



NEPS *SURVEY PAPERS*

Inga Hahn

NEPS TECHNICAL REPORT FOR
SCIENCE: SCALING RESULTS OF
STARTING COHORT 1 FOR
SEVEN-YEAR-OLD CHILDREN

NEPS *Survey Paper* No. 86
Bamberg, April 2021

Survey Papers of the German National Educational Panel Study (NEPS)

at the Leibniz Institute for Educational Trajectories (LifBi) at the University of Bamberg

The NEPS *Survey Paper* series provides articles with a focus on methodological aspects and data handling issues related to the German National Educational Panel Study (NEPS).

They are of particular relevance for the analysis of NEPS data as they describe data editing and data collection procedures as well as instruments or tests used in the NEPS survey. Papers that appear in this series fall into the category of 'grey literature' and may also appear elsewhere.

The NEPS *Survey Papers* are edited by a review board consisting of the scientific management of LifBi and NEPS.

The NEPS *Survey Papers* are available at www.neps-data.de (see section "Publications") and at www.lifbi.de/publications.

Editor-in-Chief: Thomas Bäumer, LifBi

Review Board: Board of Directors, Heads of LifBi Departments, and Scientific Management of NEPS Working Units

Contact: German National Educational Panel Study (NEPS) – Leibniz Institute for Educational Trajectories – Wilhelmsplatz 3 – 96047 Bamberg – Germany – contact@lifbi.de

NEPS Technical Report for Science: Scaling Results of Starting Cohort 1 for Seven-Year-Old Children

Inga Hahn

Leibniz Institute for Science and Mathematics Education (IPN), Kiel, Germany

Email address of the author:

hahn@ipn.uni-kiel.de

Bibliographic data:

Hahn, I. (2021). *NEPS Technical Report for Science: Scaling Results of Starting Cohort 1 for Seven-Year-Old Children* (NEPS Survey Paper No. 86). Leibniz Institute for Educational Trajectories, National Educational Panel Study. <https://doi.org/10.5157/NEPS:SP86:1.0>

NEPS Technical Report for Science: Scaling Results of Starting Cohort 1 for Seven-Year-Old Children

Abstract

The National Educational Panel Study (NEPS) examines the development of competencies across the life span and develops tests for the assessment of different competence domains. In order to evaluate the quality of these competence tests various analyses based on item response theory (IRT) were performed. This paper describes the data and scaling procedures for the scientific literacy test that was administered to seven-year-old children of starting cohort 1. The scientific literacy test contained 21 items with different response formats representing different contexts as well as different areas of knowledge. The test was administered to 1,909 students. Their responses were scaled using a partial credit model. Item fit statistics, differential item functioning, Rasch-homogeneity, the test's dimensionality, and local item independence were evaluated to ensure the quality of the test. These analyses showed that the test exhibited an acceptable reliability and that all items but one fitted the model in a satisfactory way. Furthermore, test fairness could be confirmed for different subgroups. As the correlations between the two knowledge domains were very high, the assumption of unidimensionality seems adequate. The results revealed good psychometric properties of the scientific literacy test, thus supporting the estimation of a reliable scientific literacy score. Besides the scaling results, this paper also describes the data available in the scientific use file and provides the ConQuest syntax for scaling the data.

Keywords

scientific literacy, seven-year-old children, differential item functioning item response theory, scaling, scientific use file

Content

1	Introduction.....	4
2	Testing Scientific Literacy	4
3	Data	5
3.1	The design of the study	5
3.2	Sample.....	6
4	Analyses.....	7
4.1	Missing responses	7
4.2	Scaling model	7
4.3	Checking the quality of the test	7
4.4	Software	9
5	Results	9
5.1	Descriptive statistics of the responses.....	9
5.2	Missing Responses.....	9
5.2.1	Missing responses per person.....	9
5.2.2	Missing responses per item.....	11
5.3	Parameter estimates	15
5.3.1	Item parameters.....	15
5.3.2	Person parameters	15
5.3.3	Test targeting and reliability.....	15
5.4	Quality of the test.....	17
5.4.1	Fit of the subtasks of complex multiple-choice items.....	17
5.4.2	Distractor analyses	17
5.4.3	Item fit	17
5.4.4	Differential item functioning.....	17
5.4.5	Rasch-homogeneity.....	22
5.4.6	Unidimensionality of the test.....	22
6	Discussion	23
7	Data in the Scientific Use file.....	23
7.1	Naming conventions and scientific literacy scores	23
8	References.....	25

1 Introduction

Within the National Educational Panel Study (NEPS) different competencies are measured coherently across the lifespan (Blossfeld, Roßbach, & Maurice, 2011). These include, among others, reading competence, mathematical competence, scientific literacy, information and communication literacy, metacognition, vocabulary, and domain-general cognitive functioning. An overview of the competencies measured in the NEPS is given by Weinert et al. (2011) and by Fuß, Gnamb, Lockl, and Attig (2019).

Most of the competence data are scaled using models that are based on item response theory (IRT). Because most of the competence tests were developed specifically for implementation in the NEPS, several analyses were conducted to evaluate the quality of the tests. The IRT models chosen for scaling the competence data and the analyses performed for checking the quality of the scale are described in Pohl and Carstensen (2012).

In this paper, the results of these analyses are presented for a scientific literacy test that was administered to 7-Year-Old children of starting cohort 1. First, the main concepts of the scientific literacy test are introduced. Then, the scientific literacy data of starting cohort 1 and the analyses performed on the data to estimate competence scores and to check the quality of the test are described. Finally, an overview of the data that are available for public use in the scientific use file (SUF) is presented.

Please note that the analyses in this report are based on the data available at some time before public data release. Due to ongoing data protection and data cleansing issues, the data in the SUF may differ slightly from the data used for the analyses in this paper. However, we do not expect fundamental changes in the presented results.

2 Testing Scientific Literacy

The framework and test development for the scientific literacy test are described by Weinert et al. (2011) and by Hahn et al. (2013). In the following, we point out specific aspects of the scientific literacy test that are necessary for understanding the scaling results presented in this paper.

Scientific literacy is conceptualized as a unidimensional construct comprising two sub-dimensions. These are a) the knowledge of science (KOS) and b) the knowledge about science (KAS). KOS is specified as the knowledge of basic scientific concepts and facts whereas KAS can be regarded as the understanding of scientific processes.

