



Framework For Understanding Experiment Design Validity *(..with potential applicability to OT?)*

Presentation to DOTE/OTA DoE Working Group
20 February 2009

Rick Kass
GaN Corporation
Contract Technical Support to –
Transformation Technology Directorate (TTD)
US Army Operational Test Command (USAOTC)
US Army Test and Evaluation Command (USATEC)

Army Test and Evaluation Command

Validity fair **Control** others, as you would have them control you

Freudian **Analysis**

Discovery channel “Nightmare on **Bias** street” Rebel without a **Causality**

“It was the best of **trials**. It was the worst of trials.”

Data collection “**Hypothesis** now” **Realism** **Significance** other

Who Wants to Be an Experimenter?

Bite-size sample size “I think, therefore I **experiment**”

I **experimented**, but did not inhale. **Testing**, 1, 2, 3...

In search of Bobby **Findings** **Serendipity** do da

Free play Welcome back **Concepts** **Science** fiction

Independent, Dependent, and Republican **Variables** **Risk-free** Alpha and Beta bonds

Precision E=MC² (Experiment equals Methodology Controlled by Confusion)



What is “Experiment Rigor”

To justify recommendations
...need “credible experiments”
...experiments with high degree of rigor (...validity...)

How to increase experiment rigor...?

- ...according to some..
 - ...more realistic scenarios...
 - ...better portrayal of unrestricted and adaptive threats...
 - ...more quantifiable results...
- ...according to others...
 - ...fewer variables with better control...
 - ...more trials to increase sample size
 - ...use of randomization

Experiment (..and Test ...) Rigor is all of these and more...

3

Army Test and Evaluation Command



Outline

Experiment References

Experiment Logic: 2-3-4-5-21

Experiment Rigor Requirements

- Threats to Rigor
- Good practices to counter threats

Implications of Logic

- Design of individual experiments
- Campaign of experiments

4

Army Test and Evaluation Command



Useful Definition of Experiment



35 different definitions at “[WWW. One-Look Dictionary Search](#)”

Common Themes:

A test done in order to learn something or to discover whether something works or is true (Cambridge Advanced Learning Dictionary). An operation carried out under controlled conditions in order to discover an unknown effect or law, to test or establish a hypothesis, or to illustrate a known law (Merriam-Webster Dictionary)

Experiment –

“To explore the effects of manipulating a variable.”

Shadish, Cook, & Campbell,. Experimental and Quasi-Experimental Designs for Generalized Causal Inference (p. 507)

Warfighting Experiment —

“A systematic process to explore the effects of manipulating warfighting capabilities or conditions.”

Army Test and Evaluation Command

5



Experiment Rigor References



William R. Shadish, Thomas D. Cook and Donald T. Campbell,. Experimental and Quasi-Experimental Designs for Generalized Causal Inference (Houghton Mifflin Co; 2002)

Thomas D. Cook and Donald T. Campbell,. Quasi-Experimentation: Design and Analysis Issues (Rand McNally, 1979)

Donald T. Campbell and Julian Stanley. Experimental and Quasi-Experimental Designs for Research (Rand McNally, 1963)

Experiment rigor requirements based on 40 years of writing about non-laboratory experiment requirements.
Adapted ideas and terminology for warfighting experiments

Apply traditional scientific principles to Warfighting Experimentation in innovative ways

Army Test and Evaluation Command

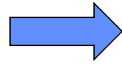
6



Outline



Experiments References



Experiment Logic: 2-3-4-5-21

Experiment Rigor Requirements

- Threats to Rigor
- Good practices to counter threats

Implications of Logic

- Design of individual experiments
- Campaign of experiments



Experiment Hypotheses

“educated guesses of what might happen”



Useful:

- Help to clarify what experiment is about
- Identify logical thread of the experiment
- Guide experiment design and data collection

Nothing magic:

If _____; **then** _____.

