

---

# Diagnostic Measurement: Theory, Methods, and Applications

Laine Bradshaw  
James Madison University  
and  
Jonathan Templin  
The University of Georgia

---

NCME 2012: Diagnostic Measurement Workshop

## Workshop Overview

---

- **Section 1: Diagnostic Measurement Introduction (8:00 – 9:30am)**
  - Conceptually, how is “diagnostic” measurement different from more traditional measurement?
  - Why/when would you use diagnostic classification models (DCMs)?

**BREAK: 9:30-9:45**
- **Section 2: Theoretical Framework of DCMs (9:45 – 11:30am)**
  - Latent variables in diagnostic measurement
  - Specification of a general DCM, the log-linear cognitive diagnosis model (LCDM)
- **Questions/Discussion (11:30-11:45am)**

**LUNCH: 11:45-12:45pm**
- **Section 3: DCMs in Practice (12:45 – 2:30 pm)**
  - How-to for estimating DCMs with Mplus
  - Applying DCMs to a test of English language grammar

**BREAK: 2:30-2:45pm**
- **Section 4: Structural Model Specifications(2:45 – 3:45pm)**
  - Specifying structural model
  - Making alterations in Mplus
- **Section 5: Model Fit (3:45 – 4:30pm)**
  - Evaluating model fit for DCMs
- **Questions/Discussion (4:30-4:45pm)**

**EVALUATE TRAINING SESSION: 4:45-5:00pm**

---

# Conceptual Foundations of Diagnostic Measurement

## Session 1

---

NCME 2012: Diagnostic Measurement Workshop

### Session Overview

---

- Key definitions
- Conceptual example
- Example uses of diagnostic models in education
  - Classroom use (formative assessment)
  - Large-scale testing use (summative assessment)
- Why diagnostic models should be used instead of traditional classification methods
- Concluding remarks

---

## Session 1: Conceptual Foundations of Diagnostic Measurement

# DEFINITIONS

---

## What are Diagnoses?

- The word and meaning of diagnosis is common in language
- The roots of the word diagnosis:
  - gnosis: to know
  - dia: from two
- Meaning of diagnoses are deeply ingrained in our society
  - Seldom merits a second thought



# Definitions

---

- *American Heritage Dictionary* definition of *diagnosis*:
  - Generally
    - ♦ (a) A critical analysis of the nature of something
    - ♦ (b) The conclusion reached by such analysis
  - Medicine
    - ♦ (a) The act or process of identifying or determining the nature and cause of a disease or injury through evaluation of a patient's history, examination, and review of laboratory data
    - ♦ (b) The opinion derived from such an evaluation
  - Biology
    - ♦ (a) A brief description of the distinguishing characteristics of an organism, as for taxonomic classification (p. 500)

## Diagnosis: Defined

---

- A diagnosis is the decision that is being made based on information
- Within psychological testing, providing a test score gives the information that is used for a diagnosis
  - BUT, the score is not the diagnosis
  - For this workshop, a diagnosis is by its nature *discrete*
    - ♦ Classification

# Day-to-Day Diagnosis

---

- Decisions happen every day:
  - Decide to wear a coat or bring an umbrella
  - Decide to study
  - Decide what to watch on TV tonight
- In all cases:
  - Information (or data) is collected
  - Inferences are made from data based on what is likely to be the true state of reality

# Diagnosis (Formalized)

---

- In diagnostic measurement, the procedures of diagnosis are formalized:
  - We make a set of observations
    - ◆ Usually through a set of test questions
  - Based on these questions we make a decision as to the underlying state (or states) of a person
    - ◆ The decision is the diagnosis

# Diagnosis (Formalized)

---

- Diagnoses featured in this workshop:
  - Educational Measurement
    - ◆ The competencies (skills) that a person has or has not mastered
      - Leads to possible tailored instruction and remediation
  - Psychiatric Assessment
    - ◆ The DSM criteria that a person meets
      - Leads to a broader diagnosis of a disorder

# Workshop Terminology

---

- **Respondents**: The people from whom behavioral data are collected
  - Behavioral data considered test item responses for workshop
  - Not limited to only item responses
- **Items**: Test items used to classify/diagnose respondents
- **Diagnostic Assessment**: The method used to elicit behavioral data
- **Attributes**: Unobserved dichotomous characteristics underlying the behaviors (i.e., diagnostic status)
  - Latent variables linked to behaviors diagnostic classification models
- **Psychometric Models**: Models used to analyze item response data
  - Diagnostic Classification Models (DCMs) is the name of the models used to obtain classifications/diagnoses

# Diagnostic Classification Model Names

---

- Diagnostic classification models (DCMs) have been called many different things
  - Skills assessment models
  - Cognitive diagnosis models
  - Cognitive psychometric models
  - Latent response models
  - Restricted (constrained) latent class models
  - Multiple classification models
  - Structured located latent class models
  - Structured item response theory

