



unesco

Institute for Statistics
Global Alliance to
Monitor Learning

PROTOCOL FOR REPORTING ON SDG GLOBAL

INDICATOR 4.1.1

March 2022

Acronyms and abbreviations

ACER	Australian Council for Educational Research
CNA	Cross-national assessments
CAT	Content Alignment Tool
EGMA	Early Grade Mathematics Assessment
EGRA	Early Grade Reading Assessment
ERCE	Estudio Regional Comparativo y Explicativo
GAML	Global Alliance for Monitor Learning
GCF	Global Content Framework
IBE	International Bureau of Education
IEA	International Association for the Evaluation of Educational Achievement
IRT	Item response theory
ISCED	International Standard Classification of Education
LLECE	Latin American Laboratory for Assessment of the Quality of Education
MICS	Multiple Indicator Cluster Survey
MILO	Monitoring Impacts on Learning Outcomes
MPL	Minimum Proficiency Level
NAEP	National Assessment of Educational Progress
NLA	National learning assessments
OECD	Organization for Economic Co-operation and Development
PAL	Network People's Action for Learning Network
PASEC	Programme d'analyse des systèmes éducatifs de la confemen
PAT	Procedural Alignment Tool
PILNA	Pacific Island Literacy and Numeracy Assessment
PIRLS	<i>Progress in International Reading Literacy Study</i>
PISA	<i>Programme for International Student Assessment</i>
PISA D	PISA for development
PLD	Proficiency Level Descriptors
SACMEQ	Southern and Eastern Africa Consortium for Monitoring Educational Quality
SDG	Sustainable Development Goal
SEA-PLM	Southeast Asia Primary Learning Metrics
SEAMEO	Southeast Asian Ministers of Education Organization
TIMSS	Trends in International Mathematics and Science Study
UIS	UNESCO Institute for Statistics

Table of Contents

Acronyms and abbreviations	1
Executive summary	4
1. OBJECTIVES AND STRUCTURE	5
2. THE CHALLENGES	5
3. REPORTING CONSISTENCY	6
4. GLOBAL CONTENT FRAMEWORK	7
4.1 <i>Why?</i>	7
4.2 <i>Outputs</i>	7
4.3 <i>Expected Outcome</i>	7
5. PROCEDURAL ALIGNMENT	7
5.1 <i>Why?</i>	7
5.2 <i>How?</i>	7
5.3. <i>Outputs</i>	8
5.4 <i>Expected Outcome</i>	8
6. PROFICIENCY FRAMEWORK AND MINIMUM LEVEL, LINKING STRATEGIES AND INTERIM REPORTING	8
6.1 <i>Why?</i>	8
6.2 <i>How?</i>	8
6.3 <i>Outputs</i>	8
6.4 <i>Expected Outputs and Outcome</i>	9
7. INTERIM REPORTING	9
8. PROTOCOL FOR REPORTING SDG INDICATOR 4.1.1 IN THE UIS DATA RELEASE	10
8.1 <i>Alignment to the global MPL:</i>	11
8.2 <i>Selection of reporting source when various sources are available</i>	11
8.3 <i>Mapping of grade to measurement point:</i>	12
8.4 <i>Mapping of domains to reading or mathematics:</i>	12
8.6 <i>Protocol for reporting with NLA</i>	12
8.7 <i>Parity index indicators (SDG indicator 4.5.1)</i>	13
8.8 <i>Summary of the reporting protocol</i>	13
ANNEX: ALTERNATIVE LINKING METHODOLOGIES	15
Pedagogical calibration or “social moderation”	15
The statistical approach.....	16

List of tables

Table 1. Summary of process and Focus of GAML	6
Table 2. Interim reporting in a nutshell	10
Table 3. 2019 Indicator 4.1.1 reporting.....	14
Table 4. Stock Taking on status and execution	18

List of figures

Figure 1. Indicator 4.1.1	11
Figure 2. Simplified illustration of the pedagogical calibration methodology	15
Figure 3. Simplified illustration of the test-based linking methodology.....	16
Figure 4. Simplified illustration of the item-based linking methodology	17

Executive summary

4.1.1 Proportion of children and young people: (a) in grades 2/3; (b) at the end of primary; and (c) at the end of lower secondary achieving at least a minimum proficiency level in (i) reading and (ii) mathematics, by sex

The reporting format of the indicator aims to communicate two pieces of information:

- I. the percentage of students meeting minimum proficiency standards for the relevant domains (mathematics and reading) for each point of measurement (grades 2/3; end of primary and end of lower secondary); and
- II. when different programs can be considered comparable, and the conditions under which the percentage can be considered comparable to the percentage reported from another country.