proposed solution(s) → **problem to be overcome**

independent variable → dependent variable

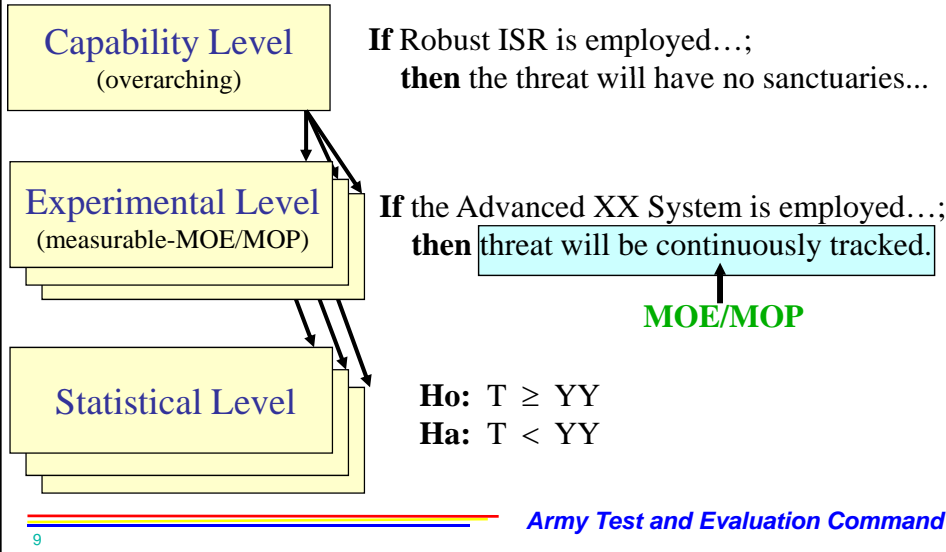
potential cause → possible effect

- Sea Basing → Rapid deployment
- Collaboration → Adaptive planning
- Global Cell → Inter-theater coordination
- Robust ISR → Deny sanctuaries

“Two parts to experiment hypotheses”



Different Levels of Hypotheses



Logic of Hypothesis Resolution



A **B**

If proposed solution : then problem to be overcome (effect)

Logic of hypothesis resolution

“Three parts to resolving hypotheses”

