

## CAEP EVALUATION TOOL FOR EPP-CREATED ASSESSMENTS USED IN ACCREDITATION

**For use with: assessments created by EPPs including observations, projects/ assignments and surveys**  
**For use by: EPPs, CAEP assessment reviewers and Site Visitors**

EXCERPT from the CAEP ACCREDITATION HANDBOOK on “Optional Early Instruments Evaluation”

Early in the accreditation process, providers can elect to submit to CAEP the generic assessments, surveys, and scoring guides that they expect to use to demonstrate that they meet CAEP standards. The purpose of this review is to provide educator preparation providers (EPPs) with formative feedback on how to strengthen assessments, with the ultimate goal of generating better information on its candidates and continuously improving its programs. This feature is a part of CAEP’s specialty/ license area review under Standard 1.

EXAMPLES OF ATTRIBUTES BELOW SUFFICIENT LEVEL	CAEP SUFFICIENT LEVEL	EXAMPLES OF ATTRIBUTES ABOVE SUFFICIENT LEVEL
<p>—</p> <ul style="list-style-type: none"> <li>Use or purpose are ambiguous or vague</li> </ul>	<p><b>1. ADMINISTRATION AND PURPOSE</b> (informs relevancy)</p> <ul style="list-style-type: none"> <li>The point or points when the assessment is administered during the preparation program are explicit</li> <li>The purpose of the assessment and its use in candidate monitoring or decisions on progression are specified and appropriate</li> <li>Evaluation categories or assessment tasks are tagged to CAEP, InTASC or state standards</li> </ul>	<p>+</p> <ul style="list-style-type: none"> <li>Purpose of assessment and use in candidate monitoring or decisions are consequential</li> </ul>
<ul style="list-style-type: none"> <li>Limited or no basis for reviewers to know what information is given to respondents</li> <li>Instructions given to respondents are incomplete or misleading</li> <li>The criterion for success is not provided or is not clear</li> </ul>	<p><b>2. INFORMING CANDIDATES</b> (informs fairness and reliability)</p> <ul style="list-style-type: none"> <li>The candidates who are being assessed are given a description of the assessment’s purpose</li> <li>Instructions provided to candidates about what they are expected to do are informative and unambiguous</li> <li>The basis for judgment (criterion for success, or what is “good enough”) is made explicit for candidates</li> </ul>	<ul style="list-style-type: none"> <li>Candidate progression is monitored and information used for mentoring</li> <li>Candidates are informed how the instrument results are used in reaching conclusions about their status and/or progression</li> </ul>
<ul style="list-style-type: none"> <li>Category or task link with CAEP, InTASC or state standards is not explicit</li> <li>Category or task has only vague relationship with content of the standards being informed</li> <li>Category or task fails to reflect the degree of difficulty described in the standards</li> </ul>	<p><b>3. CONTENT OF ASSESSMENT</b> (informs relevancy)</p> <ul style="list-style-type: none"> <li>Indicators assess explicitly identified aspects of CAEP, InTASC or state standards</li> <li>Evaluation indicators reflect the degree of difficulty or level of effort described in the standards</li> <li>Indicators unambiguously describe the proficiencies to be evaluated</li> <li>When the standards being informed address higher level functioning, the indicators require higher levels of intellectual behavior (e.g., create, evaluate, analyze, &amp; apply). For example, when a standard specifies that candidates’ students “demonstrate” problem solving, then the category or task is specific to students’ application of knowledge to solve problems</li> </ul>	<ul style="list-style-type: none"> <li>Almost all evaluation categories or tasks (at least those comprising 95% of the total score) require observers to judge consequential attributes of candidate proficiencies in the standards</li> </ul>

