# edmentum™

## Designing Score Reports to Maximize Validity and Instructional Utility

*Karen Barton & Audra Kosh*

Validity is about trust and utility.

It is a balance between purpose, defensibility,
and the decisions to be made.

# Validity in the Balance

## Why measure anything?

"We assess for two reasons:

(1) to gather evidence to inform instructional decisions and

(2) to encourage students to try to learn" (Stiggins, 2008, p.3)

## From "purpose" and intention to decisions and consequences:

- What are the instructional decisions to be made?
- Who will be making those decisions?
- What information will help them make good decisions?
- *What are the consequences?*

edmentum™

# Thomas Kuhn's
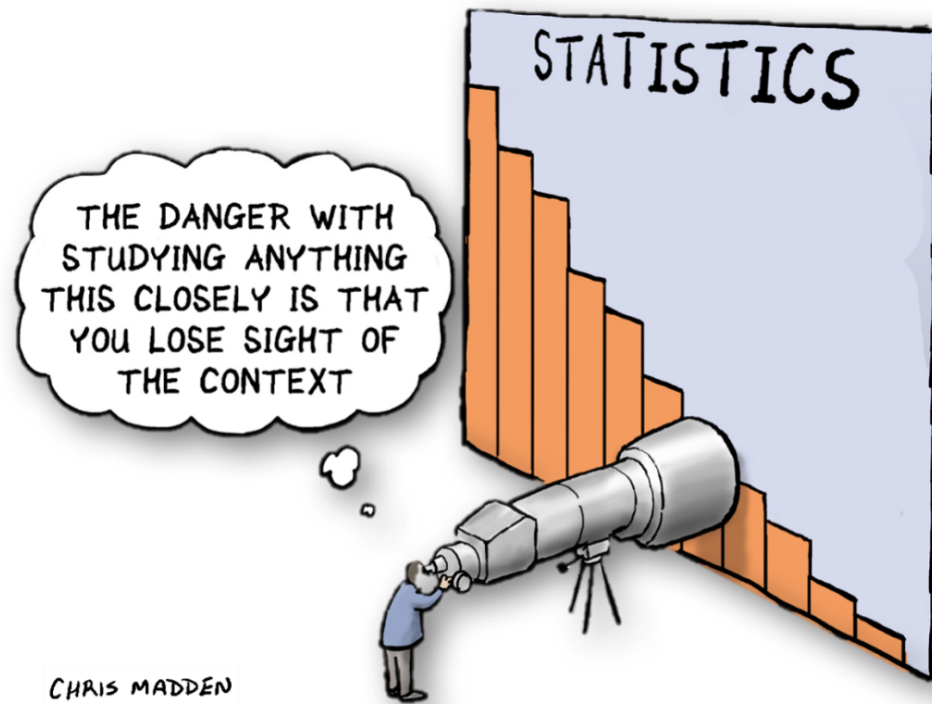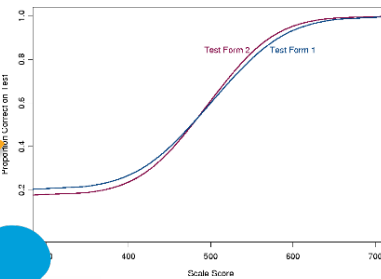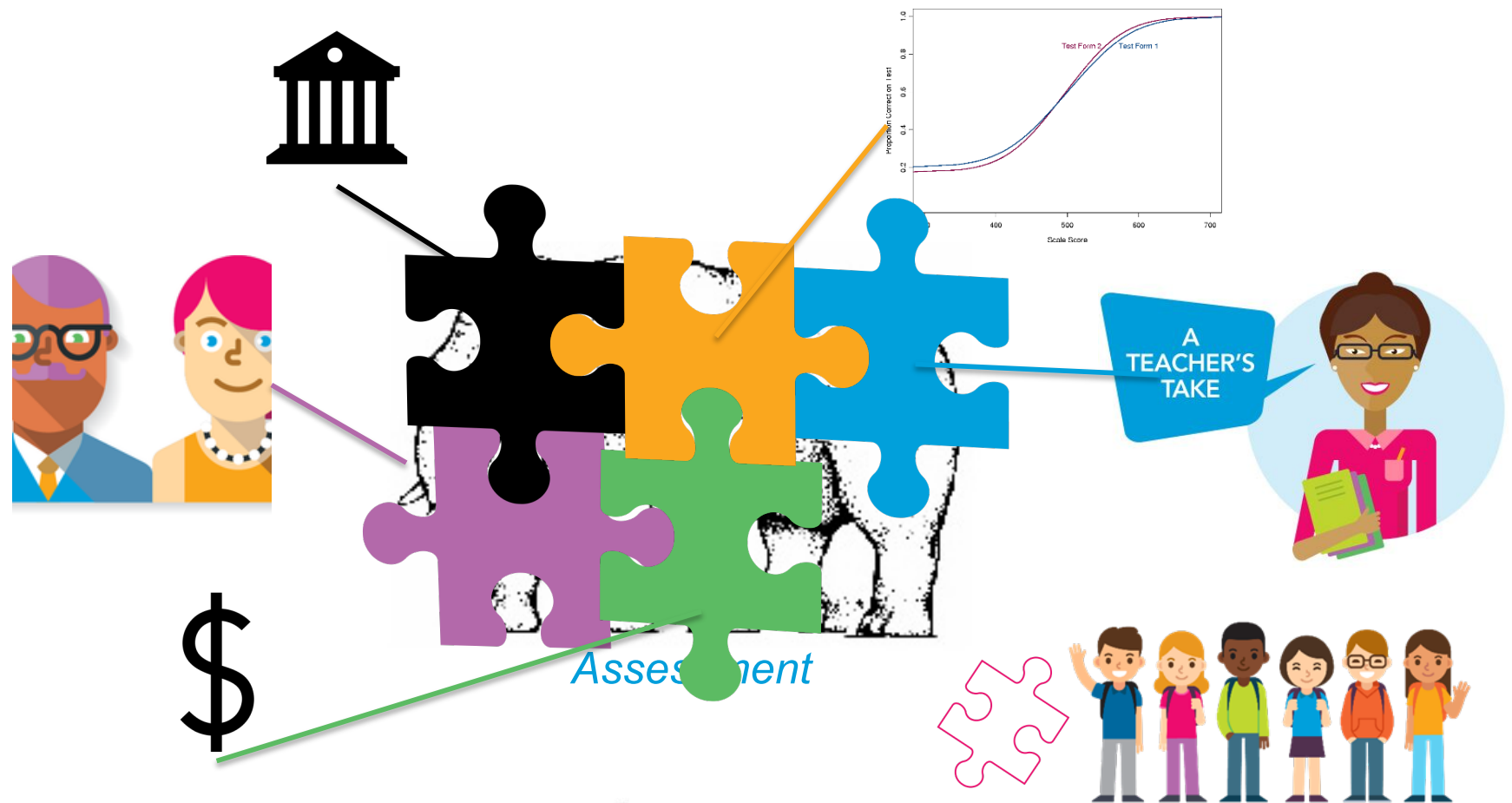# Theory Laden Perspective
# And the Impact on
# 2-Way Communications