3

1. Did A occur?
2. Did B occur?
3. Was B due to A ?

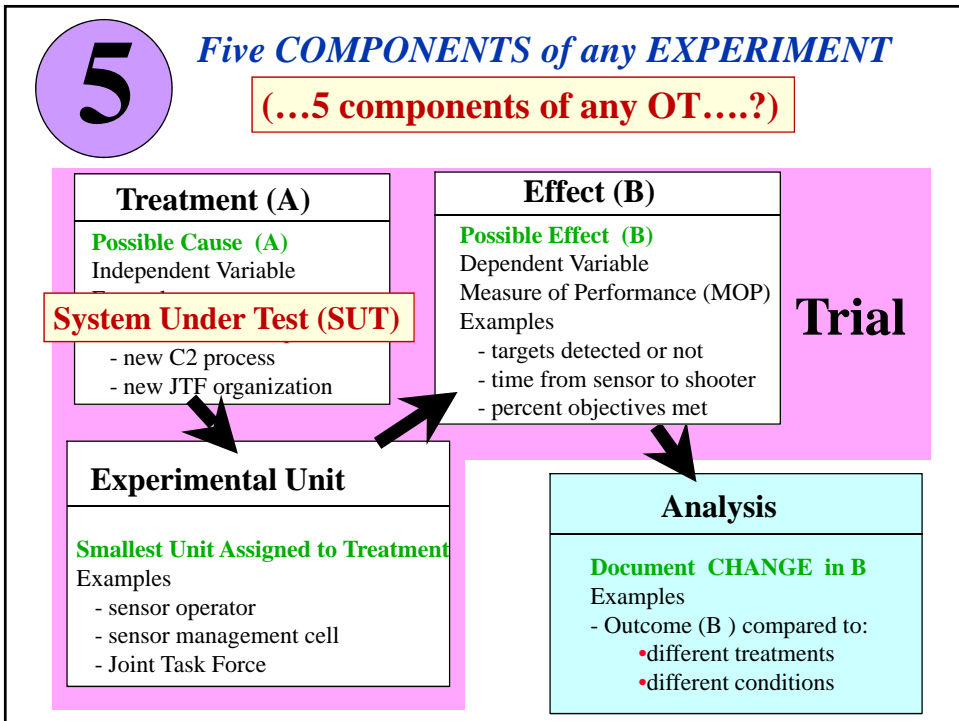
Internal Validity of an experiment

4 *Four Requirements for Rigorous (valid) Experiment*

If New Capability (A) ; Then Effect (B) .

	Requirement	Evidence for Validity	Threat to Validity
1	ability to use new capability	A occurred	Asset did not work or was not used
2	ability to detect change in effect	B ch	Too much noise, can not detect any change
3	ability to isolate for change	A alone caused B	Alternate explanations of change available
4	ability to relate results to actual operations	Change in B due to A is expected in actual operations	Observed change may not be applicable

Applicable to OT?




21 Threats to a Rigorous Warfighting Experiment


5 Experiment Components	4 Experiment Requirements				
	1. Ability to Use the Capability	2. Ability to Detect Change	3. Ability to Isolate the Reason for Change		4. Ability to Relate the Results to Operations
			Single Group	Multiple Groups	
1. Treatment	(1) Capability functionality does not work.	(5) Capability systems vary in performance.	(11) Functionality changes across trials.		(18) Functionality does not represent future capability.
2. Players	(2) Players are not adequately prepared.	(6) Experiment players vary in proficiency.	(12) Player proficiency changes across trials.		(19) Players do not represent operational unit.
3. Effect	(3) Measures are insensitive to capability impact.	(7) Data collection accuracy is inconsistent across trials.	(13) Data collection accuracy changes across trials.	(16) Data collection accuracy differs for each group.	(20) Measures do not reflect important effects.
4. Trial	(4) Capabilities are not tested in no opponent conditions.	(8) Data collection accuracy is inconsistent across trials.	(14) Trial conditions change across trials.	(17) Groups operate under different trial conditions.	(21) Scenario is not realistic.
5. Analysis	(9) Low statistical power (10) Statistical assumptions are violated.		<ul style="list-style-type: none"> •Experiment Hypothesis: if <u>A</u>, then <u>B</u>. •Purpose of an experiment: verify that <u>A</u> causes <u>B</u>. •Valid experiment allow conclusion “<u>A</u> causes <u>B</u>” to be based on evidence and sound reasoning... -By reducing or eliminating 21 threats to validity. 		

Applicable to OT?

21



Outline



Experiments References

Experiment Logic: 2-3-4-5-21

Experiment Rigor Requirements

Threats to Rigor

Good practices to counter threats

Implications of Logic

Design of individual experiments

Campaign of experiments

Army Test and Evaluation Command



Four Requirements To Design Rigorous Warfighting Experiments

Internal Validity

1. **Capability Used**
2. **Detection of Change in Effect**
3. **Isolation of Reason for Change**

External Validity

4. **Relating Results to Military Operations**

1. Ability to Use the New Capability

Most consistent “lessons learned” reported after warfighting experiments completed:


- *New Capability did not work as well as promised.*
- *Players did not know how to use it properly.*
- *The measure (instrumentation) was not sensitive to its use.*
- *The scenario play did not give the players the opportunity to use.*

...sounds familiar for OT?


Ensuring that the experimental capabilities are used and can make a difference is the first logical step in designing a valid experiment.