KOS is divided into the content-related components of matter, system, development, and interaction. KAS is divided into the process-related components of scientific enquiry and scientific reasoning. KAS and KOS are implemented in three contexts: health, environment, and technology (see Figure 1). The test items are organized as single items or as units (testlets). One unit consists of two items. Each item or unit refers to one context-component-combination.

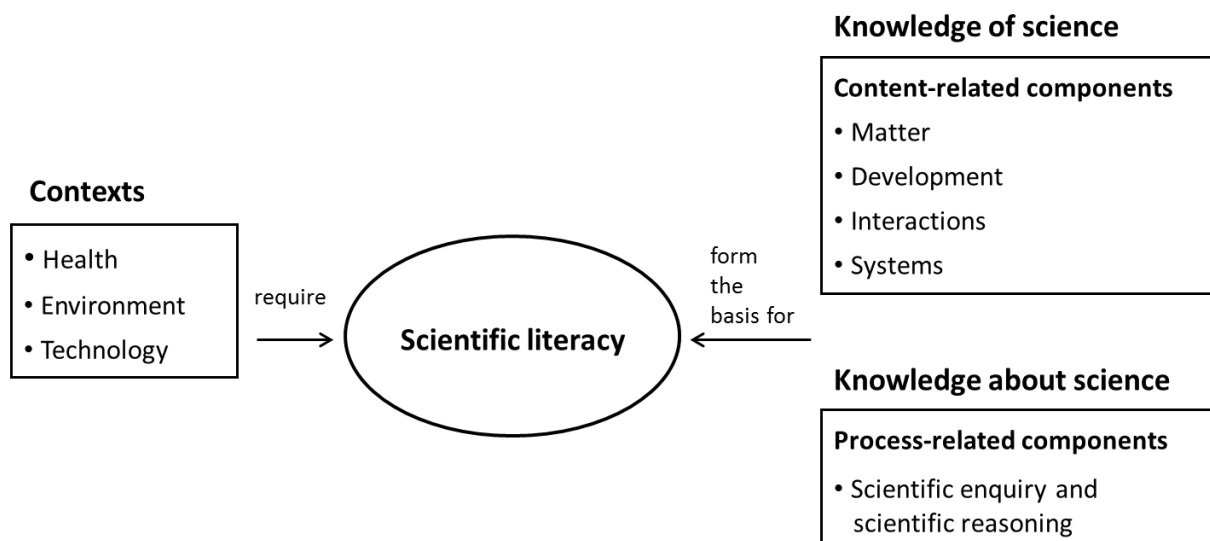


Figure 1. Assessment framework for scientific literacy (Hahn et al., 2013).

In the scientific literacy test for seven-year-old children of starting cohort 1 (Newborns), there were two types of response formats. These were simple multiple-choice (MC) and complex multiple-choice (CMC) in the special form of true-false items. In MC items the test taker had to identify the correct answer out of four response options. The three incorrect response options functioned as distractors. In CMC items four subtasks with two response options each (e.g., yes/ no) were presented.

3 Data

3.1 The design of the study

The study assessed different competence domains in the following order: listening comprehension at word level, phonological working memory, scientific literacy, metacognition, and executive control (delay of gratification). The scientific literacy test was administered after the phonological working memory test. Note that there was no multi-matrix design regarding the choice and the order of the items within a specific test. All children received the same science items in the same order. The testing time for the scientific literacy test was 20 minutes.

The allocation of the 21 items to the content areas (KOS and KAS) is summarized in Table 1. Table 2 shows how the items cover the different contexts of the scientific literacy framework (Hahn et al., 2013), whereas Table 3 gives an overview of the response formats.

Table 1:

Classification of Items into Knowledge Domains

Knowledge domains	Number of Items
Knowledge of Science (KOS)	14
Knowledge about Science (KAS)	7
Total number of items	21

Table 2:

Number of Items by Different Contexts

Context	Number of Items
Health	5
Environment	9
Technology	7
Total number of items	21

Table 3:

Number of Items by Response Formats

Response format	Number of Items
Simple Multiple-Choice	19
Complex Multiple-Choice (True-false items)	2
Total number of items	21

3.2 Sample

A total of 1,909 individuals received the scientific literacy test. For six participants less than three valid item responses were available. Because no reliable ability scores can be estimated based on such few valid responses, these cases were excluded from further analyses (Pohl & Carstensen, 2012). Thus, the analyses presented in this paper are based on a sample of 1,903 individuals (50.0 % girls). A detailed description of the study design, the sample, and the administered instrument is available on the NEPS website (<http://www.neps-data.de>).

4 Analyses

4.1 Missing responses

There are different kinds of missing responses. These are a) invalid responses, b) omitted items, c) items that test-takers did not reach, d) items that have not been administered, and e) multiple kinds of missing responses within CMC items that are not determined. In this study, all subjects received the same set of items so there are no missing responses due to items not being administered and since the test was tablet-based and the response formats were forced (SMC or CMC) there were also no invalid responses (a) or other multiple kinds of missing responses (e).

Missing responses provide information on how well the test worked (e.g., time limits, understanding of instructions, handling of different response formats) and need to be accounted for in the estimation of item and person parameters. We, therefore, thoroughly investigated the occurrence of missing responses in the test. First, we looked at the occurrence of the different types of missing responses per person. This indicated how well the persons were coping with the test. We then looked at the occurrence of missing responses per item to obtain some information on how well the items worked.

4.2 Scaling model

To estimate item and person parameters for scientific literacy, a partial credit model was used (PCM; Masters, 1982) that estimates item difficulties for dichotomous variables and location parameters for polytomous variables. Ability estimates for scientific literacy were estimated as weighted maximum likelihood estimates (WLEs). Item and person parameter estimation in NEPS is described in Pohl and Carstensen (2012), whereas the data available in the SUF are described in Section 7.

CMC items consisted of a set of subtasks that were aggregated to a polytomous variable for each CMC item, indicating the number of correctly solved subtasks within that item. If at least one of the subtasks contained a missing response, the whole CMC item was scored as missing. Categories of polytomous variables with less than $N = 200$ responses were collapsed to avoid possible estimation problems. This usually occurred for the lower categories of polytomous items; especially when the item consisted of many subtasks. In these cases, the lower categories were collapsed into one category. For both CMC items categories were collapsed (see Appendix A). To estimate item and person parameters, a scoring of 0.5 points for each category of the polytomous items was applied, while simple MC items were scored dichotomously as 0 for an incorrect and as 1 for the correct response (see Pohl & Carstensen, 2013, for studies on the scoring of different response formats).