edmentum

# How close are we?

# Paradoxical Perspectives
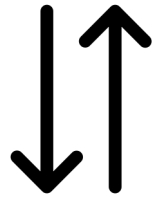


A TEACHER'S TAKE

*Assessment*

edmentum™

# Validity in the Balance

According to the Standards (2014):

**Test score reporting** is a developer responsibility.
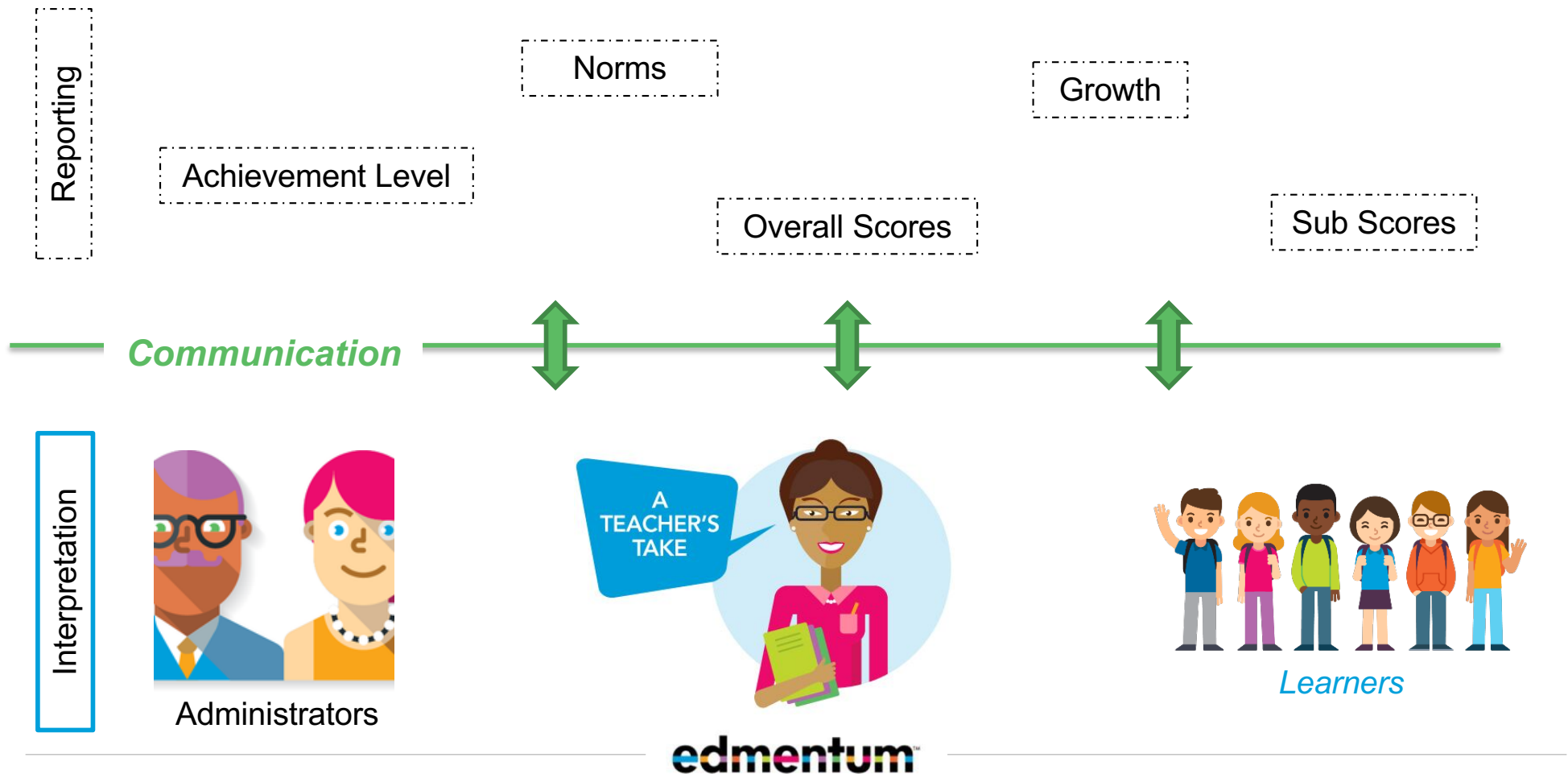
*providing the information*

**Interpretation** is the test user responsibility.

*understanding, communicating and making decisions*

To increase validity of reporting requires attending to information and how it is communicated, as well as greater awareness of context, decisions, and consequences.