This requires the following inputs to frame the indicator:

- What contents should be measured and what is the percentage of coverage to be covered by a given assessment to be comparable to others?
- What procedures are good enough to ensure quality of the data collected? and
- A proficiency scale where all assessments could be informed (and its conversion function or the linking procedure), and a definition of the minimum level for each domain that would allow the estimation of the percentage of students achieving the minimum proficiency level.

An ideal program for reporting on SDG4.1.1 will have gone through three steps: Conceptual Framework, Methodological Framework, and a Reporting Framework. Each of these contains several complex sub-steps. For various levels and types of assessment, UIS had completed most of this work before accepting the responsibility of being custodian of reporting on SDG4.1.1.

Acknowledging that much work had already been done, UIS has prioritized and motivated others to carry out work that had not yet been done. The table below, and this document in general, summarize the work to date. This is represented in the second column of the table. The note discusses the focus of UNESCO Institute for Statistics' (UIS) work (second column), and the columns to the right; and presents a protocol that guides how reporting for Indicator 4.1.1 is now implemented.

1. Objectives and Structure

The document will explain the flow of work, the activities and the outputs in the context of Global Alliance for Monitor Learning's (GAML) broader work program for Indicator 4.1.1. We present them in a logical rather than chronological order.

Each of the activities and outputs help build the tools to generate a minimum level of consistency of education systems' reporting against Indicator 4.1.1, while retaining sufficient flexibility for education systems to pursue assessment programs appropriate to their context and needs.

The reporting format aims to communicate two pieces of information:

- I. the percentage of students meeting minimum proficiency standards for the relevant domains (mathematics and reading) for each point of measurement (grades 2/3; end of primary and end of lower secondary); and
- II. when different programs can be considered comparable and the conditions under which the percentage of programs can be considered comparable to the percentage reported from another country.

This requires the following inputs to frame the indicator:

- What contents should be measured and what is the percentage of coverage covered by a given assessment to be comparable to others? ;
- What procedures are adequate to ensure the quality of the data collected? ; and
- A proficiency scale where all assessments could be informed (and its conversion function or the linking procedure), and a definition of the minimum level for each domain that would allow the estimation of the percentage of students achieving the minimum proficiency level.

Next section defines challenges and Section 3 provides deeper context and sets the logic of workflow. Sections 4, 5 and 6 go deeper in each of the stages of process following same logic and format. Sections 7 and 8 describes the protocol for reporting the indicator to the UIS.

2. The Challenges

The challenges of achieving consistency in global reporting go far beyond the definition of the indicators themselves. In many cases, there is no "one-stop shop" or single source of information for a specific indicator consistent across international contexts. Even when there is agreement on the metric to be used in reporting, a harmonising process may still be necessary to ensure that coverage of the data is consistent.

There are two extremes: at least in theory, the greatest confidence would arise by reporting using a perfectly equated assessment program while, again in theory, the greatest flexibility would arise if reporting could happen with minimal alignment. Both extremes are unsatisfactory for reasons too complex for this document. UIS's approach is a middle one: allow flexibility of reporting, but with growing alignment and comparability over time, without ever necessarily reaching the extreme of a perfectly equivalent assessment or set of assessments. This would allow any assessment program that follows certain comparability guides ahead of time, as well as certain

quality assurance and procedural guides, to report in the relevant domains. This flexible approach implies developing tools to guide countries' work that, if complemented by capacity development activities, will ensure that Indicator 4.1.1 reporting drives knowledge sharing, and growth in global capacity to use assessment programs as levers for system improvement.

3. Reporting Consistency

The objective is to define the criteria and generate the tools that could serve as reference points, transparency tools and normative references.

The tools to be generated have the potential to become a standard against which countries, regions, institutions, international agencies and professionals can benchmark their programs and certificates, and make international comparisons, if they choose to do so. This process already takes place informally in many ways and/or it is now de facto embedded into the various international (and national) assessments.

The workflow is designed following the structure of the implementation of any learning assessment. **Table 1** summarizes the relevant areas of work and contextualizes the work that has taken place and is taking place, with regard to the three main steps in developing a means of reporting on SDG4. This table provides deeper and more detailed context to the introductory materials presented thus far, and highlights the focus of the current work of the UIS and its partners in the last column. It is exactly the same as column 2 in the introductory table above.