4.3 Checking the quality of the test

The scientific literacy test was specifically constructed to be implemented in the NEPS. To ensure appropriate psychometric properties, the quality of the test was evaluated in several pretests and analyses.

Before aggregating the subtasks of CMC items to a polytomous variable, this approach was justified by preliminary psychometric analyses. For this purpose, the subtasks were analyzed together with the MC items in a Rasch model (Rasch, 1980). The fit of the subtasks was

evaluated based on the weighted mean square (WMNSQ), the respective t -value, point-biserial correlations of the correct responses with the total correct score, and the item characteristic curves. Only if the subtasks exhibited a satisfactory item fit, they were used to construct polytomous CMC variables that were included in the final scaling model.

The MC items consisted of one correct response and one or more distractors (i.e., incorrect response options). The quality of the distractors within MC items was examined using the point-biserial correlation between an incorrect response and the total score. Negative correlations indicate good distractors, whereas correlations between .00 and .05 are considered acceptable and correlations above .05 are viewed as problematic distractors (Pohl & Carstensen, 2012).

After aggregating the subtasks to polytomous variables, the fit of the dichotomous MC and polytomous CMC items to the partial credit model (Masters, 1982) was evaluated using three indices (Pohl & Carstensen, 2012). Items with a WMNSQ > 1.15 (t -value $> |6|$) were considered as having a noticeable item misfit, and items with a WMNSQ > 1.20 (t -value $> |8|$) were judged as having a considerable item misfit and their performance was further investigated. Correlations of the item score with the corrected total score (equal to the corrected discrimination as computed in ConQuest) greater than .30 were considered as good, greater than .20 as acceptable, and below .20 as problematic. The overall judgment of the fit of an item was based on all fit indicators.

Scientific literacy should measure the same construct for all children. If any items favored certain subgroups (e.g., if they were easier for boys than for girls), measurement invariance would be violated and a comparison of competence scores between these subgroups (e.g., boys and girls) would be biased and thus unfair. For the present study, test fairness was investigated for the variables gender, the number of books at home (as a proxy for socioeconomic status), and migration background (see Pohl & Carstensen, 2012, for a description of these variables). Differential item functioning (DIF) analyses were estimated using a multigroup IRT model, in which the main effects of the subgroups as well as differential effects of the subgroups on item difficulty were modeled. Based on experiences with preliminary data, we considered absolute differences in estimated difficulties between the subgroups that were greater than 1 logit as very strong DIF, absolute differences between 0.6 and 1 as noteworthy of further investigation, differences between 0.4 and 0.6 as considerable but not severe, and differences smaller than 0.4 as negligible DIF. Additionally, the test fairness was examined by comparing the fit of a model including differential item functioning to a model that only included main effects and no DIF.

The scientific literacy test was scaled using the PCM (Masters, 1982), which assumes Rasch-homogeneity. The PCM was chosen because it preserves the weighting of the different aspects of the framework as intended by the test developers (Pohl & Carstensen, 2012). Nonetheless, Rasch-homogeneity is an assumption that might not hold for empirical data. To test the assumption of equal item discrimination parameters, a generalized partial credit model (GPCM; Muraki, 1992) was also fitted to the data and compared to the PCM.

The science test was constructed to measure a unidimensional scientific literacy score (Hahn et al., 2013). The assumption of unidimensionality was, nevertheless, tested by specifying a two-dimensional model with process-related items (KAS) representing one and content related items (KOS) the other dimension. The correlation between the subdimensions as well

as differences in model fit between the unidimensional model and the two-dimensional model were used to evaluate the unidimensionality of the test.

Moreover, we examined whether the residuals of the one-dimensional model exhibited approximately zero-order correlations as indicated by Yen's Q3 (Yen, 1984). Because in the case of locally independent items, the Q3 statistic tends to be slightly negative, we report the corrected Q3 that has an expected value of 0. Following prevalent rules-of-thumb (Yen, 1993) values of Q3 falling below .20 indicate that the assumption of local item dependence (LID) is essentially met.

4.4 Software

The IRT models were estimated in ConQuest version 4.2.5 (Adams, Wu, & Wilson, 2015).

5 Results

All but one of the 21 items (including all subtasks for the polytomous items) were included in the analyses. Item scg31610_sc1n8_c was excluded from the analyses due to an insufficient discrimination, t-value, and differential item functioning.

5.1 Descriptive statistics of the responses

To a) get a first rough descriptive measure of the item difficulties and b) check for possible estimation problems, before performing IRT analyses we evaluated the relative frequency of the responses given. The percentage of persons correctly responding to an item (relative to all valid responses) ranged from 14.0% to 88.6% for the MC items. For the CMC items, the percentage of persons who correctly answered all subtasks varied between 12.3% and 41.5%.

5.2 Missing Responses

5.2.1 Missing responses per person

Figure 2 shows the number of omitted responses per person. As illustrated in Figure 2 most respondents, 88.9%, did not skip any item, and less than 0.2% omitted more than two items.