<i>Threats to the Ability to Use the Capability</i>	
THREAT	PREVENTION
<p>Treatment</p> <p>1. Capability functionality does not work. Does the HS/SW work?</p>	<ul style="list-style-type: none"> • Ensure functionality of experimental capability is present.
<p>Unit</p> <p>2. Players are not adequately prepared. Do the players have the training and TTP to use the capability?</p>	<ul style="list-style-type: none"> • Ensure player <u>organized</u>, <u>equipped</u>, and <u>trained</u> for capability use. • Provide sufficient SOPs for capability use. • Provide pre-experiment "practice time."
<p>Effect</p> <p>3. Measures are insensitive to impact Is the output sensitive to capability use?</p>	<ul style="list-style-type: none"> • Pilot-test impact on experiment outcome • "Verify" model input-output logic
<p>Trial</p> <p>4. Capability has no opportunity to perform. Does the scenario and MESL call for capability use?</p>	<ul style="list-style-type: none"> • Pilot-test scenario and MSEL • "White cell" specific scenario injects and monitor for use

Applicable to OT?



Four Requirements To Design Rigorous Warfighting Experiments



Internal Validity

1. Capability Used
2. Detection of Change in Effect
3. Isolation of Reason for Change

External Validity

4. Relating Results to Military Operations

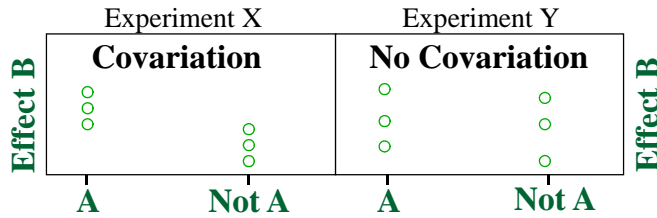
18
Army Test and Evaluation Command

2. Ability to Detect Change in the Effect

- Given that A was employed
- Next Question: **Did B (effect) change when A was applied ?**

Ability to detect change in B: Statistically Valid Experiment

“Detect Change” = “Detect COVARIATION:” B changes when A applied



Two Groups of Threats to Detecting Change

- Fail to Detect Real Change
 - Incorrectly see no covariation (**Type II Error, Producer Risk, Beta Error**)
- Incorrectly Detect Change--
 - Incorrectly see covariation (**Type I Error, Consumer Risk, Alpha Error**)

2. Ability to Detect Change-- statistical validity

Threats

Ability to detect change is enhanced as variability is reduced

Fail to Detect Change	<p>Treatment 5. Capability Systems vary in performance</p> <ul style="list-style-type: none"> • Continual fluctuation in reliability <p>Unit 6. Players vary in performance</p> <ul style="list-style-type: none"> • Different levels of training • Different reasons for use <p>Effect 7. Data collection accuracy inconsistent</p> <ul style="list-style-type: none"> • Variation in collectors <p>Trial 8. Trial conditions</p> <ul style="list-style-type: none"> • Inadvertent <p>Analysis 9. Inadequate power</p> <ul style="list-style-type: none"> • Inconsistent alpha risk (1% or 5%) • Inefficient statistical test 	<ul style="list-style-type: none"> • Hold constant • Examine only subset of population • Inconsistent data collection versus data collectors • Experienced data collectors • Set boundary conditions • More repetitions • Increase alpha risk (to 10%) • Use paired comparisons
Incorrectly Detect Change	<p>10. Statistical Test Assumptions Violated</p> <ul style="list-style-type: none"> • Some statistical techniques have sensitive assumptions • Error rate problem (fishing) • Large number of statistical tests 	<ul style="list-style-type: none"> • Use appropriate statistical tool • Select fewer, more meaningful MOPs

Applicable to OT?



Four Requirements To Design Rigorous Warfighting Experiments

Internal Validity

1. Capability Used
2. Detection of Change in Effect
3. Isolation of Reason for Change

External Validity

4. Relating Results to Military Operations

3. Isolating the Reason for Change

- Given that **A** was employed
- Given that **B** changed as **A** was applied
- Next Question: **What really produced the change in B?**

Design Validity -- A alone caused change?