## Psychometric Soapbox

---

- DCMs are but a small set of tools that must be adapted for a common purpose
  - Part of a methodological toolbox that is used to classify respondents
  - Should also include content experts and end-users of the diagnoses
- DCMs link empirical observations and respondents characteristics
  - The models are only as good as underlying theories

## CONCEPTUAL EXAMPLE

## Diagnostic Modeling Concepts

---

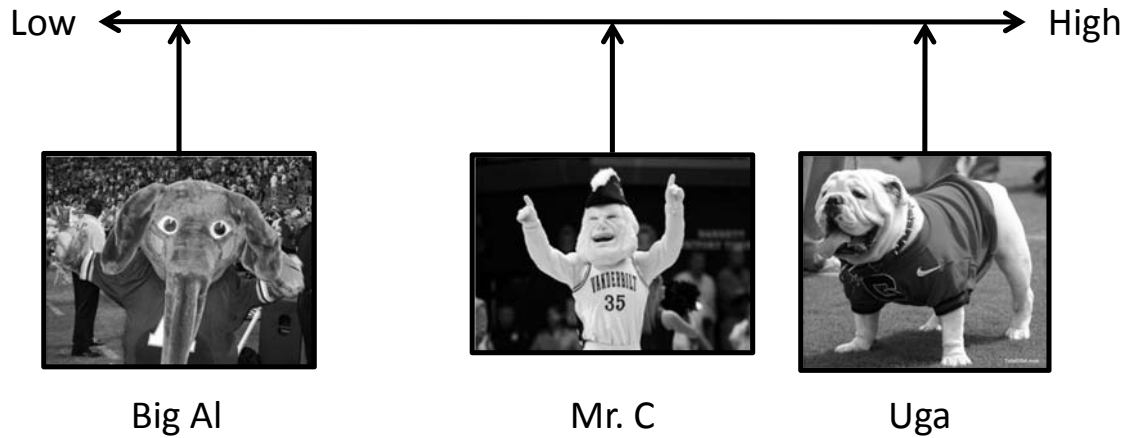
- Imagine that an elementary teacher wants to test basic math ability
- Using traditional psychometric approaches, the teacher could estimate an ability or test score for each respondent
  - Classical Test Theory: Assign respondents a test score
  - Item Response Theory: Assign respondents a latent (scaled) score
- By knowing each respondent's score, the students are ordered along a continuum



# Traditional Psychometrics

---

## Mathematics Ability of SEC Mascots



# Traditional Psychometrics













---

- What results is a (weak) ordering of respondents
  - Ordering is called weak because of error in estimates
  - $Uga > Mr. C > Big Al$
- Questions that traditional psychometrics cannot answer:
  - Why is Big Al so low?
    - ◆ How can we get him some help?
  - How much ability is “enough” to pass?
    - ◆ How much is enough to be proficient?
  - What math skills have the students mastered?

# Multiple Dimensions of Ability

- As an alternative, we could have expressed math ability as a set of basic skills:
  - Addition
  - Subtraction
  - Multiplication
  - Division

## Ability from a Diagnostic Perspective

	<u>Has Mastered</u>	<u>Has Not Mastered</u>
Addition	  	
Subtraction	 	
Multiplication		 
Division		 

# Multiple Dimensions of Ability

---

- The set of skills represent the multiple dimensions of elementary mathematics ability
- Other psychometric approaches have been developed for multiple dimensions
  - Classical Test Theory - Scale Subscores
  - Multidimensional Item Response Theory (MIRT)
- Yet, issues in application have remained:
  - Reliability of estimates is often poor for most practical test lengths
  - Dimensions are often very highly correlated
  - Large samples are needed to calibrate item parameters in MIRT

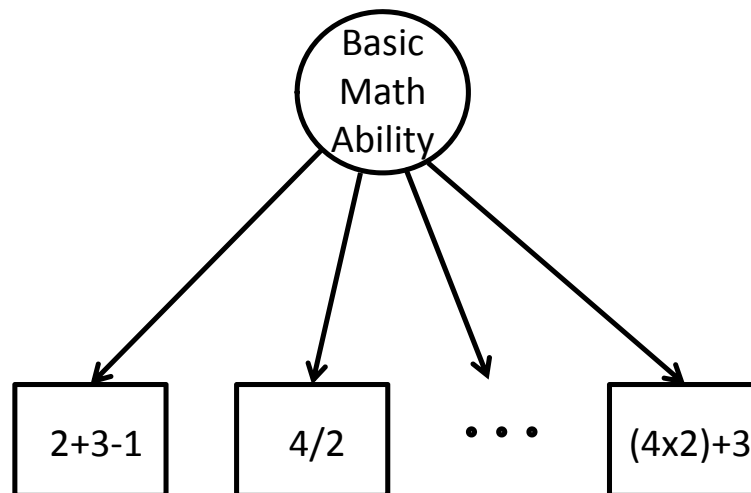
## DCMs as an Alternative

---

- DCMs do not assign a single score
- Instead, a **profile** of **mastered** attributes is given to respondents
  - Multidimensional models
- DCMs provide respondents valuable information with fewer data demands
  - Higher reliability than comparable IRT/MIRT models
  - Complex item structures possible