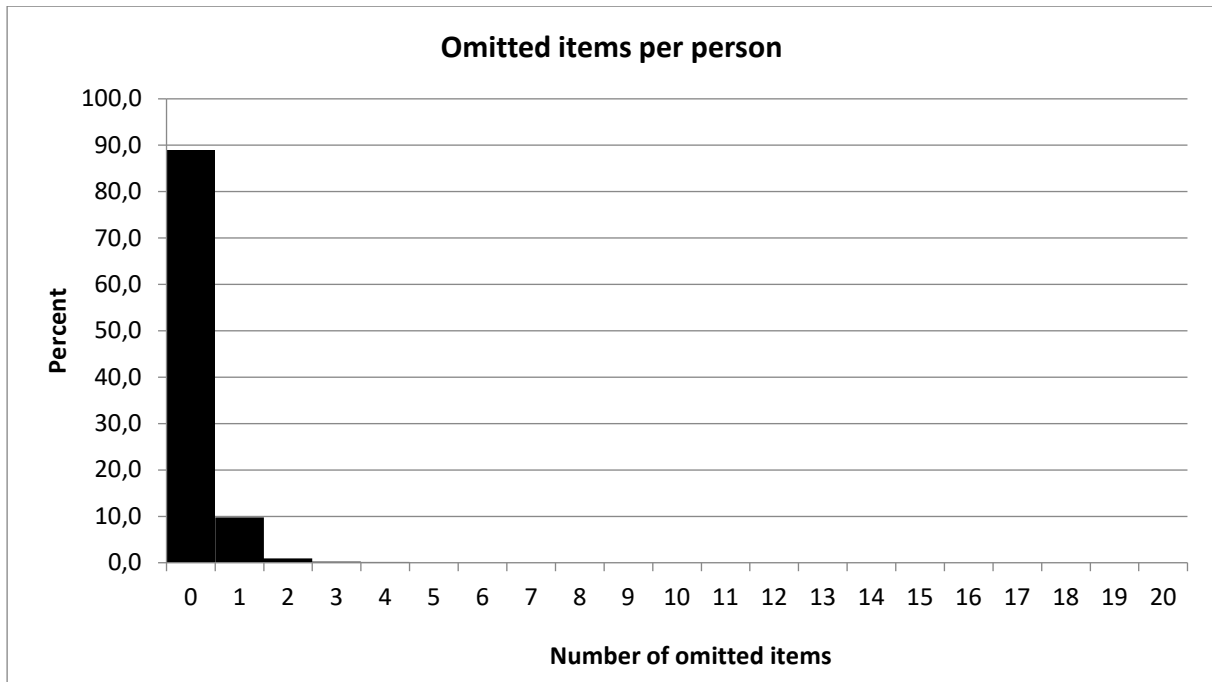


Figure 2. Number of omitted responses per person.

Another source of missing responses are items that were not reached by the respondents; these are all missing responses after the last valid response. The number of not-reached items was very low, about 99.7% of the respondents were able to finish the test within the allocated time limit (Figure 3).

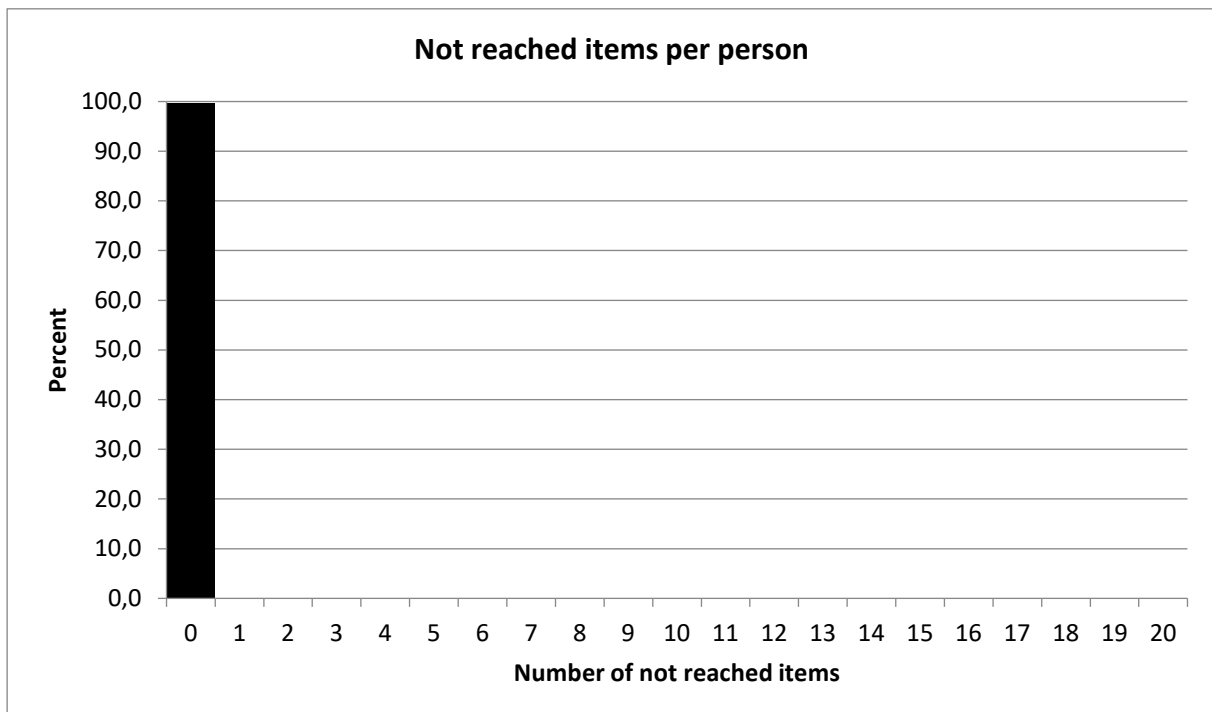


Figure 3. Number of not reached items per person.

The total number of missing responses, aggregated over omitted and not-reached missing responses, is illustrated in Figure 4. 88.9% of the students answered all questions and,

consequently, had no missing responses. Only 0.2% of the students had 5 or more missing responses. Hence, the number of missing responses per person can be classified as very small.