- Threat -- Something else

...also a good 2-type division for OT?

depends on type of experimental design

Single Group Design

One unit receives all treatment conditions

		Scenario 1	Scenario 2
Same unit ↓	Unit C with Current		
	Unit C with Future		

Compare group under different conditions

Multiple Group Design

Different units receive different treatment conditions

		Scenario 1	Scenario 2
Different units ↓	Unit C with Current		
	Unit D with Future		

Compare group to another group

- Side-by-side baseline
- Side-by-side "shoot off"

3. Isolating the Reason for Change in SINGLE-GROUP DESIGNS

Sequence of trial presentation
is critical consideration

Sequence 1: Unbalanced

Mon	Tue	Wed	Thu
Current	Current	Future	Future

(1+0=1) (1+1=2) (1+2=3) (1+3=4)
Current=3 Future=7

Sequence 2: Balanced

Mon	Tue	Wed	Thu
Current	Future	Current	Future

(1+0=1) (1+1=2) (1+2=3) (1+3=4)
Current=4 Future=6

Sequence 3: Counterbalanced

Mon	Tue	Wed	Thu
Current	Future	Future	Current

(1+0=1) (1+1=2) (1+2=3) (1+3=4)
Current=5 Future=5

$(1 + 0 = 1)$
Treatment Effect Learning Effect Observed Effect

**In single-group design,
order effect generates greatest threat
to Isolating Reason for Change**

3. Isolating the Reason for Change SINGLE-GROUP DESIGN ORDER EFFECTS

	THREAT		PREVENTION
Treatment	<p>11. Capability Functionality changes across trials System or process improves or degrades over time</p>		<ul style="list-style-type: none"> • Use fixed configuration
Unit	<p>12. Player Proficiency changes across trials Performance improves during later trials due to experience rather than treatment present</p>		<ul style="list-style-type: none"> • Train to maximum prior to start
Effect	<p>13. Data Collection Accuracy across trials Data collection improve or degrade over time or changing results</p>		<ul style="list-style-type: none"> • Train data collectors to maximum performance prior to start • Check and recalibrate instrumentation after each trial
Trial	<p>14. Trial conditions change across trials Weather, OPFOR, and simulations improve or degrade over time</p>		<ul style="list-style-type: none"> • Train OPFOR to maximum performance prior to start
		<p>General prevention/check</p> <ul style="list-style-type: none"> • Counterbalance presentation sequence • Check for increase/decrease over time 	

Single-group design validity is enhanced as **unintended changes over time are controlled**

Applicable to OT?

3. Isolating the Reason for Change in MULTIPLE-GROUP DESIGNS

- Different player units receive different treatments

	Scenario 1	Scenario 2
Unit C with Current		↓ B ₁
Unit D with Future		↑ B ₂

- Previous Order-Effect threats are neutralized
 - if same sequence given to both groups, and
 - all comparisons are between groups

(Compare Unit C with current systems to Unit D with future systems)

- While Multiple-Group designs alleviate Order-Effect threats
 - A new set of threats arise
 - ...because now, different treatments are intertwined with different groups
 - ...difficult to separate treatment effects from group effects
 - ...now the differences between capability's might be due to...
 - inherent personnel differences
 - differences in data collection accuracies
 - differences of trial conditions
- ...between groups

3. Isolating the Reason for Change MULTIPLE-GROUP DESIGN UNINTENDED DIFFERENCES

	THREAT	PREVENTION
Unit	<p>15. Player Groups differ in Proficiency</p> <ul style="list-style-type: none"> • Initial group differences <ul style="list-style-type: none"> • nonrandomized assignment • Evolving group differences <ul style="list-style-type: none"> • drop-out differences between groups • Design group differences <ul style="list-style-type: none"> • change after assigning individuals to groups based on past scores • Dominator group differences <ul style="list-style-type: none"> • one individual can influence score • Motivational differences <ul style="list-style-type: none"> • initiation • compensation • resentment 	<ul style="list-style-type: none"> • Use randomization or matching. Report similarities and differences. • Monitor drop outs. • Use no-treatment control groups. • Use statistical methods, analyze data with and without treatment. • Distribute information flow between groups.
Effect	<p>16. Data Collection Accuracy differs for each Player Group Different instrumentation, SMEs, or data collectors</p>	<ul style="list-style-type: none"> • Conduct pretrial and posttrial comparability. • Rotate data collectors between groups.
Trial	<p>17. Player Groups operate under different Trial Conditions Different OPFOR tactics or environmental conditions</p>	<ul style="list-style-type: none"> • Use simultaneous presentation when possible. • Measure trial conditions for comparability.