Table 1. Summary of process and Focus of GAML

Phase /Tools	What It Addresses	Main Components	Focus of UIS Work
Conceptual Framework	What to assess? - Concept Who to assess? – Population: in and out of school? What contextual information to collect?	<ul style="list-style-type: none"> • Domain and subdomain: Minimum coverage • Target population • Background Questionnaire 	<ul style="list-style-type: none"> • Global Content Framework (GCF) • Content Alignment Tool (CAT)
Methodological Framework	What are the procedures for data integrity	<ul style="list-style-type: none"> • Test design • Sampling frame • Operational design • Data generation • Data analysis 	<ul style="list-style-type: none"> • Good practices guidance • Procedural Alignment Tool
Reporting Framework	What format to report? What is the minimum level? How to link or “harmonize”?	<ul style="list-style-type: none"> • Reporting model • Scale or proficiency framework • Linking • Definition of an interim reporting strategy 	<ul style="list-style-type: none"> • Proficiency Framework and minimum level • Linking strategies • Interim Reporting strategy

Source: UNESCO Institute for Statistics (UIS).

4. Global Content Framework

This section describes in more detail the work that needs to be done, or is being done, for **row 1 of column 3, in Table 1** above.

4.1 Why?

Assessment programs differ in their conceptual frameworks. For example, depending on the curriculum in a country, national assessments usually have different content coverage for a given grade. Domains can be defined differently and, in some cases, programs assess different skills, use different content to assess the same domain, and do both differently, even for the same grade.

To assess the degree of alignment among various assessments and to begin to lay out the basis for a global comparison, UIS and the International Bureau of Education (IBE-UNESCO) collaboratively developed a Global Content Framework (GCF) for the domains of Mathematics and Reading.

4.2 Outputs

There are three final products:

1. Global Content Framework (GCF) of [Mathematics](#) and [Reading](#) to serve as reference;
2. [Content Alignment Tool](#) (CAT) including alignment criteria;
3. A [platform](#) to help countries self-assess.

4.3 Expected Outcome

To ensure data integrity with respect to minimum comparability.

5. Procedural alignment

This section describes in more detail the work that needs to be done, or is being done, for **row 2 of column 3, in Table 1** above.

5.1 Why?

Robust, consistent operations and procedures are an essential part of any large-scale assessment, to maximise data quality and minimise the impact of procedural variation on results. Examples of procedural standards may be found in all large-scale international assessments, and for many large-scale assessments at regional level, where the goal is to establish procedural consistency across international contexts. Many national assessments also set out clear procedural guidelines, to support consistency in their operationalization.

Assessment implementation faces many methodological decisions including test formats and sampling decisions. There is no need for identical procedures and format across assessments. However, there is a need for a minimum set of procedures so data integrity is protected, and results are robust as well as reasonably comparable for any given country over time, but also across countries at any given point in time.

5.2 How?

By defining the minimum procedures that ensure data integrity sufficient enough to report and compare results from different assessment programs.

5.3. **Outputs**

1. [Manual of Good Practices in Learning Assessment](#)
2. [Quick Guide: Making the Case for a Learning Assessment](#)
3. [Quick Guide: Implementing a National Learning Assessment](#)
4. [Procedural Alignment Tool](#)
5. [Online procedural alignment tool platform](#)

5.4 **Expected Outcome**

Comparable data from the point of view of procedures.

6. Proficiency Framework and Minimum Level, Linking Strategies and Interim Reporting

6.1 **Why?**

This section describes in more detail the work that needs to be done, or is being done, for **row 3 of column 3, in Table 1** above.

Assessment programs typically report using different scales. Analysis of results therefore remains contained to one particular test, methodology and scale. While methodologies tend to converge between international and regional assessments, it is still difficult to situate assessments in a common reference continuum of learning outcomes for each level and domain.

The most important issue in the definition of the scales are the proficiency benchmarks or levels embedded within the numerical scale and their cut points on that numerical scale. These benchmarks are typically associated with Proficiency Level Descriptors, which describe in some detail the skills that are typical of students at any given cut point in the scale. Typically, an overarching policy statement or policy definition gives meaning to the succession of cut scores and the proficiency levels but most importantly for defining what constitutes a *minimum* (which is what the SDG4.1.1. indicators call for) proficiency level that has reference to the content.¹

6.2 **How?**

A scale where all the learning assessment programs could be located and the definition of a linking strategy to that scale. The definition of a scale implies:

- A metric that is arbitrary
- The definition of a set of proficiency levels or benchmark including the minimum level
- The policy statements associated to the sets of benchmarks

6.3 **Outputs**

The final products are:

1. [A scale for each domain and point of measurement](#) (benchmarks and definition of the minimum proficiency level or each domain and point of measurement)

¹ Taking from the [National Assessment of Educational Progress \(NAEP\)](#) on policy statement: "Policy definitions are general statements to give meaning to the levels."

2. [A portfolio of linking strategies](#) and the tools that allow to locate assessments proficiency levels in a scale
3. [An Interim reporting](#) on strategy and protocol

6.4 Expected Outputs and Outcome

A proficiency scale that involves the definition of performance levels that are required of students to be proficient, the definition of the number of performance levels, determining the labels and writing descriptions for the levels of the proficiency metric.²

7. Interim Reporting

Definition of activity: To provide a reporting strategy until the content and procedural alignment are finished.