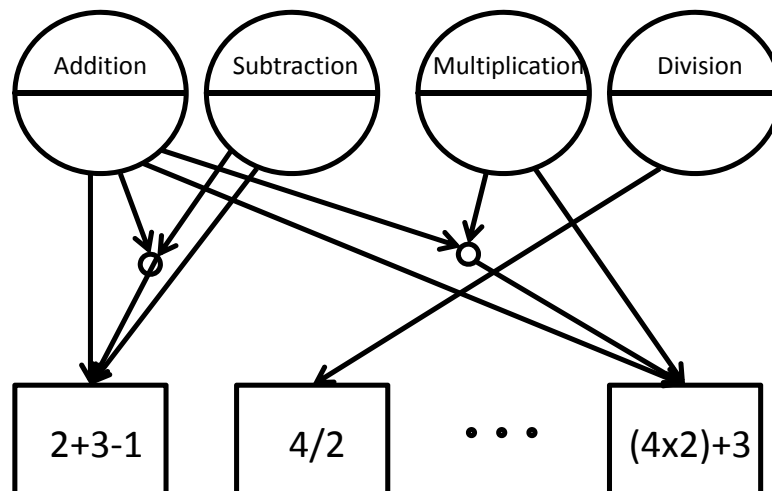
## Path Diagram of Traditional Psychometrics

---



## Path Diagram of Diagnostic Models

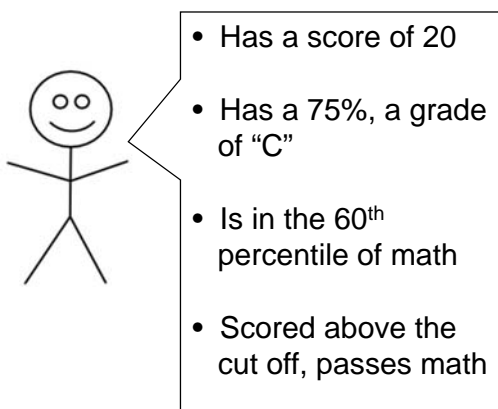
---



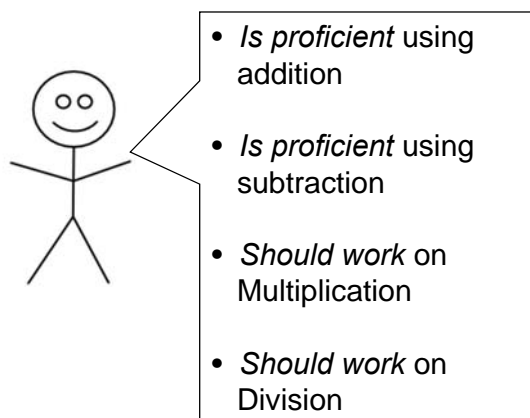
# Psychometric Model Comparison

---

## Using Traditional Models



## Using Diagnostic Models



## DCM Specifics

---

- Let's expand on the idea of the basic math test
- Possible items may be:
  - $2+3-1$
  - $4/2$
  - $(4 \times 2) + 3$
- Not all items measure all attributes
- A Q-matrix is used to indicate the attributes measured by each item
  - This is the ***factor pattern matrix*** that assigns the loadings in ***confirmatory factor analysis***

# The Q-Matrix

- An example of a Q-matrix using our math test

	Add	Sub	Mult	Div
<b>2+3-1</b>	1	1	0	0
<b>4/2</b>	0	0	0	1
<b>(4 x 2)+3</b>	1	0	1	0

## Respondent Profiles

- Respondents are characterized by profiles specifying which attributes have been mastered
  - Numeric values are arbitrary, but for our purposes
    - ♦ Mastery given a 1
    - ♦ Non-mastery given a 0
- For example:

	Add	Sub	Mult	Div
<b>Respondent A</b>	1	1	0	0

- Respondent profile estimates are in the form of ***probabilities of mastery***

# Expected Responses to Items

Q-matrix

	Add	Sub	Mult	Div
<b>2+3-1</b>	1	1	0	0
<b>4/2</b>	0	0	0	1
<b>(4 x 2)+3</b>	1	0	1	0

By knowing which attributes are measured by each item and which attributes have been mastered by each respondent, we can determine the items that will likely be answered correctly by each respondent

Respondent Mastery

	Add	Sub	Mult	Div
<b>Respondent 1</b>	1	1	0	0
<b>Respondent 2</b>	0	1	0	1
<b>Respondent 3</b>	1	0	1	0
<b>Respondent 4</b>	1	1	1	0