Applicable to OT?

Multiple-group design validity unintended differences between treatments are controlled



Four Requirements To Design Rigorous Warfighting Experiments

Internal Validity

1. **Capability Used**
2. **Detection of Change in Effect**
3. **Isolation of Reason for Change**

External Validity

4. **Relating Results to Military Operations**

4. Ability to Relate Results to Actual Operations

DEFINITION

- Given that **A** was employed
- Given that **B** changed as **A** was applied
- and **A** alone probably caused change in **B**
- Next Question: **Are these findings related to actual operations?**

Operational Validity:
*Experiment effects can be expected
in actual combat operations.*

Threat - - Amount of change in the outcome measure (**B**) may not occur in actual combat

Realism in conducting experiment is key to eliminating threat

4. Ability to Relate Results to Actual Operations

Experiment Operational Realism Validation similar to M&S validation

Validation of M&S

“...determining the degree to which a model is an accurate representation of the real world...” (DOD VVA Recommended Practice Guide, 1996)

Techniques

Face Validation- experts provide subjective assessments

Operational Validation of Warfighting Experiments

...determining the degree to which an experiment is an accurate representation of the real world.

Techniques

Prototype Validation

Threat Validation

Scenario Validation

Exercise Simulation Accreditation

Predictive Validation

- comparison to training exercise results (UJTL tasks, conditions, standards)
- comparisons to actual operations

Experts provide subjective assessment

OT

OT

4. Ability to Relate Results to Actual Operations

THREAT

PREVENTION

<p>18. Functionality does not represent future capability</p> <p>Treatment Not functionally representative</p>	<ul style="list-style-type: none"> • Ensure functionality of experimental “surrogate” capability is present.
<p>19. Players do not represent operational warfighters</p> <p>Unit</p> <ul style="list-style-type: none"> • Level of training --undertrained or overtrained (golden crew) • Nonrepresentative players 	<ul style="list-style-type: none"> • Use actual end use • Provide sufficient “practice time.” • Use “...ats
<p>20. Measures do not represent operational effect</p> <p>Effect</p> <ul style="list-style-type: none"> • Use of approxi • Inadequate measure • Single data or • Qualitative measures only 	<ul style="list-style-type: none"> • Use simulation to address complex measure based on component measure input (model-test-model). • Use multiple data collectors. • Show correlation to related quantitative measures
<p>21. Unrealistic scenario</p> <p>Trial</p> <ul style="list-style-type: none"> • Blue operations inappropriate • Threat unrealistic • Unrealistic setting • Player familiarity with scenario 	<ul style="list-style-type: none"> • Provide combat developer accreditation • Provide adaptive independent accredited threat • Provide appropriate political and military background • Adaptive “free play” threat enhances scenario setting and uncertainty

Applicable to OT?



Outline

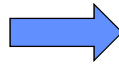


Experiment References

Experiment Logic: 2-3-4-5-21

Experiment Rigor Requirements

- Threats to Rigor
- Good practices to counter threats



Implications of Logic

- Design of individual experiments
- Campaign of experiments

Understanding 4 Experiment Requirement provides insights into Experiment Design TRADEOFFS

All Experiments are tradeoffs: -can not eliminate all threats to validity
The 100% valid Experiment does not exist