edmentum

# Perspectives and Communication
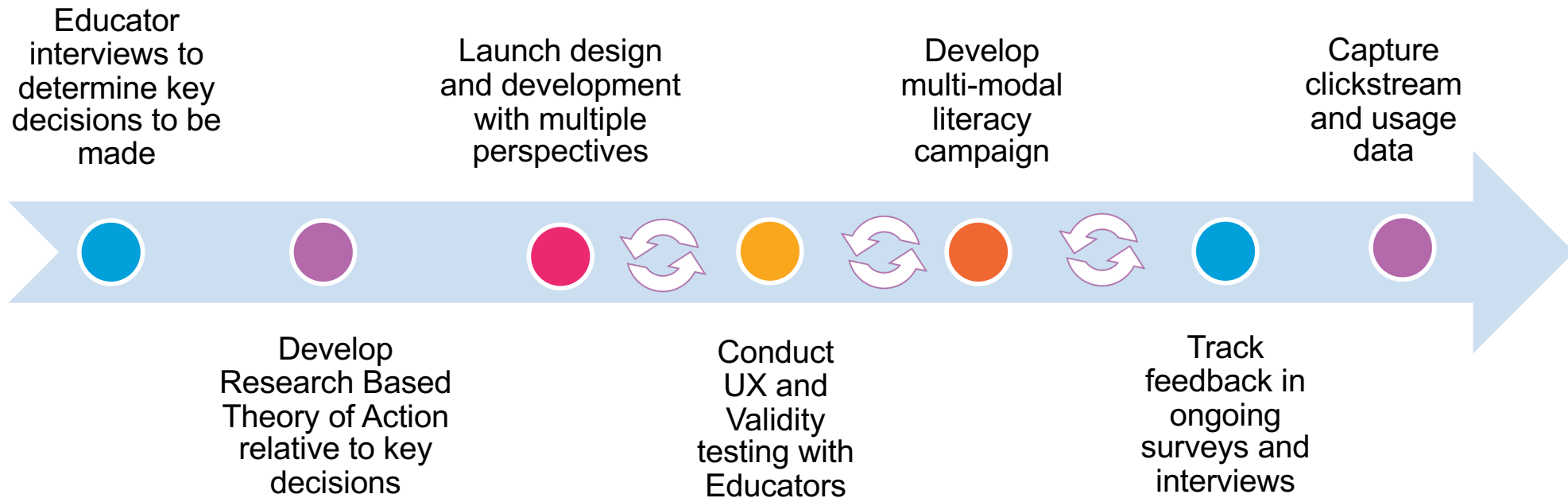
# Educator First Approach

**Prior research and guidance**

- Principled approach (Lewis, 2019)
- Design recommendations (Zenisky & Hambleton, 2012)

**Key elements**

- Validity of the design, not just the information
- Reliability or consistency of the interpretations
- Literacy of the information conveyed
- Transparency
- Ease of use
- Actionable

**edmentum**

# Educator First Workflow Example

Educator interviews to determine key decisions to be made

Launch design and development with multiple perspectives

Develop multi-modal literacy campaign

Capture clickstream and usage data

Develop Research Based Theory of Action relative to key decisions

Conduct UX and Validity testing with Educators

Track feedback in ongoing surveys and interviews

edmentum

# Usability vs Validity Testing

|  | Usability | Validity |
| --- | --- | --- |
| **Driver** | UI/UX | Research |
| **Format** | Focus Group | Individual Educators |
| **Tasks** | Open feedback, prompting questions | Locate information, true/false questions |
| **Interaction** | Highly conversational | Mostly listening |

edmentum™

# Usability Testing

# Applying Feedback

# Applying Feedback

# The Final Design

Mathematics Diagnostic 2 Experience

ⓘ Learn More

Hover over each item from Lucy's diagnostic test to reveal information about the domain, time on item, correct/incorrect response, and skill grade level. Notice how the estimate of Lucy's mathematics ability bounces up and down and the confidence bars tighten as the test narrows in on Lucy's precise mathematics ability.



page_quality score="4"

# In-Product Support

Mathematics Diagnostic 2 Experience

Hover over each item from Lucy's diagnostic test to reveal information about the domain, time on item, correct/incorrect response, and skill grade level. Notice how the estimate of Lucy's mathematics ability bounces up and down and the confidence bars tighten as the test narrows in on Lucy's precise mathematics ability.



### Mathematics Diagnostic Experience - Learn More

**Item:** An individual question.

**Time on item:** The amount of time a student spent on each item. If time on item is only several seconds for many items, this may indicate the score is invalid and that the student was simply clicking through quickly and not trying to answer the questions. Consider talking with the student about their effort on the diagnostic and potentially having the student retake the diagnostic.