Scope: The UIS has defined an interim reporting strategy that lies within the long-term vision of the UIS reporting strategy.

Currently, the UIS is accepting all national and international assessment data with footnotes and qualifiers to explain where the data come from and to help the users in understanding the limitations of these data.

Table 2. Assessments currently used for reporting

Assessments	Grade 2/3	End of primary	End of lower secondary
ERCE/UNESCO	x	x	
PASEC	x	x	
PILNA	x	x	
PIRLS	x	x	
PISA			x
PISA-D			x
SACMEQ IV		x	
SEA-PLM		x	
TIMSS 4 th grade - Math	x	x	
TIMSS 8 th grade- Math			x
National assessments to be suitable for reporting need either	Subject to Statistical linking Or sound policy linking <u>On an interim basis</u> if a proficiency level has been identified whose descriptor is aligned with the global MPL		
Modules that measure only one Proficiency Level			
MICS6**	x*		
MILO		X	
PAL Network**	X	X	
EGRA.EGMA**	Only if the level aligned is administered in the field and subject to representativeness		

² The initial development of the reporting proficiency scale would draw from both experts' opinion and analysis of existing data and policy level descriptors.

* interim and subject to alignment with the MPL and assessment of all assessment characteristics

** Population Based

Table 2.a. Mapping considerations

Minimum Proficiency Level	As defined by the global minimum proficiency level
Assessments mapped to grade 2 or 3	Assessments administered in grade 2 or 3 according to ISCED level in the country plus one year when primary lasts more than 4 years, except for TIMSS/PIRLS grade 4, which are mapped at the end of primary when primary lasts six years or less.
Assessments mapped at the end of Primary	Last grade of Primary according to ISCED level in the country plus or minus one year from the last year, except for TIMSS/PIRLS grade 4, which are mapped at the end of primary when primary lasts six years or less.
Assessments mapped at the end of Lower Secondary	Last year/grade of Lower Secondary according to ISCED level in the country Plus two minus one grade with respect to the last grade of Lower Secondary

Note: TIMSS/PIRLS Grade 4: these results are allocated to the end of primary when, according to the International Standard Classification of Education (ISCED) levels in a given country, there are 6 or less grades in primary. When primary has more than 6 grades, they are allocated to grade 2/3. * only if they could report to the requested level.

EGMA: Early Grade Mathematics Assessment

EGRA: Early Grade Reading Assessment

LLECE: Latin American Laboratory for Assessment of the Quality of Education

MILO: Monitoring Impacts on Learning Outcomes

MICS 6: Multiple Indicator Cluster Surveys PAL Network: People's Action for Learning- Network

PASEC: Programme d'analyse des systèmes éducatifs de la CONFEMEN (Programme of Analysis of Education Systems of CONFEMEN)

PILNA: Pacific Islands Literacy and Numeracy Assessment

PIRLS: Progress in International Reading Literacy Study

PISA: Programme for International Student Assessment

PISA D: Programme for International Student Assessment for Development

SACMEQ: Southern and Eastern Africa Consortium for Monitoring Educational Quality

SEA-PLM: Southeast Asia Primary Learning Metrics

TIMSS: Trends in International Mathematics and Science Study

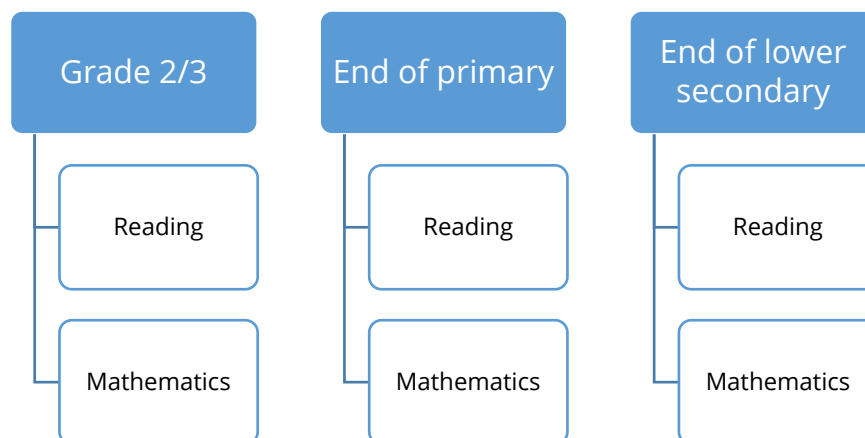
Source: UNESCO Institute for Statistics (UIS).