4 Requirements
3 Components

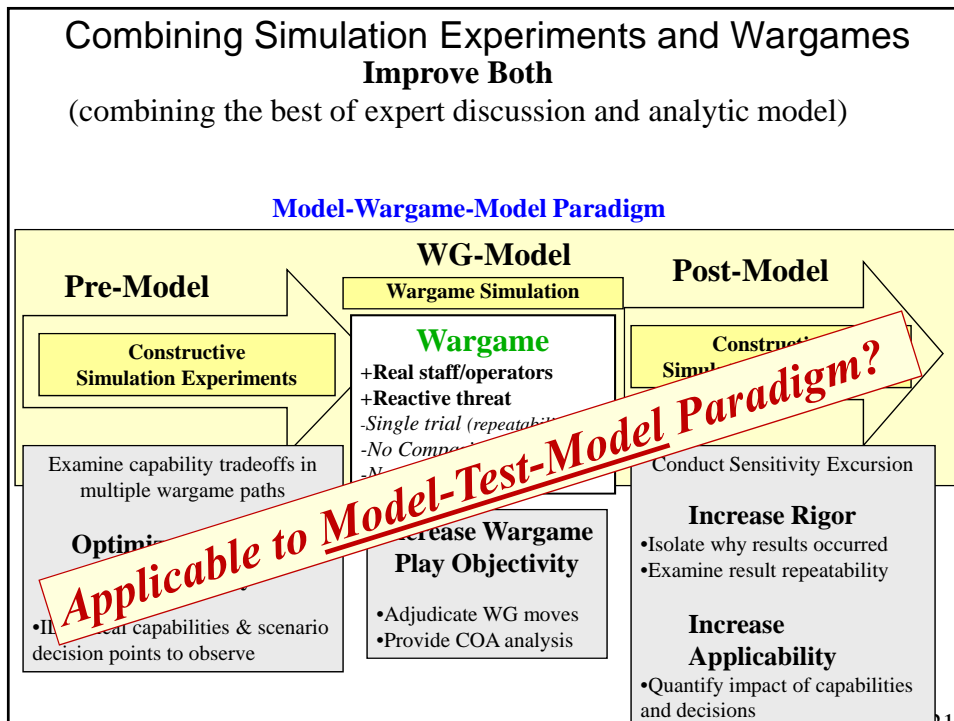
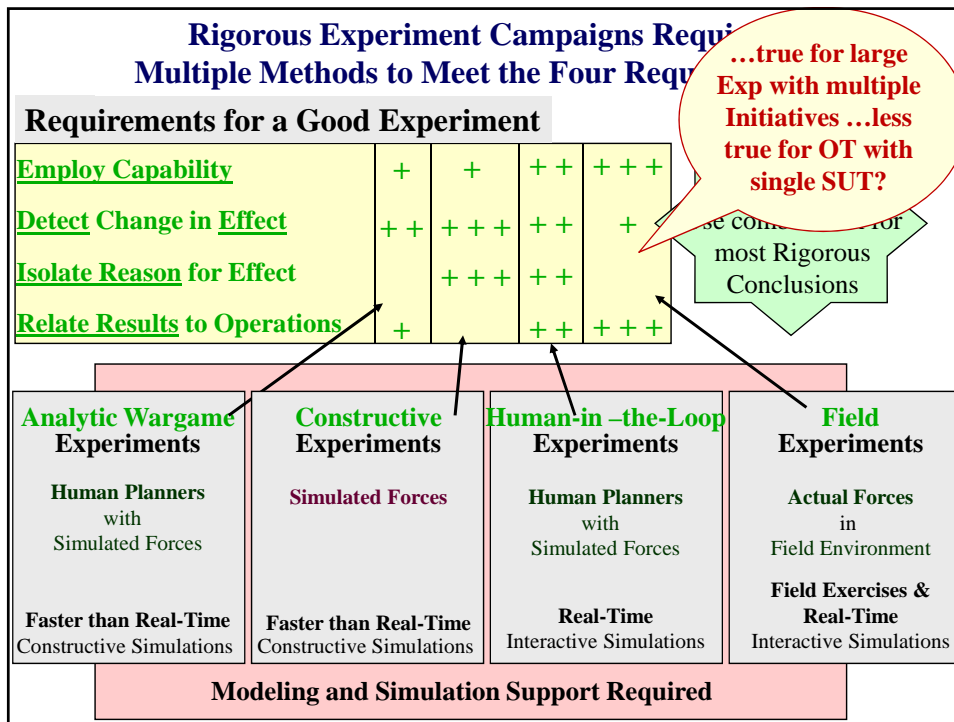
21 Threats to a Good Warfighting Experiment