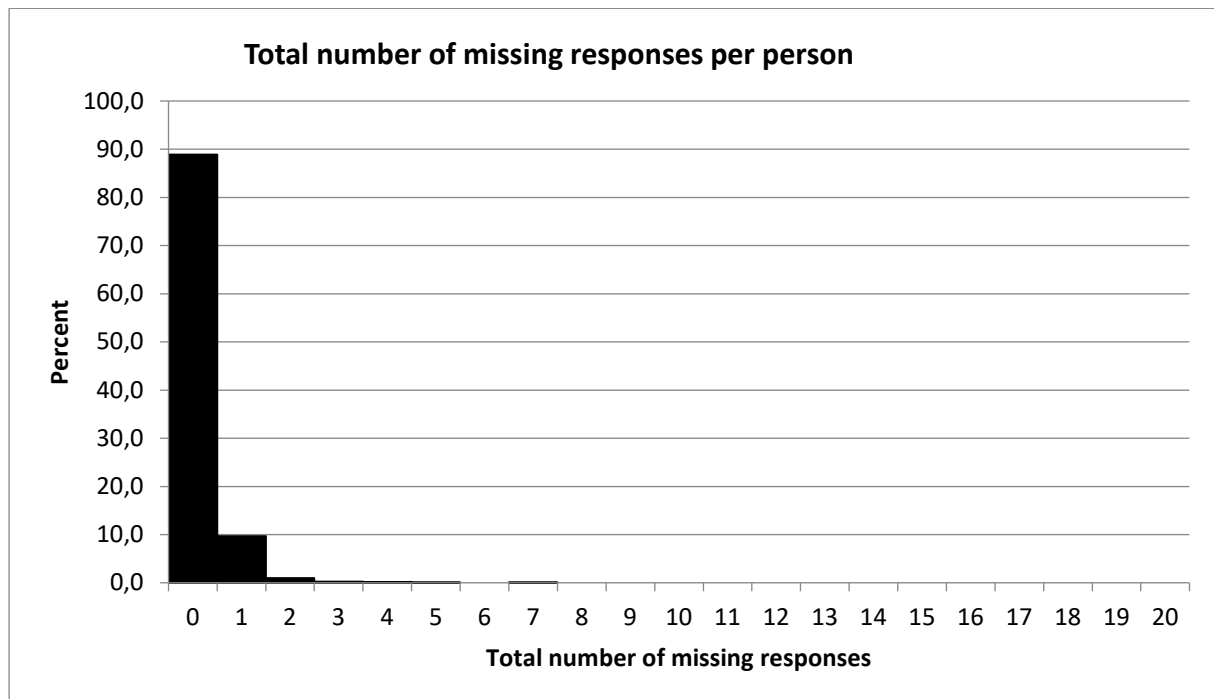


Figure 4. Total number of missing responses per person.

5.2.2 Missing responses per item

Table 4 shows the number of valid responses for each item as well as the percentage of missing responses. Overall, omission rates were very low, varying across items between 0.1% and 2.5%. Thus, there was no item with an omission rate exceeding 10.0%. The number of missing responses was uncorrelated ($r = -.381$, $p = .097$) with the difficulty of the item. This result indicates that the test-takers did not omit more difficult items. The relative frequency of not reached items increased towards the end of the test. Eventually, 0.2% of the students did not reach the last item and, thus, did not complete the test. The total number of missing responses per item varied between 0.1% and 2.5%.

Table 4:

Valid Responses and Missing Values

Item	Position in the test	Number of valid responses	Not reached items (%)	Omitted items (%)
scg10820_sc1n8_c	1	1881	0.0	1.2
scg10840_sc1n8_c	2	1855	0.0	2.5
scg11510_sc1n8_c	3	1888	0.0	0.8
scg16510_sc1n8_c	4	1901	0.0	0.1
scg1652s_sc1n8_c	5	1884	0.0	1.0
scg10920_sc1n8_c	6	1880	0.0	1.2
scg1011s_sc1n8_c	7	1897	0.0	0.3
scg11110_sc1n8_c	8	1883	0.0	1.1
scg11130_sc1n8_c	9	1881	0.0	1.2
scg16530_sc1n8_c	10	1900	0.0	0.2
scg16020_sc1n8_c	11	1899	0.0	0.2
scg16030_sc1n8_c	12	1893	0.0	0.5
scg11610_sc1n8_c	13	1901	0.0	0.1
scg10310_sc1n8_c	14	1899	0.0	0.2
scg10520_sc1n8_c	15	1881	0.1	1.1
scg16310_sc1n8_c	16	1899	0.1	0.1
scg16220_sc1n8_c	17	1893	0.1	0.2
scg33710_sc1n8_c	18	1893	0.1	0.4
scg31010_sc1n8_c	19	1893	0.2	0.4
scg30109_sc1n8_c	21	1896	0.2	0.2

Table 5:

Item parameters

No.	Item	Percentage correct	Difficulty/location parameter	SE (difficulty/location parameter)	WMNSQ	t-value for WMNSQ	Pt.-bis. Corr. of correct response	Discrimination (GPCM)	Yens Q3
1	scg10820_sc1n8_c	78.2	-1.509	0.063	0.98	-0.50	0.39	0.92	0.09
2	scg10840_sc1n8_c	48.9	-0.002	0.053	0.97	-2.10	0.47	1.03	0.11
3	scg11510_sc1n8_c	30.0	0.956	0.057	0.98	-0.70	0.42	0.90	0.06
4	scg16510_sc1n8_c	71.8	-1.069	0.057	1.00	-0.20	0.41	0.86	0.05
5	scg1652s_sc1n8_c	n.a.	-0.831	0.069	0.97	-1.20	0.38	1.06	0.09
6	scg10920_sc1n8_c	87.5	-2.294	0.078	1.00	-0.10	0.31	0.86	0.12
7	scg1011s_sc1n8_c	n.a.	1.590	0.077	0.99	-0.20	0.33	0.91	0.12
8	scg11110_sc1n8_c	80.2	-1.638	0.065	1.00	0.00	0.37	0.86	0.10
9	scg11130_sc1n8_c	77.1	-1.433	0.062	0.93	-2.30	0.48	1.38	0.07
10	scg16530_sc1n8_c	19.2	1.618	0.064	1.05	1.50	0.26	0.49	0.08
11	scg16020_sc1n8_c	13.9	2.038	0.072	1.01	0.30	0.27	0.66	0.08
12	scg16030_sc1n8_c	73.6	-1.186	0.059	0.92	-3.00	0.51	1.46	0.09
13	scg11610_sc1n8_c	49.7	0.018	0.053	1.03	2.00	0.37	0.60	0.06

14	scg10310_sc1n8_c	49.3	0.030	0.053	1.09	5.40	0.29	0.35	0.07
15	scg10520_sc1n8_c	71.4	-1.084	0.058	0.99	-0.20	0.39	0.82	0.12
16	scg16310_sc1n8_c	65.2	-0.720	0.055	1.00	0.20	0.42	0.85	0.15
17	scg16220_sc1n8_c	53.8	-0.177	0.053	0.96	-2.30	0.48	1.05	0.12
18	scg33710_sc1n8_c	37.5	0.578	0.054	1.02	0.80	0.38	0.70	0.15
19	scg31010_sc1n8_c	41.8	0.371	0.053	1.07	4.50	0.31	0.41	0.19
20	scg30109_sc1n8_c	65.8	-0.759	0.055	0.99	0.50	0.43	0.85	0.19

Note. SE = Standard error of item difficulty / location parameter, WMNSQ = Weighted mean square, t = t -value for WMNSQ. Percent correct scores are not informative for polytomous CMC (denoted by n.a.) For the dichotomous and polytomous items, the item-total correlation corresponds to the point-biserial correlation between the correct response and the total score (discrimination value as computed in ConQuest).