8. Protocol for reporting SDG indicator 4.1.1 in the UIS data release

Indicator 4.1.1 has three points of measurement and two domains resulting in 6 sub-indicators as reflected in the figure below.

4.1.1 Indicator comprises the following 6 indicators:

Figure 1. Indicator 4.1.1



Source: UNESCO Institute for Statistics (UIS)

8.1 Alignment to the global MPL:

1. For each assessment, programme alignment should be completed using:
 - a. CNA: according to [agreed alignment from the consensus meeting and its revision](#)
 - b. National learning assessment (NLA): provided country does not have any CNA. Use statistical linking or policy linking.

8.2 Selection of reporting source when various sources are available

For each of the indicators listed above for global reporting, the sources of data selected should be prioritized according to the following order of assessments, providing that a mapping of grade has guided a first selection of sources:

- i. International assessments
- ii. Regional assessments
- iii. National assessments
- iv. Population-based assessments

The assessment which maps the best to the grade will always be the preferred choice. Until the equating is completed a cross national assessment within the window of reporting would be the second preferred choice as it allows the comparison with other countries. The national assessment will be chosen only if no alternative program for reporting is available within the window of time required.

Here is a concrete example to illustrate the selection process: Honduras has learning assessments data for Indicator SDG4.1.1a in different years- Reading. These data have been generated from a national, regional and international learning assessments programs in different years:

Year	Source of data
2011	PIRLS 2011 (Grade 4)
2013	ERCE 2013 (Grade 3)
2016	National Learning Assessment (NLA) (Grade 3)

The assessment program to use for reporting on SDG 4.1.1.a - Reading for Honduras will be ERCE 2013. The alternatives PIRLS targets grade 4. Grade 4 is one grade above the intended grade of the indicator –grade 2/3. Thus, ERCE and NLA both assessing grade 3 students, are better option than PIRLS. However, ERCE prevails over NLA, according to the priority of sources of information indicated above as it facilitates the comparison between different countries under the same tool (ERCE).

Until the process of equating international, regional and national assessments will be produced, it is important to use only one source of information so that progress can be tracked on a comparable basis. The restriction is one level and domain. The country could have different reporting sources for a given year for different levels (for instance ERCE for 2/3 grade and PISA/TIMSS for end of lower secondary).

8.3 Mapping of grade to measurement point:

1. Grade 2 or 3: Plus one year when primary lasts more than 4 years according to ISCED levels in the country, except for TIMSS/PIRLS grade 4, which are mapped to the end of primary education when primary education lasts six years or less.
2. End of primary education: Plus or minus one year from the last year of primary education, except for TIMSS/PIRLS grade 4, which are mapped to the end of primary education when primary education lasts six years or less, according to the ISCED level mapping in the country .
3. End of lower secondary education: Plus two or minus one of last year of lower secondary according to ISCED level mapping in the country.

8.4 Mapping of domains to reading or mathematics:

1. Reading:
 - a. The country has no assessment in reading alternative domains, like language or writing, should be used if reading is not available for reporting;
 - b. Different languages could be used for reporting? when results are available in different languages, the official or most relevant language in the country should be used.
2. Mathematics: alternative domains can also be considered and used for reporting.

8.6 Protocol for reporting with NLA

To report on the Indicator for SDG 4.1.1 using NLA, countries must ensure compliance with the following principles:

- a. NLA uses Item response theory (IRT)

- b. Results are available as the percentage of students by proficiency levels
- c. A proficiency level descriptor is aligned to the [global minimum proficiency level](#)
- d. The content of NLA sufficiently covers the Global content framework in the relevant domain, either reading and/or mathematics. The content alignment tool allows this mapping and the guidance could be found in the ([CAT](#) report)
- e. NLA complies with the minimum level of good practices. As discussed in previous section the procedural alignment tool allows countries to assess the level of compliance with a set of good practices and, most importantly, if they reach a minimum level. A report is produced at the end of the process. ([PAT](#) report)
- f. Footnotes are added to the data points (e.g. name of the national assessment, minimum proficiency level, and grade)

8.7 Parity index indicators (SDG indicator 4.5.1)

Parity index indicators are calculated using the [adjusted parity index](#) formula.