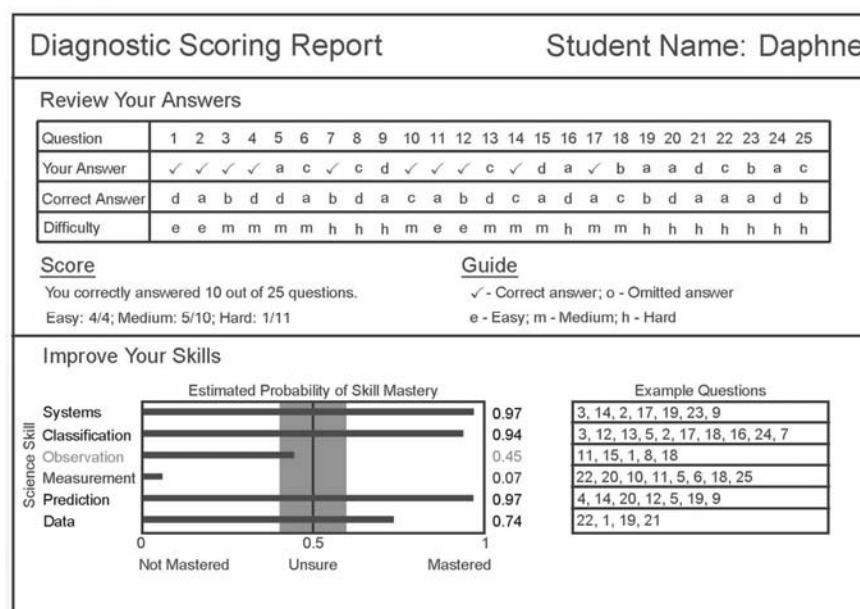
Prob Ans #1

Prob Ans #2

Prob Ans #3

Prob Ans #1 & #3

# DCM Scoring and Score Reporting



from Templin (2007)

# DCM Conceptual Summary

---

- DCMs focus on **WHY** a respondent is not performing well as compared to only focusing on **WHO**
- The models define the chances of a correct response based on the respondent's attribute profile
- Many models have been created ranging in complexity
  - In Session #2 we discuss a general DCM
  - The general model subsumes all other latent-variable DCMs
- The model predicts how respondents will answer each item
  - Also allows for classification/diagnoses based on item responses

## How do DCMs Produce Diagnoses?

---

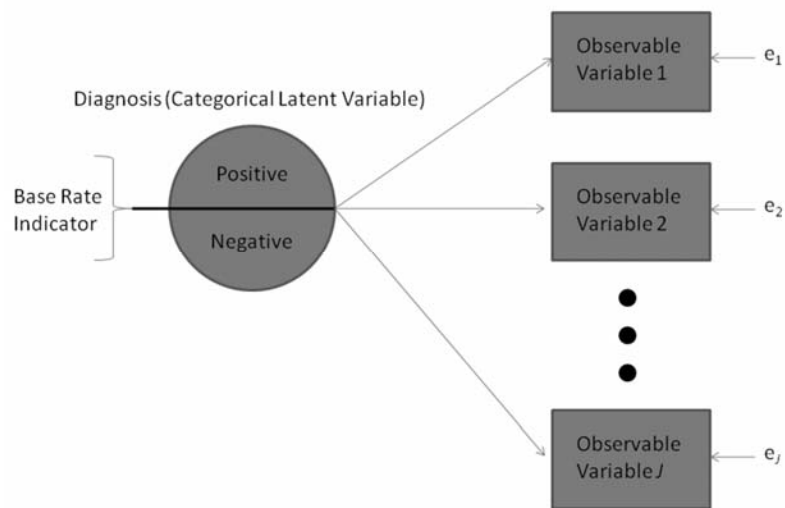
- Diagnostic decisions come from comparing observed behaviors to two parts of the psychometric model:

1.	Item/variable information (item parameters)	Measurement Model
	<ul style="list-style-type: none"><li>♦ How respondents with different diagnostic profiles perform on a set of test items</li><li>♦ Helps determine which items are better at discriminating between respondents with differing diagnostic profiles</li></ul>	
2.	Respondent information pertaining to the base-rate or proportion of respondents with diagnoses in the population	Structural Model
	<ul style="list-style-type: none"><li>♦ Provides frequency of diagnosis (or diagnostic profile)</li><li>♦ Helps validate the plausibility of the observed diagnostic profiles</li></ul>	



# Conceptual Model Mapping in DCMs

---



---

Session 1: Conceptual Foundations of Diagnostic Measurement

## USES OF DIAGNOSTIC MODEL RESPONDENT ESTIMATES

# DCMs In Practice

---

- To demonstrate the potential benefits of using DCMs, we present a brief example of their use
  - From Henson & Templin (2008); Templin & Henson (2008)
- An urban county in a southern state wanted to improve student's End-Of-Course (EOC) scores on the state's 10<sup>th</sup> grade Algebra 2 exam
- A benchmark test was given in the middle of a semester
  - Formative test designed to help teachers focus instruction
- Respondents and their teachers received DCM estimates
  - Used these to characterize student proficiency levels with respect to 5 state-specified goals for Algebra 2 (standards)