# CAEP INSTRUMENT RUBRIC

June 2016

EXAMPLES OF ATTRIBUTES BELOW SUFFICIENT LEVEL	CAEP SUFFICIENT LEVEL	EXAMPLES OF ATTRIBUTES ABOVE SUFFICIENT LEVEL
<ul style="list-style-type: none"> <li>• Evaluation categories or tasks not described or ambiguous</li> <li>• Many evaluation categories or tasks (more than 20% of the total score) require judgment of candidate proficiencies that are of limited importance in CAEP, InTASC or state standards</li> </ul>	<ul style="list-style-type: none"> <li>• Most indicators (at least those comprising 80% of the total score) require observers to judge consequential attributes of candidate proficiencies in the standards</li> </ul>	
<ul style="list-style-type: none"> <li>• Rating scales are used in lieu of rubrics; e.g., “level 1= significantly below expectation” . . . “level 4 = significantly above expectation”.</li> <li>• Levels do not represent qualitative differences and provide limited or no feedback to candidates specific to their performance.</li> <li>• Proficiency level attributes are vague or not defined, and may just repeat from the standard or component</li> </ul>	<p>4. <b>SCORING</b> (informs reliability and actionability)</p> <ul style="list-style-type: none"> <li>• The basis for judging candidate work is well defined</li> <li>• Each proficiency level is qualitatively defined by specific criteria aligned with indicators</li> <li>• Proficiency level descriptions represent a developmental sequence from level to level (to provide raters with explicit guidelines for evaluating candidate performance and candidates with explicit feedback on their performance)</li> <li>• Feedback provided to candidates is actionable</li> <li>• Proficiency level attributes are defined in actionable, performance-based, or observable behavior terms. NOTE: If a less actionable term is used such as “engaged”, criteria are provided to define the use of the term in the context of the indicator</li> </ul>	<ul style="list-style-type: none"> <li>• Higher level actions from Bloom’s taxonomy are used such as “analysis” or “evaluation”</li> </ul>
<ul style="list-style-type: none"> <li>• Plan to establish validity does not inform reviewers whether validity is being investigated or how</li> <li>• The instrument was not piloted prior to administration</li> <li>• Validity is determined through an internal review by only one or two stakeholders.</li> <li>• Described steps do not meet accepted research standards for establishing validity.</li> <li>• Plan to establish reliability does not inform reviewers</li> </ul>	<p><b>5.a DATA VALIDITY</b></p> <ul style="list-style-type: none"> <li>• A description or plan is provided that details steps the EPP has taken or is taking to ensure the validity of the assessment and its use</li> <li>• The plan details the types of validity that are under investigation or have been established (e.g., construct, content, concurrent, predictive, etc.) and how they were established</li> <li>• If the assessment is new or revised, a pilot was conducted.</li> <li>• The EPP details its current process or plans for analyzing and interpreting results from the assessment</li> <li>• The described steps generally meet accepted research standards for establishing the validity of data from an assessment</li> </ul> <p><b>5.b DATA RELIABILITY</b></p> <ul style="list-style-type: none"> <li>• A description or plan is provided that details the type of reliability that is being investigated or has been established (e.g., test-retest, parallel forms, inter-rater, internal</li> </ul>	<ul style="list-style-type: none"> <li>• A validity coefficient is reported</li> <li>• types of validity investigated go beyond content validity and move toward predictive validity</li> <li>• A reliability coefficient is reported</li> <li>• Raters are initially, formally calibrated to</li> </ul>