8.8 Summary of the reporting protocol

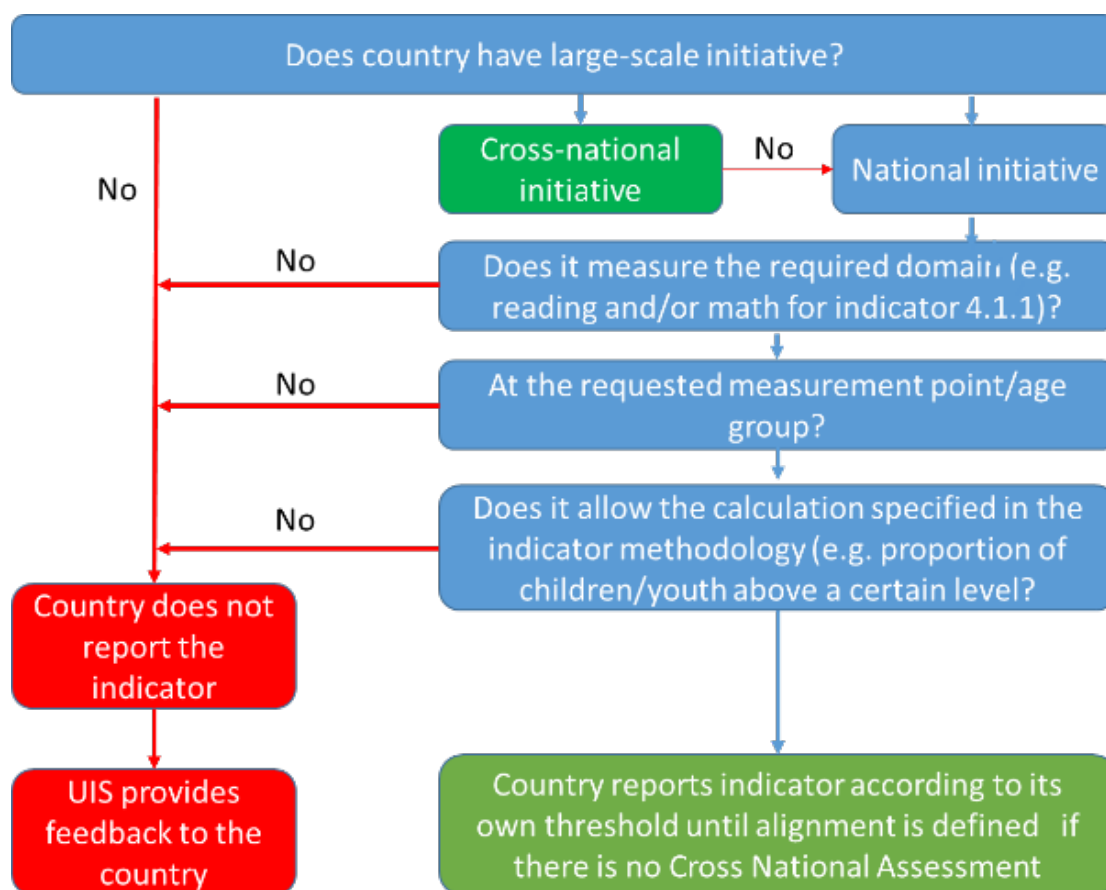


Table 3. 2019 Indicator 4.1.1 reporting

Assessment	Sources of data	Alignment to the Global MPL	Priority for reporting by education level	Footnotes
International assessments: PISA, PISA4D, TIMSS/PIRLS	Yes	According to Consensus Meeting	1	[Assessment name and year]
Regional assessments: LLECE, PASEC, SACMEQ, PILNA	Yes	According to Consensus Meeting	2	[Assessment name and year]
National learning assessments (NLA)	Yes	Suggest and validate alignment	3	“National Learning Assessment (NLA):” [assessment name] & “; Grade ” [grade number] & “; Minimum proficiency level: “ [name or number of the MPL used for reporting]
EGRA/EGMA	Yes	According to Consensus Meeting	4	[Assessment name and year]
MICS	Yes	According to Consensus Meeting	4	[Assessment name and year]
PAL NETWORK (e.g. ASER, UWESO, etc.)	Yes	According to Consensus Meeting	4	[Assessment name and year]

Source: UNESCO Institute for Statistics (UIS)

