade barriers	outside of range with natural barriers			
Mortar threats	2		1 present	2 present					
Mortar threat location	4		within range with manmade barriers	outside of range with natural barriers	within range with natural barriers	outside of range with manmade barriers			
MANPAD threats	2		1 present	2 present					
			outside of range	within range with	within range with	outside of range			

- Test parameters include sensor configurations, threat profiles, interfacing systems, output types, and environmental variables

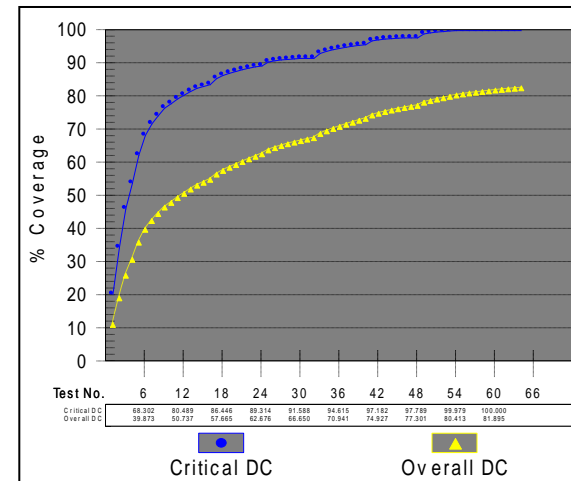
WFDS Program Benefits using Statistical Test Optimization

PHADKE

Baseline – Original 71 tests



Statistically Opt. – 64 tests



- Improved test effectiveness (Identified and eliminated 750+ test gaps)
- Reduced test cost and schedule by 10%
- Provided a template for future deployments

What Metrics Work
for Statistical Test Optimization

Criteria for Evaluating Test Metrics

- **Metrics should be related to Program Objectives**
 - **Operational Effectiveness**
 - **Cost**
 - **Schedule**
- **If the Statistical Metrics are not correlated with Program Objectives, we have a problem!**

What Happens When We Use Statistical Test Power

- This metric is mentioned often in DOT&E literature
- Key factors that influence statistical test power
 - Acceptable error (the alpha)
 - The magnitude of the effect size you want to measure
 - Standard deviation for the population
 - Sample size
- Do we have a good handle on any of these ahead of time?
- For the WFDS example, the test power (assuming alpha = 0.10, effect size/standard deviation = 0.33) was:
 - Original 71 Test Scenarios: 0.97 ← Metric looks better for original tests
 - Optimized 64 Test Scenarios: 0.95
 - Test Power calculations conducted using the JMP® Software tool from the SAS corporation

- For WFDS, statistical test power did not correlate with operational effectiveness, cost or schedule

What happens when we use Domain Coverage metric?

- Domain coverage is a measure of how well tests cover requirements and interop of requirements
 - Commonly broken into two pieces: Critical Domain Coverage (CDC) and Overall Domain Coverage (ODC)
 - CDC Definition: How well you are covering individual requirements and interop of pairs of requirements
 - ODC Definition: How well you are covering individual requirements, interop of pairs, interop of triples, etc
 - For WFDS, the Domain Coverage:
 - Original 71 Test Scenarios: CDC: 84%, ODC: 55.8%
 - Optimized 64 Test Scenarios: **CDC: 100%, ODC: 82.6%**
- Metric looks better for optimized tests

Comparing Statistical Test Power and Domain Coverage

Test Design 1

		Threat location			
		within range with natural barriers	outside of range with natural barriers	outside of range with manmade barriers	within range with manmade barriers
50 Cal threats	1 present	XXXXX X	XXXXX X		
	2 present	XXXXX X	XXXXX X		
	4 present				

Test Design 2

		Threat location			
		within range with natural barriers	outside of range with natural barriers	outside of range with manmade barriers	within range with manmade barriers
50 Cal threats	1 present	X	X	X	X
	2 present	X	X	X	X
	4 present	X	X	X	X

What happens when we calculate test power?

Test Power of Design 1 = Test Power of Design 2

Two completely different test designs, same test power

What happens when we calculate Domain Coverage?

DC of Design 1 << DC of Design 2

DC discriminates between the test designs and favors greater test coverage and effectiveness



Defense Industry has seen significant benefits using Domain Coverage metric



Deployment Results Summary

<u>Test</u>	<u>Original Test Plan</u>	<u>Optimized Test Plan</u>
Subsystem Testing	28 Tests	8 Tests (71% reduction)
Systems Mission Testing	25 Missions	18 Missions (28% reduction)
Subsystem Simulation	100 Runs	40 Runs (60% reduction)
Range Testing	1036 Tests	632 tests (39% reduction)
Software Subsystem Testing	90 Tests	63 Tests (30% reduction)
System Scenario Generation	8 Missions	6 Missions (25% reduction)
System MOE Testing	1600 Tests	885 tests (45% reduction)
System Testing	246 Tests	48 tests (80% reduction)

In each case, the reduction in number of test cases was achieved while maintaining or improving upon existing test coverage.

Takeaways

- **Improved Test Effectiveness and Efficiency for WFDS using Statistical Test Optimization**
 - 71 tests reduced to 64 (10% reduction in cost and schedule)
 - Eliminated 750+ test gaps
- **Use statistical metrics that correlate with program metrics: test effectiveness, cost, and schedule**
 - If they don't correlate or discriminate properly, they risk driving your programs in the wrong direction.

References

- **Statistical Test Power calculations conducted using SAS Institute's JMP® Software. JMP is a registered trademark of SAS.**
- **Domain Coverage calculations conducted using Phadke Associates rdExpert™ Test Suite Software. rdExpert is a trademark of Phadke Associates, Inc.**
- **Reference material and articles available on <http://education.phadkeassociates.net>**