# CAEP INSTRUMENT RUBRIC

June 2016

EXAMPLES OF ATTRIBUTES BELOW SUFFICIENT LEVEL	CAEP SUFFICIENT LEVEL	EXAMPLES OF ATTRIBUTES ABOVE SUFFICIENT LEVEL
<p>whether reliability is being investigated or how.</p> <ul style="list-style-type: none"> <li>• Described steps to not meet accepted research standards for reliability.</li> <li>• No evidence, or limited evidence, is provided that scorers are trained and their inter-rater agreement is documented.</li> </ul>	<p>consistency, etc.) and the steps the EPP took to ensure the reliability of the data from the assessment</p> <ul style="list-style-type: none"> <li>• Training of scorers and checking on inter-rater agreement and reliability are documented</li> <li>• The described steps meet accepted research standards for establishing reliability</li> </ul>	<p>master criteria and are periodically formally checked to maintain calibration at levels meeting accepted research standards</p>
<p><b>WHEN THE INSTRUMENT IS A SURVEY:</b> Use Sections 1 and 2, above, as worded and substitute 6.a and 6.b, below for sections 3, 4 and 5.</p>		
<ul style="list-style-type: none"> <li>• Individual item are ambiguous or include more than one subject</li> <li>• Items are stated as opinions rather than as behaviors or practices</li> <li>• Dispositions surveys provide no explanations of their purpose</li> <li>• Scaled choices are numbers only, without qualitative description linked with the item under investigation</li> <li>• Limited or no feedback provided to candidates</li> <li>• No evidence that questions are piloted</li> </ul>	<p><b>6.a. SURVEY CONTENT</b></p> <ul style="list-style-type: none"> <li>• Questions or topics are explicitly aligned with aspects of the EPP’s mission and also CAEP, InTASC or state standards</li> <li>• Questions have a single subject; language is unambiguous</li> <li>• Leading questions are avoided</li> <li>• Items are stated in terms of behaviors or practices instead of opinions, whenever possible</li> <li>• Surveys of dispositions make clear to candidates how the survey is related to effective teaching</li> </ul> <p><b>6.b DATA QUALITY</b></p> <ul style="list-style-type: none"> <li>• An even number of scaled choices helps prevent neutral (center) responses</li> <li>• Scaled choices are qualitatively defined using specific criteria aligned with key attributes identified in the item</li> <li>• Feedback provided to the EPP is actionable</li> <li>• EPP provides evidence that questions are piloted to determine that candidates interpret them as intended and modifications are made, if called for</li> </ul> <p><b>Criteria listed below are evaluated on site:</b></p> <ul style="list-style-type: none"> <li>• <i>EPP provides evidence that candidate responses are compiled and tabulated accurately</i></li> <li>• <i>Interpretations of survey results are appropriate for the items and resulting data</i></li> <li>• <i>Results from successive administrations are compared (for evidence of reliability)</i></li> </ul>	<ul style="list-style-type: none"> <li>• Scoring is anchored in performance or behavior demonstrably related to teaching practice</li> <li>• Dispositions surveys make an explicit connection to effective teaching</li> <li>• EPP provides evidence of survey construct validity derived from its own or accessed research studies</li> </ul>

# CAEP INSTRUMENT RUBRIC

June 2016

## CHECKLIST

Item Category	Below Adequate	CAEP Adequate Level	Above Adequate	N/A
<p><b>1. ADMINISTRATION AND PURPOSE:</b> Point when instrument is administered in the program; its purpose, and standards addressed (informs relevance). Evaluation categories or assessment tasks are tagged to CAEP, InTASC or state standards.</p>				
<p><b>2. INFORMING RESPONDENTS:</b> Information given to respondent before and at the administration of the instrument (informs fairness and reliability); basis for judging candidate performance is explicit.</p>				
<p><b>3. CONTENT OF ASSESSMENT:</b> evaluation categories explicitly linked with standards, reflect degree of difficulty in standards, and unambiguously describe proficiencies to be evaluated; when standards include higher level functioning, the evaluation categories explicitly require higher levels of intellectual behavior; most evaluation categories require judgment of consequential candidate proficiencies (informs relevancy).</p>				
<p><b>4. SCORING:</b> Basis for judging candidate work is well defined; each proficiency level is qualitatively defined by criteria aligned with the category; proficiency descriptions represent a developmental sequence from level to level and are defined in actionable, performance-based or observable behavior terms; feedback for candidates is actionable (informs reliability and actionability).</p>				
<p><b>5. a. DATA VALIDITY:</b> EPP provides a description or plan that details steps to ensure validity of the assessment and its use; assessment was piloted prior to administration; EPP details process or plans for analyzing and interpreting results.</p> <p><b>b. DATA RELIABILITY:</b> EPP provides a description or plan that details steps to ensure reliability of the assessment; training of scorers and checking inter-rater agreement and reliability are documented.</p>				
<p><b>WHEN THE INSTRUMENT IS A SURVEY:</b> Use sections 1 and 2, above, and substitute 6.a and 6.b for sections 3, 4 and 5.</p>				
<p><b>6. a. SURVEY CONTENT:</b> survey items explicitly aligned with EPP mission and CAEP, InTASC or state standards; questions have a single subject, use unambiguous language; leading questions are avoided; items stated as behaviors or practices rather than opinions.</p> <p><b>b. DATA QUALITY:</b> Even number of scale choices prevents neutral responses; scaled choices are qualitatively defined using criteria aligned with key attributes identified in the item; feedback provided to the EPP is actionable; questions are piloted to ensure intended interpretation; interpretations of results appropriate for items and data.</p>				
<p><b>OVERALL – How would you rate this assessment?</b></p>				
<p><i>Provide a rationale for your overall rating:</i></p>				