Annex: Alternative linking methodologies

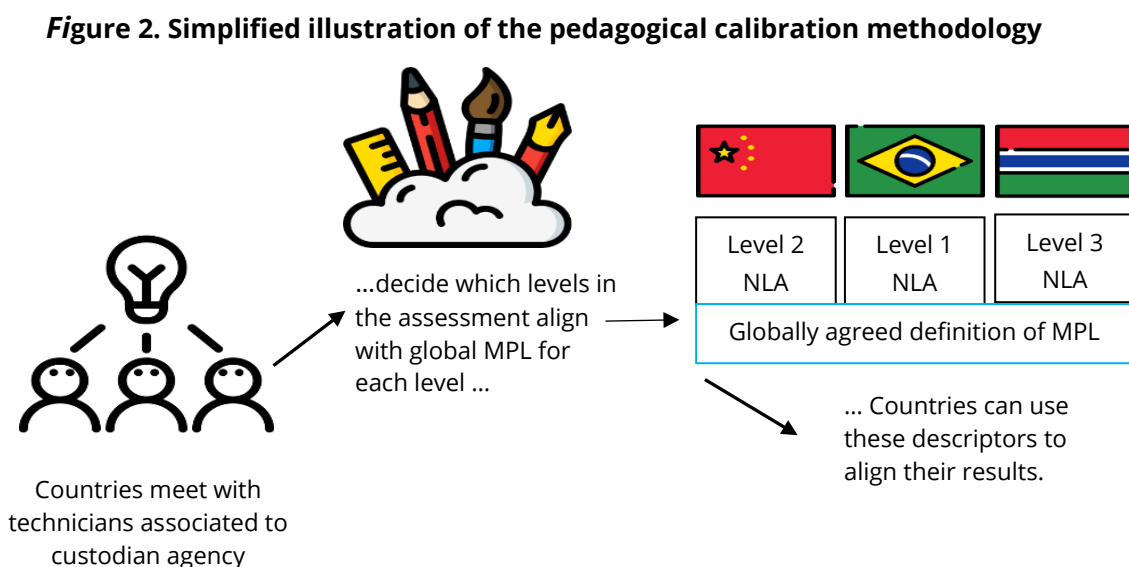
Pedagogical calibration or “social moderation”

What? The approach involves using the definitions of the Proficiency Level Descriptors (PLDs) embedded in a given assessment, which are based on their own proficiency framework and items, to identify the PLD aligned with the MPL global definition and sample items, sourced from various assessments that have been judged by experts to be on either side of that global MPL. The alignment of PLDs (the one of the alternative assessment to be linked to the global MPL) to the global MPL is guided by experts’ judgement, but using a rigorous convergence process amongst the experts.

To increase accuracy, we propose that the majority of the items and the richness of the descriptors be around the global MPL. Thus in a sense most of the analytical and empirical “supports” would be around only one break point (two PLDs), but more PLDs could be used, in order to provide context and increase utility to the countries (especially if some countries would bottom out or top out if only one break point is used).

How? Operationally, for each assessment (national or international) a group of 8 to 10 subject matter experts will convene for a socially moderated workshop. During the workshop, the experts provide individual and independent judgements about each item on the test to define alignment and set initial cut scores based on their understanding of the proficiency level descriptors and experience with the student populations and how they aligned and iterate until a desirable level of convergence is reached.

Which levels of education? Pedagogical calibration could be applied to all levels of education but given the lower feasibility of using a pure psychometric linking method in early grades and the absence of international and regional assessments at the end of primary education in many regions make this methodology the preferred method for alignment in both data points in primary education.



Note: Country flags are used for illustrative purposes only and do not necessarily implicate the participation of these countries in the project.

Source: UNESCO Institute for Statistics (UIS).

The statistical approach

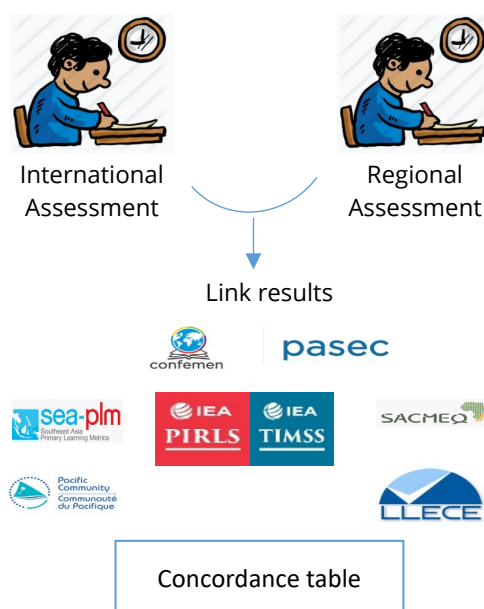
A. Recalibration by running a parallel test on a representative sample of students

What? The test-based linking approach relies on the participation of countries in cross-national assessments (both at the global and regional levels). The methodology allows two assessments, one international and the other regional, to be expressed on the same scale (and that scale to be aligned to the global one, or becoming, de facto, the key to the global one). The alignment will be possible through the creation of concordance tables which will provide a translation between from the countries' regional assessment results to the TIMSS and PIRLS achievement scales.

How? Students will take two assessments, an international and a regional one, and their results on both tests will be aligned in such a way that a link is established between the results from regional assessments conducted at the primary level and the TIMSS and PIRLS International Benchmarks for numeracy and literacy.