## DCM Study

---

- The benchmark test was developed for use with a DCM
  - Characteristics of the test were fixed via standard setting
- Five attributes were measured
  - Mastery was defined as meeting the proficient level for each attribute
  - Attributes were largest represented in EOC exam
- Respondents then took the EOC exam
  - 50 item test:
    - ♦ Score of 33+ considered proficient
  - Benchmark estimates linked to EOC estimates
- Next slides describe how DCMs can help guide instruction

## Descriptive Statistics of Attribute Patterns

- First, the basic descriptive statistics for each possible pattern
- What we expect a respondent with a given attribute pattern to score on the EOC test

Skill Pattern	Expected Score
[00000]	22.9
[00001]	26.0
[00011]	29.3
[00111]	31.4
[01111]	34.8
[11111]	41.9

## Gain by Mastery of Each Attribute

- The difference in test score between masters and non-masters of an attribute can be quantified
- Correlation between attribute and EOC score indicates amount of gain in EOC score by mastery of each attribute

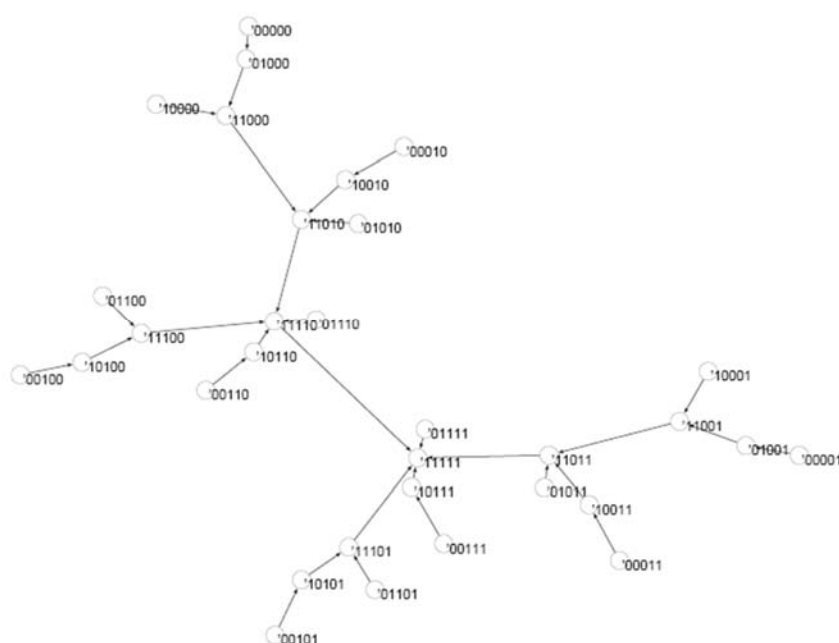
Skill	Gain in Score	Ability Correlation
1	2.61	0.81
2	2.50	0.81
3	1.15	0.63
4	1.19	0.63
5	0.75	0.45