📄 Guide to the Student Summary Report ⧉
▶ Video: How does a computer adaptive diagnostic test work? ⧉

Close

edmentum™

# Validity Testing Results

| Correct interpretation with ease | Correct interpretation with struggle | Incorrect interpretation |
|:---:|:---:|:---:|
| ✓ | — | ✗ |

| Report | Concept Assessed | Educator A | Educator B | Educator C | Educator D | Educator E |
|---|---|:---:|:---:|:---:|:---:|:---:|
| Student Report | National Percentile Rank | ✓ | ✗ | ✗ | ✓ | ✓ |
| | CAT Visual right/wrong indicators | ✓ | ✓ | — | ✓ | ✓ |
| | Growth (to get to certain NPR) | — | ✗ | — | ✓ | — |
| | Zoomed-in view of scale | | ✗ | ✓ | ✓ | ✓ |
| | SEM | ✗ | ✗ | ✗ | ✗ | ✓ |
| | Skill-level raw score information | ✓ | | ✓ | ✗ | ✓ |
| | Growth (gain score) | — | ✗ | ✓ | ✗ | ✗ |
| | CAT Visual (domain color coding) | — | ✓ | ✓ | ✓ | |
| | CAT Visual (number of test sessions) | ✓ | ✓ | ✓ | ✓ | ✓ |
| | CAT Visual (adaptive nature) | — | ✗ | ✓ | ✗ | ✓ |
| | CAT Visual (SEM) | ✓ | ✗ | — | ✗ | ✗ |
| Class Report | High/low overall students | ✓ | ✓ | ✗ | ✓ | — |
| | Learning Path Entry Grade by domain | — | ✗ | ✓ | ✓ | ✓ |
| | Growth | ✓ | ✓ | ✓ | — | ✓ |
| | Learning Path Entry Grade Overall | ✓ | ✓ | — | — | — |
| | Scale Score Standard Deviation | — | ✗ | ✗ | — | ✓ |
| | National percentile rank | ✗ | ✗ | ✗ | ✗ | ✓ |

edmentum™

# Educating and Strengthening Receivers

Technical Documents

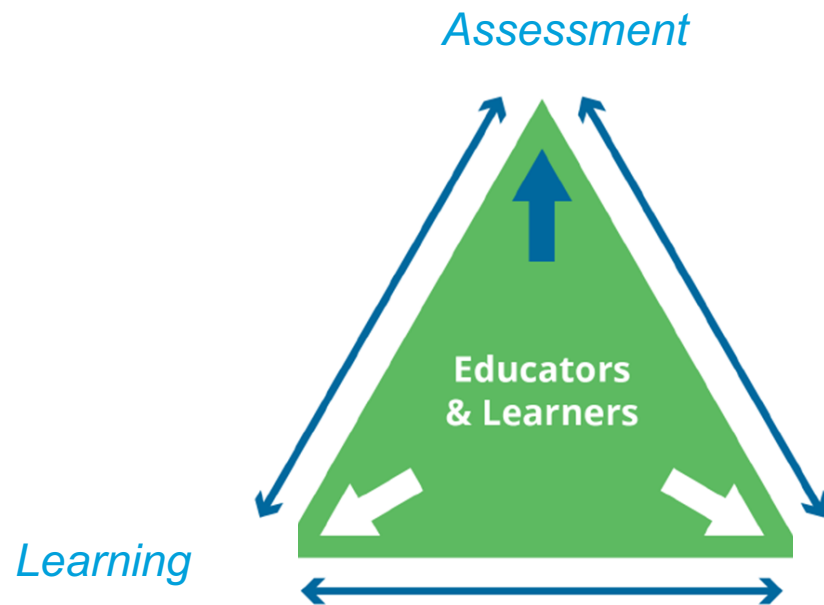Blogs, Marketing, Messaging

Videos

Training Professional Development

In Product Support

edmentum™

# Triangulation and Context in Reporting

Measures of Growth, beyond a single metric.

Assessment

Educators & Learners

Learning

Instruction

edmentum

# Who is the ultimate stakeholder?

How might reports go even further by encouraging student agency and building communications between students and teachers?

edmentum

# Purpose to Impact

- Consider purpose in context

- Establish trust and transparency

- Increasing literacy and impact of actions during testing

- Don't overestimate utility – *ask*

- Don't underestimate responsibility – *go beyond*

edmentum™

*"You can have brilliant ideas*
*(or psychometrics and assessment designs),*
*but if you can't get them across,*
*your ideas won't get you anywhere."*
*~Lee Iacocca*

Be Valid – Be Useful

edmentum™