5.3 Parameter estimates

5.3.1 Item parameters

Column 3 in Table 5 shows the percentage of correct responses in relation to all valid responses for each item. Note that although the amount of missing responses was low, this probability cannot be interpreted as an index for item difficulty. The percentage of correct responses within items varied between 13.9% and 87.5% with an average of 53.4% ($SD = 22.9$) correct responses.

The estimated item difficulties (for dichotomous items, MC items) and location parameters (for polytomous variables, CMC items) are also given in Table 5. The step parameters (for polytomous variables) are depicted in Table 6. All CMC items showed less than $N = 200$ participants in the two lowest categories, thus the three lowest categories were collapsed. These items were scaled using a scoring of 0, 0.5, and 1. The item difficulties were estimated by constraining the mean of the ability distribution to be zero. The estimated item difficulties (or location parameters for polytomous variables) ranged between -2.29 (scg10920_sc1n8_c) and 2.04 (scg16020_sc1n8_c). In total, the estimated item difficulties had a mean of -0.28 ($SD = 1.19$). Due to the large sample size, the standard errors of the estimated item difficulties were very small ($SE(\beta) \leq 0.078$). Overall, the item difficulties fitted the person abilities very well. The test was only slightly too easy.

Table 6:

Step parameters for the CMC items

Item	Step 1 (SE)	Step 2
scg1652s_sc1n8_c	-0.401 (0.049)	0.401
scg1011s_sc1n8_c	-0.381 (0.053)	0.381

Note. The last step parameters are not estimated and have, thus, no standard error because they are constrained parameters for model identification.

5.3.2 Person parameters

Person parameters are estimated as WLEs (Pohl & Carstensen, 2012). A description of the data in the SUF can be found in section 7. An overview of how to work with competence data is given in Pohl and Carstensen (2012).

5.3.3 Test targeting and reliability

Test targeting focuses on comparing the item difficulties with the person abilities (WLEs) to evaluate the appropriateness of the test for the specific target population. In Figure 5, the difficulties of the scientific literacy items and the ability of the test takers are plotted on the same scale. The distribution of the estimated test takers' ability is mapped onto the left side whereas the right side shows the distribution of item difficulties.

The mean of the ability distribution was constrained to be zero. The variance was estimated to be 0.662, indicating a somewhat limited variability between subjects. The reliability of the test (EAP/PV reliability = .695; WLE reliability = .670) was acceptable. Although the items

covered a wide range of the ability distribution, there could have been one or two more items covering the upper peripheral ability areas. Overall all ability regions seem to be covered quite well.

Scale in logits	Person ability	Item difficulty
2	X	11
	X	
	X	
	X	
	X	7 10
	X	
	XX	
	XXX	
	XXXX	
	XXXXX	
	XXXXX	3
1	XXXXX	
	XXXXXXXX	18
	XXXXXXXX	
	XXXXXXXXXX	19
	XXXXXXXXXX	
	XXXXXXXXXX	13 14
	XXXXXXXXXX	2
	XXXXXXXXXX	17
	XXXXXXXXXX	
	XXXXXXXXXX	
	XXXXXXXXXX	
0	XXXXXXXXXX	16 20
	XXXXXX	5
	XXXXX	
	XXXXX	4 15
	XXX	12
	XXX	
	XX	
	XX	1 9
	X	
	X	8
	X	
-1		
-2		6

Figure 5. Test targeting. The distribution of person abilities in the sample is depicted on the left side of the graph. Each 'X' represents 18.1 cases. The difficulty of the items is depicted on the right side of the graph. Each number represents an item (see Table 5).

5.4 Quality of the test

5.4.1 Fit of the subtasks of complex multiple-choice items

Before the subtasks of the CMC item were aggregated and analyzed via a partial credit model, the fit of the subtasks was checked by analyzing the single subtasks together with the MC items in a Rasch model. Counting the subtasks of the CMC item separately, there were 26 items. The percentage of a correct response ranged from 14.0% to 88.6% across all items (*Mdn* = 69.0%). Thus, the number of correct and incorrect responses was reasonably large. All subtasks of the CMC items showed a satisfactory item fit with a good WMNSQ, ranging from 0.91 to 1.07. The respective *t*-values ranged from -3.5 to 4.7, and there were no noticeable deviations of the empirically estimated probabilities from the model-implied item characteristic curves. Due to the good model fit of the subtasks, their aggregation to a polytomous variable seemed justified.

5.4.2 Distractor analyses

In addition to the overall item fit, we specifically investigated how well the distractors performed in the test by evaluating the point-biserial correlation between each incorrect response (distractor) and the students' total score. There were two items with distractors showing a point-biserial correlation with the total scores above zero: *scg11510_sc1n8* and *scg16530_sc1n8*. For all the other items (including the CMC items) the results indicate that the distractors worked well.

5.4.3 Item fit

The evaluation of the item fit was performed based on the final scaling model, the partial credit model, using the MC items and the CMC items. Altogether, the item fit can be considered to be very good (see Table 5). Values of the WMNSQ ranged from 0.92 (item *scg16030_sc1n8*) to 1.09 (item *scg10310_sc1n8*). All *t*-values of the WMNSQ ranged below 6. Thus, there was no indication of a severe item over- or underfit. Point-biserial correlations between the item scores and the total scores ranged from .26 (item *scg16530_sc1n8*) to .51 (items *scg16030_sc1n8*) and had a mean of .38. All item characteristic curves showed a good fit of the items to the PCM.

5.4.4 Differential item functioning

Differential item functioning (DIF) was used to evaluate test fairness for several subgroups (i.e., measurement invariance). For this purpose, DIF was examined for the variables gender, the number of books at home (as a proxy for socioeconomic status) and migration background, and school type (see Pohl & Carstensen, 2012, for a description of these variables). Table 7 shows the absolute difference between the estimated item difficulties in different groups. Male vs. female, for example, indicates the difference in difficulty $\beta(\text{male}) - \beta(\text{female})$. A positive value indicates a higher difficulty for males, a negative value a lower difficulty for males as opposed to females. Also, Table 8 shows the main effect for the examined subgroups (inclusive Cohen's *d*).

