

Behavior Screening Tools Chart Rating Rubric

The following rubrics are applied separately for each subscale, grade level/span, and informant targeted by the screening tool.

1. Classification Accuracy

Classification accuracy is rated separately for each criterion measure and time of administration (e.g., fall, winter, spring). Ratings are provided for up to two criterion measures and up to three different time points. Data for additional criterion measures or administration times may be reported but will not be rated. The key questions for classification accuracy are as follows:

- Q1. Was an appropriate criterion measure of social, emotional, or behavioral skills used as an outcome?
- Q2. Was a convincing rationale provided for selecting the comparison point against which the screener was judged (e.g., percentile, cut score)?
- Q3. Were the classification analyses and cut points adequately performed?

Rating	Definition
Full bubble	All of Q1–Q3 rated as yes. and The median estimate of the area under the curve ^a (AUC) is ≥ 0.75 . and Sensitivity ≥ 0.70 and specificity ≥ 0.70 .
Half bubble	All of Q1–Q3 rated as yes. and The median estimate of the AUC is ≥ 0.70 . or Sensitivity ≥ 0.60 and specificity ≥ 0.60 .
Empty bubble	Does not meet full or half bubble.

^a AUC statistic: an overall indication of the diagnostic accuracy of a receiver operating characteristic curve. This curve is a generalization of the set of potential combinations of sensitivity and specificity possible for predictors. AUC values closer to 1 indicate the screening measure reliably distinguishes among students with satisfactory and unsatisfactory reading performance, whereas values at 0.50 indicate the predictor is no better than chance.



2. Reliability

Rating	Definition
Full bubble	At least two types of reliability were reported that are appropriate ^a for the purpose of the tool. and The analyses are drawn from at least two samples that are representative of students across all performance levels. and The median of the estimates for each type met or exceeded 0.70.
Half bubble	At least two types of reliability were reported that are appropriate ^a for the purpose of the tool. and The analyses are drawn from one sample representative of students across all performance levels, and the median of the estimates met or exceeded 0.70. or The analyses are drawn from at least two samples representative of students across all performance levels, and the median of the estimates for each type met or exceeded 0.60.
Empty bubble	Does not meet full or half bubble.
Dash	Reliability data not provided.

^a Tests that require human judgment must report interrater reliability to be eligible for a full or half bubble rating. Other types of reliability must include justification of appropriateness given the purpose of the tool.

3. Validity

Rating	Definition
Full bubble	At least two types of appropriately justified ^a validity analyses are reported. and The analyses are drawn from at least one sample representative of students across all performance levels. and The median of the estimates for each met or exceeded 0.60 (or was within an acceptable range given the expected relationship with the criterion measure(s)).
Half bubble	One type of appropriately justified ^a validity analysis is reported. and The analysis is drawn from a sample representative of students across all performance levels. and The median of the estimates met or exceeded 0.60 (or was within an acceptable range given the expected relationship with the criterion measure(s)).
Empty bubble	Does not meet full or half bubble.



Rating	Definition
Dash	Validity data not provided.

^a Appropriately justified types of validity must include at least one criterion measure that is external to the screening system and theoretically linked to the underlying construct measured by the tool.

4. Sample Representativeness

Description	Definition
National with cross-validation	At least one classification accuracy analysis was conducted using a national sample. ^a and At least one cross-validation study was conducted.
National without cross-validation	At least one classification accuracy analysis was conducted using a national sample. ^a There is no cross-validation study.
Regional with cross-validation	At least one classification accuracy analysis was conducted using one or more state or regional samples. and At least one cross-validation study was conducted.
Regional without cross-validation	At least one classification accuracy analysis was conducted using one or more state or regional samples. There is no cross-validation study.
Local with cross-validation	At least one classification accuracy analysis was conducted using one or more local district samples. and At least one cross-validation study was conducted.
Local without cross-validation	At least one classification accuracy analysis was conducted using one or more local district samples. There is no cross-validation study.

^a A national sample consists of at least 150 students across at least three of nine geographical divisions defined by U.S. Census Bureau: https://www2.census.gov/geo/pdfs/maps-data/maps/reference/us_regdiv.pdf



5. Bias Analysis Conducted

Rating	Definition
Yes	One or more of the following types of analyses were conducted: Multiple-group confirmatory factor models for categorical item responses Explanatory group models, such as multiple-indicators, multiple-causes or explanatory item response theory (IRT) with group predictors Differential item functioning from IRT Testing differential classification accuracy across demographic groups
No	Does not meet "Yes"

This resource was produced under U.S. Department of Education, Office of Special Education Programs, Award No. H326Q210001. Celia Rosenquist serves as the project officer. The views expressed herein do not necessarily represent the positions or policies of the U.S. Department of Education. No official endorsement by the U.S. Department of Education of any product, commodity, service, or enterprise mentioned in this document is intended or should be inferred.