Note: 50 item test

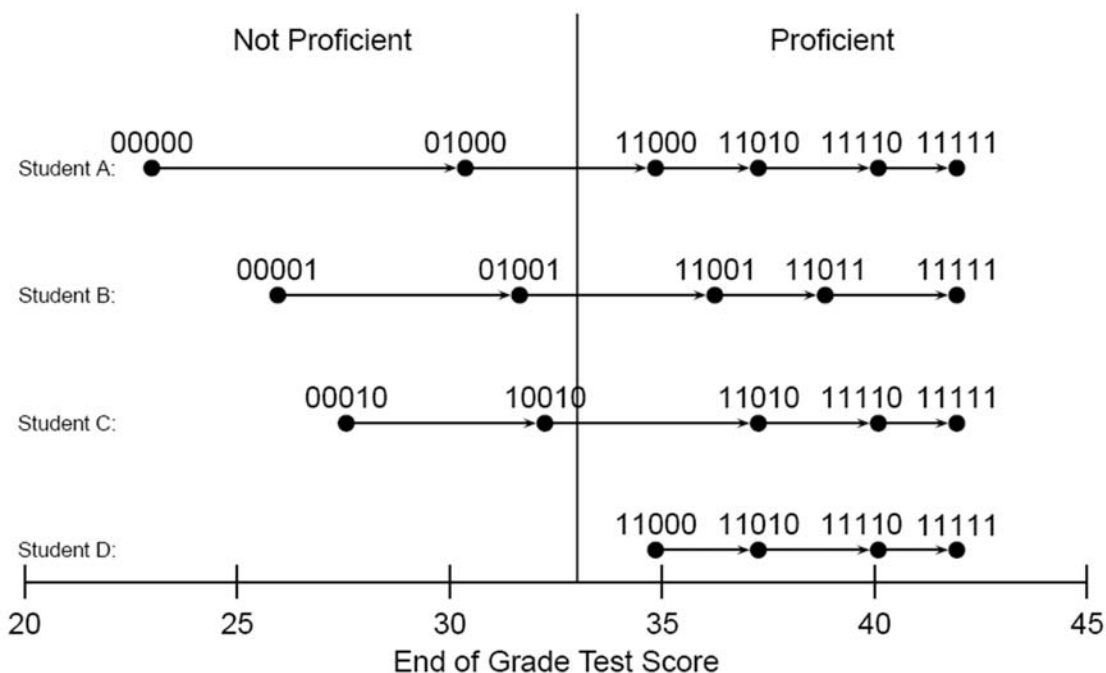
## Pathways to Proficiency

- DCMs can be used to form of a “learning path” a respondent can follow that would most quickly lead to proficiency on the EOC test
- The pathway tells the respondent and the teacher the sequence of attributes to learn next that will provide the biggest increase in test score
- This mechanism may help teachers decide focus on when teaching a course
  - Balances time spent on instruction with impact on test score
- Provides a practical implementation of DCMs in today’s classroom testing environment