Which levels of education? It could be done at the end of primary education where regional and global assessments are administered and global capacity and interest exists. It would take advantage of the proficiency levels set forth by the International Association for the Evaluation of Educational Achievement (IEA) for its two assessments: Trends in International Mathematics and Science Study (TIMSS) and Progress in International Reading Literacy Study (PIRLS). TIMSS has been measuring trends in mathematics and science at four-year intervals since 1995. PIRLS has measured trends in reading literacy at five-year intervals since 2001. With 50 to 70 countries participating in each assessment cycle, the TIMSS and PIRLS achievement scales and their International Benchmarks are well established and used by countries all around the world.

Figure 3. Simplified illustration of the test-based linking methodology



Source: UNESCO Institute for Statistics (UIS)

There are five regional assessments conducting reading and mathematics assessments, and these represent the potential of the application of the methodology:

- SACMEQ – Southern and Eastern Consortium for Monitoring Educational Quality
- PASEC – Programme for the Analysis of Educational Systems
- LLECE – Latin American Laboratory for the Assessment for the Quality of Education
- SEA-PLM – Southeast Asia Primary Learning Metrics
- PILNA – Pacific Island Literacy and Numeracy Assessment
- MILO - Monitoring Impacts on Learning Outcomes

However, for the first application of the methodology, UIS and IEA will be implementing it in the Latin American region through a collaboration with LLECE. Because the concordance tables provide a projected TIMSS or PIRLS score for all possible regional assessment scores, it will be possible to determine the regional assessment scores equivalent to each of the TIMSS and PIRLS International Benchmarks.

B. Psychometrically informed recalibration based on common items

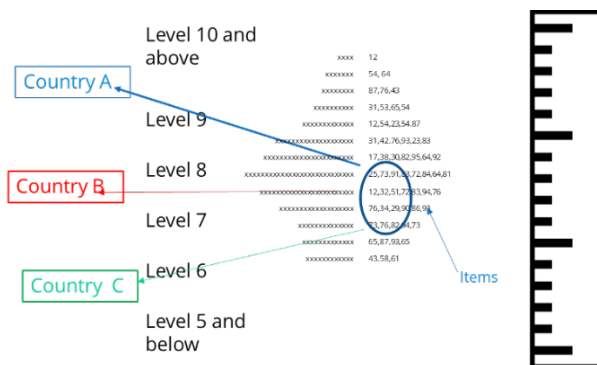
What? Item-based linking is a statistical method that relies on common assessment items given to different students, and use them to calibrate the level of difficulty of items that are not in common. It relies on psychometric calibration of these common items to achieve an alignment of standards and levels of difficulties.

- ✓ implies the use of common items in different assessment programmes
- ✓ one version has been proposed by the Australian Council for Educational Research (ACER) as part of an overall proposal of progression in learning but options are not exhausted.³

How? The methodology is based on the collection of hundreds of test questions from a wide range of assessment programs and then by ordering them by level of difficulty using statistical methods and experts' judgment to compare data from the different tests to construct a single scale for each domain. The test questions are analyzed to identify the kinds of skills required to answer each of the questions correctly. Each level of the reporting scales describes the skills for each segment of based on empirical test data. The reporting scales stem from these descriptions and are based on empirical evidence of skills demonstrated at similar points in the development of reading and mathematics, combined with expert judgments. All the items, or a useful subset of them, would then be calibrated using a set of common items at each "grade level." That being said, it is important to highlight that the methodology does not measure achievement. It is rather an approach to report learning outcomes in a meaningful and consistent way.

³ Note that the reference scale is built from items coming from various assessments.

Figure 4 - Simplified illustration of the item-based linking methodology



Source: UNESCO Institute for Statistics (UIS)

Table 4. Stock Taking on status and execution

	<i>Pedagogical calibration</i>	<i>Test-based</i>	<i>Item-based</i>
<i>Minimum Proficiency Level</i>			
<i>Accomplished thus far</i>	<ul style="list-style-type: none"> • Agreement on definition • Alignment of international assessments and population based • Unpacking of tasks • Compilation of items 		
<i>Remaining steps</i>			
• What?	Refinement of PLD and unpacking of tasks		
• When?	Early July		
<i>Linking Methodologies</i>			
<i>Accomplished thus far</i>	Piloting toolkit underdevelopment	<ul style="list-style-type: none"> ✓ Countries identified ✓ Regional organization ✓ Instruments elaborated ✓ Pilot Finished in LAC 	Compilation of items
<i>Remaining steps</i>			
• To be executed by	<ul style="list-style-type: none"> ✓ Second half 2019 ✓ 2 countries (USAID) ✓ 2 countries (UIS) ✓ 1 regional assessment (UIS) 	<ul style="list-style-type: none"> ✓ Main field of both regional and global assessments ✓ Concordance tables 	Bank of items and exchange platform
• When?	2019-2020	2019-2021	2019-

Source: UNESCO Institute for Statistics (UIS)