Table 7:

Differential item functioning (differences between difficulties)

Item	Gender		Books		Migration status		
	Male vs. female	<100 vs. >100	<100 vs. missing	>100 vs. missing	Without vs. With	Without vs. Missing	With vs. Missing
scg10820_sc1n8_c	0.194	0.232	-0.052	-0.290	-0.112	0.716	0.822
scg10840_sc1n8_c	0.124	0.090	0.030	-0.062	-0.230	0.104	0.326
scg11510_sc1n8_c	-0.458	0.030	-0.148	-0.174	0.062	-0.284	-0.354
scg16510_sc1n8_c	-0.058	0.278	0.118	-0.166	0.078	-0.178	-0.264
scg1652s_sc1n8_c	0.376	0.220	0.102	-0.112	-0.100	-0.188	-0.080
scg10920_sc1n8_c	0.218	0.306	0.122	-0.192	-0.068	-0.020	0.040
scg1011s_sc1n8_c	0.412	0.082	0.000	-0.084	-0.192	-0.600	-0.416
scg11110_sc1n8_c	-0.662	0.158	-0.218	-0.378	-0.038	0.106	0.138
scg11130_sc1n8_c	-0.068	0.428	0.192	-0.246	-0.434	-0.452	-0.028
scg16530_sc1n8_c	0.198	0.422	-0.008	0.418	0.264	0.324	0.054
scg16020_sc1n8_c	-0.058	0.396	-0.604	-0.198	-0.208	-0.472	-0.274
scg16030_sc1n8_c	0.176	0.194	0.088	-0.112	0.010	-0.298	-0.316

scg11610_sc1n8_c	-0.524	0.168	-0.114	0.054	0.174	-0.006	-0.188
scg10310_sc1n8_c	0.348	0.208	-0.098	0.108	0.242	0.246	-0.004
scg10520_sc1n8_c	0.546	0.056	-0.328	-0.272	-0.028	-0.008	0.012
scg16310_sc1n8_c	-0.128	0.110	0.092	0.196	0.150	-0.050	-0.206
scg16220_sc1n8_c	0.008	0.052	0.066	0.116	-0.384	-0.234	0.140
scg33710_sc1n8_c	-0.250	0.062	0.260	0.318	0.136	0.512	0.368
scg31010_sc1n8_c	0.088	0.170	0.150	0.318	0.490	0.604	0.108
scg30109_sc1n8_c	-0.146	0.020	0.058	0.072	-0.194	-0.350	-0.164

Gender

The sample included 952 (50.1%) male test-takers (coded 0) and 951 (49.9%) female test-takers (coded 1). On average, male students had slightly higher scores in scientific literacy than female students (main effect = 0.152 logits, Cohen's $d = 0.188$). There was one item showing considerable gender DIF up to -0.662 (item `scg11110_sc1n8_c`). Since this item displayed a good item fit and showed no other DIF it remained in the analysis.

Books

The number of books at home was used as a proxy for socioeconomic status. There were 533 (28.0%) test takers with 0 to 100 books at home (coded 0), 1,261 (66.3%) test takers with more than 100 books at home (coded 1), and 109 (5.7%) test-takers did not give a valid response (coded 9). DIF was investigated using these three groups. There were considerable average differences between these three groups. Participants with 100 or fewer books at home on average showed lower scientific literacy scores than participants with more than 100 books (main effect = -0.606 logits, Cohen's $d = -0.810$). Participants with up to 100 books performed worse than participants without a valid response on the variable 'books at home' (main effect = -0.148 logits, Cohen's $d = -0.179$) while participants with more than 100 books at home scored higher than participants without a valid response (main effect = 0.458 logits, Cohen's $d = 0.607$). There was no considerable DIF comparing participants with many or fewer books (highest DIF = 0.428) and there was also no considerable DIF comparing the group without valid responses to the two groups with valid responses (highest DIF = 0.418 logits).

Migration background

There were 1,371 (72.0%) participants without a migration background (coded 0) and 442 (23.2%) participants with a migration background (coded 1). A total of 90 (4.7%) students could not be allocated to either group. These groups were used for investigating DIF of migration. There was a considerable difference in the average performance of participants with or without migration background. Participants without a migration background showed higher scientific literacy scores than participants with a migration background (main effect = 0.474 logits, Cohen's $d = 0.605$) and also higher scores than students with an unknown background on migration (main effect = 0.738 logits, Cohen's $d = 0.973$). Furthermore, students with a migration background scored higher than those with an unknown background on migration (main effect = 0.268 logits, Cohen's $d = 0.342$). There was no considerable DIF comparing participants with and without a migration background (highest DIF = 0.490). Comparing the group without valid responses to the two groups with valid responses, there was one item with high DIF (items `scg10820_sc1n8_c`). Since this item displayed a good item fit and showed no other DIF it remained in the analysis.

Table 8:

Main effects and Cohen's d of the examined subgroups

Variables	Subgroups	Main effect	Cohen's d
Gender	Male (0)	0.152	0.188
	Female (1)		
Books	0 to 100 books at home (0)	-0.606	-0.810
	More than 100 books at home (1)		
	0 to 100 books at home (0)	-0.130	-0.179
	Invalid response (9)		
	More than 100 books at home (1)	0.458	0.607
	Invalid response (9)		
Migration background	Without migration background (0)	0.474	0.605
	With migration background (1)		
	Without migration background (0)	0.748	0.973
	Invalid response (9)		
	With migration background (1)	0.268	0.342
	Invalid response (9)		

Note. The numbers behind the subgroups display their coding.

Besides investigating DIF for every single item, an overall test for DIF was performed by comparing models that allow for DIF with those that allow only for main effects. In Table 9, the models including only the main effects are compared with those that additionally estimate DIF. Akaike's (1974) information criterion (AIC) and the Bayesian information criterion (BIC, Schwarz, 1978) were used for comparing the models. The AIC favored the model considering DIF only for the migration background. For the DIF variables books and gender, the AIC favored the model which allows only for main effects. The BIC takes the number of estimated parameters into account and, thus, prevents from overparameterization of models. Using BIC, the more parsimonious model including only the main effect is preferred over the more complex DIF model for all DIF variables.