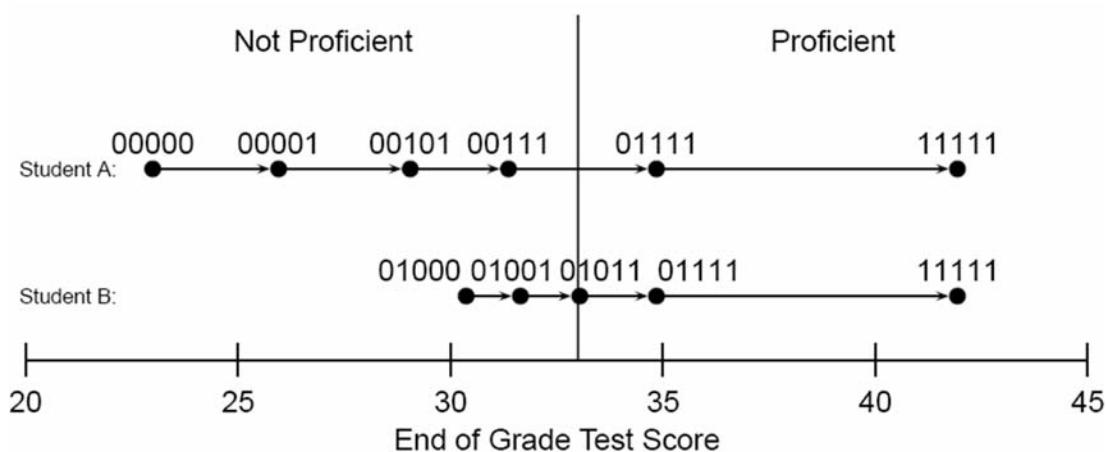
# Proficiency Road Map



## Fast Path to Proficiency



## Harder Paths to Proficiency



- Some paths are less efficient at increasing EOC test scores

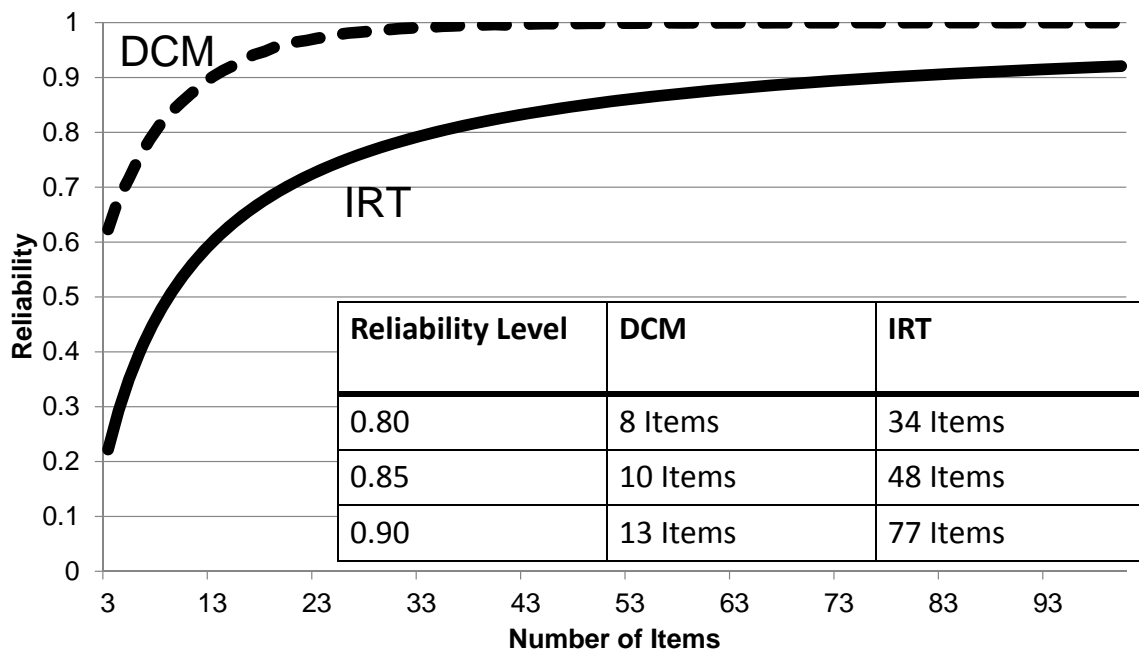
## **IMPLICATIONS FOR LARGE SCALE TESTING PROGRAMS**

### **DCM Characteristics**

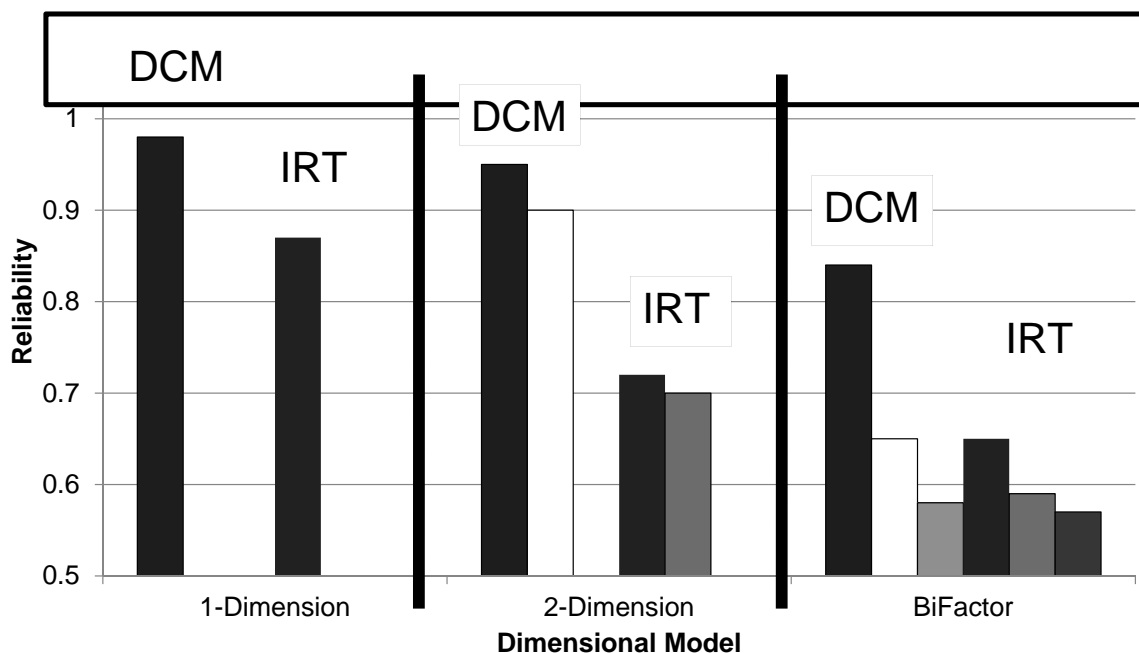
---

- As mentioned previously, DCMs provide a higher level of reliability for their estimates than comparable IRT or CTT models (Templin & Bradshaw, in press)
  - It is easier to place a respondent into one of two groups (mastery or non-mastery) than to locate them on a scale
- Such characteristics allow DCMs to potentially change how large scale testing is conducted
  - Most EOC-type tests are for classification
    - ◆ Proficiency standards
  - DCMs provide direct link to classification
    - ◆ And direct access to standards