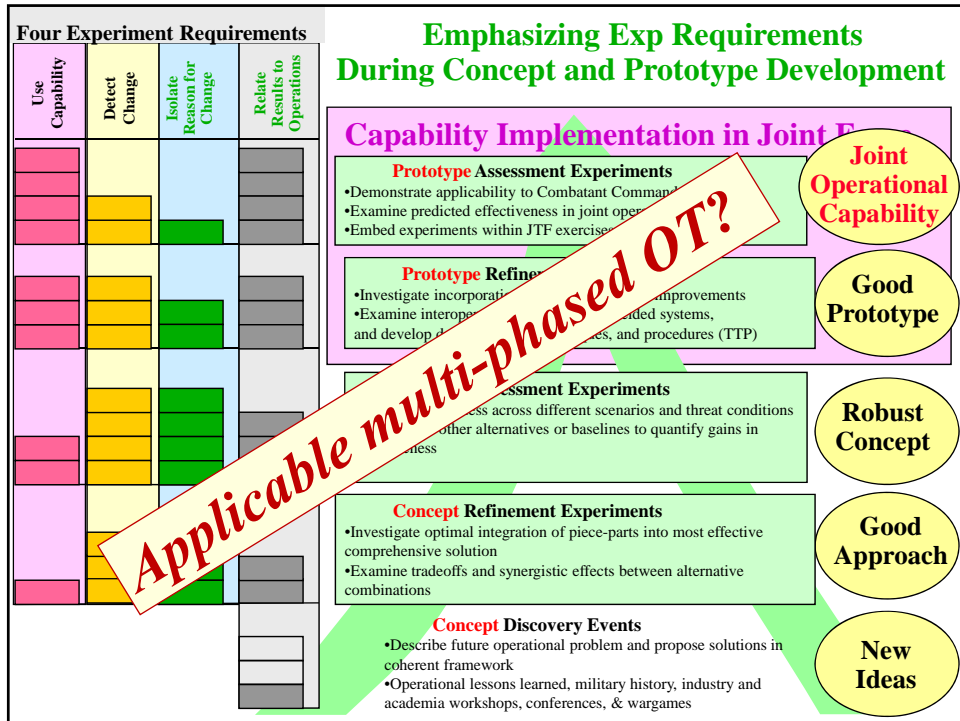
	① Ability to Use Cap Ability	② Ability to Detect Results	③ Ability to Isolate Results for Results (Single Cause)	④ Ability to Relate Results to Hypothesis
1. Requirement	1. Capability not available when needed	2. Capability not available when needed	3. Capability changes over time	4. Measurement capability not available when needed
2. Requirement	5. Paper not available when needed	6. Paper not available when needed	7. Paper not available when needed	8. Paper not available when needed
3. Requirement	9. Data not available when needed	10. Data not available when needed	11. Data not available when needed	12. Data not available when needed
4. Requirement	13. Data not available when needed	14. Data not available when needed	15. Data not available when needed	16. Data not available when needed
5. Requirement	17. Data not available when needed	18. Data not available when needed	19. Data not available when needed	20. Data not available when needed


Is this true for OT also?

and experiment provides sufficient validity to support the pending decision


- A valid experiment is a balance between - -
- **Internal validity** - - precision and control
 - **External validity** - - representativeness and realism
 - Example: increasing repetitions for precision, also increases scenario familiarity thus decreasing realism







Summary



How to Design a Rigorous Experiment

Understand Experiment Logic: “2, 3, 4, 5 and 21”

Focus individual experiments on –

4 OT Validity Requirements

... meeting the **4 experiment requirements**

... by eliminate/controlling the **21 Threats to Validity**

Embed individual experiments within an experiment campaign

Army Test and Evaluation Command

36