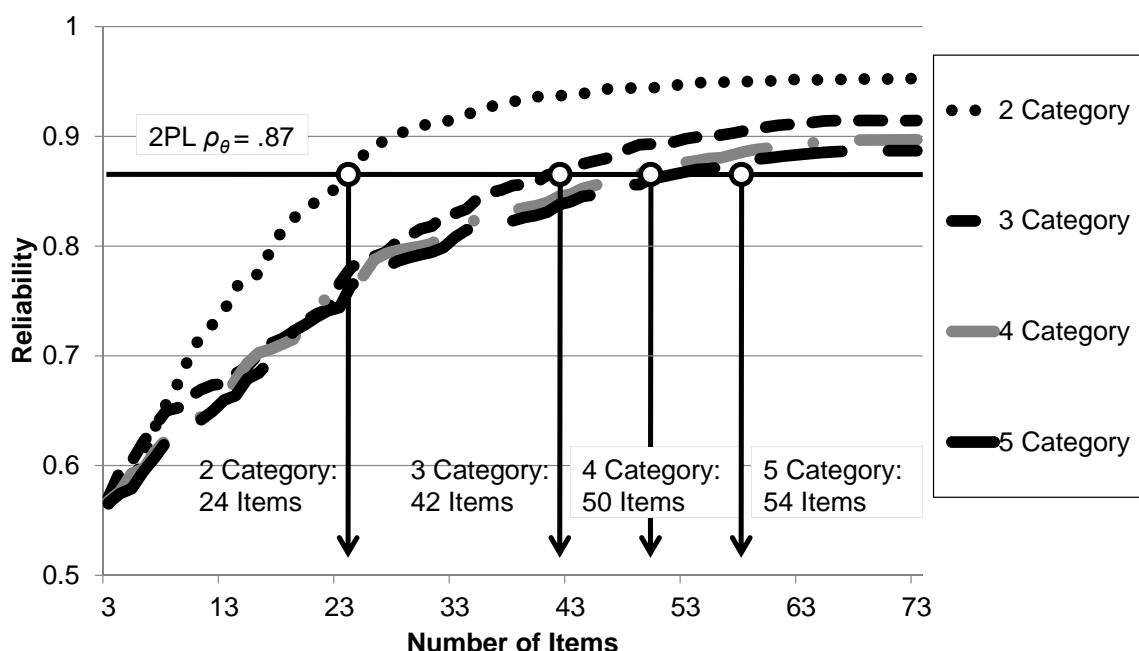
# Theoretical Reliability Comparison



# Uni- and Multidimensional Comparison



## DCMs for an EOC Test



## Ramifications for Use of DCMs

- Reliable measurement of multiple dimensions is possible
  - Two-attribute DCM application to empirical data:
    - ♦ Reliabilities of 0.95 and 0.90 (compared to 0.72 and 0.70 for IRT)
  - Multidimensional proficiency standards
    - ♦ Respondents must demonstrate proficiency on multiple areas to be considered proficient for an overall content domain
  - “Teaching to the test” would therefore represent covering more curricular content to best prepare respondents
- Shorter unidimensional tests
  - Unidimensional DCM application to empirical data:
    - ♦ Test needed only 24 items to have same reliability as IRT with 73 items



# The Paradox of DCMs

---

- DCMs are often pitched as models that allow for measurement of “fine-grained” skills (e.g., Rupp & Templin, 2008)
- Paradox of DCMs:
  - Sacrifice fine-grained measurement of a latent trait for only several categories
  - Increased capacity to measure ability multidimensionally

# When Are DCMs Appropriate?

---

- Which situations lend themselves more naturally to such diagnosis?
  - The *purpose* of the diagnostic assessment matters most
  - DCMs provide classifications directly
    - ♦ Optimally used when tests are used for classification
      - EOC Tests
      - Licensure/certification
      - Clinical screening
      - College entrance
      - Placement tests
  - DCMs *can* be used as coarse approximations to continuous latent variable models
    - ♦ i.e., EOG example (2-5 category levels shown)

---

Session 1: Conceptual Foundations of Diagnostic Measurement

## **BENEFITS OF DCMS OVER TRADITIONAL CLASSIFICATION METHODS**

### **Previous Methods for Classification**

---

- Making diagnoses on the basis of test responses is not a new concept
  - Classical test theory
  - Item response theory
  - Factor analysis
- Process is a two-stage procedure
  1. Scale respondents
  2. Find appropriate cut-scores
- Classify respondents based on cut-scores

## Problems with the Two-Stage Approach

---

- The two-stage procedure allows for multiple sources of error to affect the results
- 1. The latent variable scores themselves: estimation error
  - Uncertainty is typically not accounted for in the subsequent classification of respondents (i.e., standard errors)
  - The classification of respondents at different locations on the score continuum with multiple cut-scores is differentially precise
    - ◆ Uncertainty of the latent variable scores varies as a function of the location of the score

## Problems with the Two-Stage Approach

---

- 2. Latent variable assumptions: that latent variable scores follow a continuous, typically normal, distribution
  - Estimates reflect the assumed distribution
  - Can introduce errors if the assumption is incorrect
- 3. Cut-score determination
  - Standard setting is imprecise when used with general abilities
    - ◆ Standard setting methods can be directed to item performance
  - Some theoretical justification needs to be provided for such a cut-off

# Why are DCMs Better for Classification?

---

- The need for a two-stage procedure to set cut-scores for classification is eliminated when DCMs are used
  - Reduces classification error
- Quantifies and models the measurement error of the observable variables
  - Controlling for measurement error when producing the diagnosis
- DCMs have a natural and direct mechanism for incorporating base-rate information into the analysis
  - No direct way to do so objectively in two-stage procedures
- Item parameters provide information as to the diagnostic quality of each item
  - Not directly estimable in two-stage approaches
  - Can be used to build tests that optimally separate respondents

---

Session 1: Conceptual Foundations of Diagnostic Measurement

## CONCLUDING REMARKS

# Session 1 – Take-home Points

---

- DCMs provide direct link between diagnosis and behavior
  - Provide diagnostic classifications directly
  - Diagnoses set by psychometric model parameters
- DCMs are effective if classification is the ultimate purpose
  - Reduce error by removing judgments necessary in two-stage approach
- DCMs can be used in many contexts
  - Can be used to create highly informative tests
  - Can be used to measure multiple dimensions