Table 9:

Comparison of models with and without DIF

DIF variable	Model	Deviance	N	Number of parameters	AIC	BIC
Gender	main effect	43287.44	1813	24	43335.44	43467.51
	DIF	43314.45	1813	44	43402.45	43644.57
Books	main effect	42835.60	1794	24	42883.60	43015.42
	DIF	42868.79	1794	44	42956.79	43198.45
Migration background	main effect	43287.44	1813	24	43335.44	43467.51
	DIF	43226.45	1813	44	43314.45	43556.57

Note. All analyses are based on cases without missing values on the grouping variable.

5.4.5 Rasch-homogeneity

An essential assumption of the Rasch (1980) model is that all item-discrimination parameters are equal. To test this assumption, a generalized partial credit model (GPCM; Muraki, 1992) that estimates discrimination parameters was fitted to the data. The estimated discriminations differed moderately among items (see Table 5), ranging from 0.35 (item scg10310_sc1n8) to 1.46 (item scg16030_sc1n8). The average discrimination parameter fell at 0.85. Model fit indices suggested a better model fit of the GPCM (AIC = 45,469.96, BIC = 45,703.10) as compared to the PCM model (AIC = 45,669.74, BIC = 45,797.41). Despite the empirical preference for the GPCM, the PCM model matches the theoretical conceptions underlying the test construction more adequately (see Pohl & Carstensen, 2012, 2013, for a discussion of this issue). For this reason, the partial credit model was chosen as our scaling model to preserve the item weightings as intended in the theoretical framework.

5.4.6 Unidimensionality of the test

The dimensionality of the test was investigated by specifying a one- and a two- dimensional model. The first model is based on the assumption that scientific literacy is a one-dimensional construct that measures one distinct competence whereas the second model distinguishes between the two sub-competencies: the process-related components (knowledge about science – KAS) and the content-related components (knowledge of science – KOS; for more details see Hahn et al., 2013). For estimating a two-dimensional model Gauss' Hermite quadrature estimation in ConQuest was used (nodes were chosen in such a way that stable parameter estimation was obtained). The unidimensional model (BIC = 45,797.41, number of parameters = 23) fitted the data slightly better than the two-dimensional model (BIC = 45,799.15, number of parameters = 25). Also, the correlation between the two dimensions was very high ($r = .94$). So the one-dimensional measurement model was used to estimate a single competence score for scientific literacy.

6 Discussion

The analyses in the previous sections aimed at providing detailed information on the quality of the science test for seven-year-old children of starting cohort 1 and at describing how scientific literacy was estimated.

We investigated different kinds of missing responses and examined the item and test parameters. We checked item fit statistics for simple MC items, subtasks of CMC items, as well as the polytomous CMC items and examined the correlations between correct and incorrect responses and the total score. Further quality inspections were conducted by examining differential item functioning, testing Rasch-homogeneity, investigating the tests' dimensionality as well as local item dependence.

Various criteria indicated a good fit of the items and measurement invariance across various subgroups. The number of missing responses was very small.

The test had an acceptable reliability and distinguished well between test-takers. The test's variance was acceptable.

Indicated by various fit criteria – WMNSQ, t -value of the WMNSQ – the items exhibited a good item fit. Also, discrimination values of the items (either estimated in a GPCM or as a correlation of the item score with total score) were acceptable. Different variables were used for testing measurement invariance across various subgroups. Only two items showed considerable DIF for the examined variables, indicating that the test was fair to the considered subgroups.

Fitting a two-dimensional partial credit model (the dimensions being the “content-related components” and the “process-related components”) yielded no better model fit than the unidimensional partial credit model. Also, the high correlation between the two dimensions indicates that a unidimensional model describes the data reasonably well.

Summarizing the results, the test had good psychometric properties that facilitated the estimation of a unidimensional scientific literacy score.

7 Data in the Scientific Use file

7.1 Naming conventions and scientific literacy scores

There are 21 items in the data set that are either scored as dichotomous variables (MC items) with 0 indicating an incorrect response and 1 indicating a correct response or scored as a polytomous variable (CMC items) indicating the (partial) credit. The dichotomous variables are marked with a ‘_c’ at the end of the variable name, the CMC items are marked with a ‘s_c’ at the end of the variable name. Note that the value of the polytomous variable does not necessarily indicate the number of correctly responded subtasks (see section 4.2 aggregation of CMC items). In the scaling model, each category of CMC items was scored with 0.5 points. Please note that when categories of the polytomous variables had less than 200 valid responses, the categories were collapsed. For the science test, this concerned the three lowest categories of all of the polytomous items (see section 5.3.1 on the aggregation of CMC items). In the scaling model, the collapsed polytomous item was scored in steps of 0, 0.5, 1.0 (denoting the highest).

Manifest scale scores are provided in form of WLE estimates (scn8_sc1) including the respective standard error (scn8_sc2). These WLE estimates can only be used for cross-sectional analyses because the study which was supposed to link the competence scores of this study to the competence scores of the preceding study (B102, five-year-old children) had to be postponed due to the corona pandemic. The linked WLE estimates will be available at a later time.

The ConQuest Syntax for estimating the WLE scores from the items is provided in Appendix A. For students who either did not take part in the science test or who did not give enough valid responses, no WLE is estimated. The value on the WLE and the respective standard error for these persons are denoted as not-determinable missing values. Alternatively, users interested in examining latent relationships may either include the measurement model in their analyses or estimate plausible values. A description of these approaches can be found in Pohl and Carstensen (2012). Plausible values for the literacy test can be estimated using the R package *NEPSScaling* (Scharl, Carstensen, & Gnams, 2020).

8 References

- Adams, R. J., Wu, M. L., & Wilson, M. R. (2015). ACER ConQuest: Generalised Item Response Modelling Software (Version 4) [Computer software]. Camberwell: Australian Council for Educational Research.
- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, *19*, 716–722.
- Blossfeld, H.-P., Roßbach, H.-G., & Maurice, J. v. (2011). Education as a Lifelong Process – The German National Educational Panel Study (NEPS). *Zeitschrift für Erziehungswissenschaft, Sonderheft 14*.
- Fuß, D., Gnamb, T., Lockl, K., & Attig, M. (2019). *Competence data in NEPS: Overview of measures and variable naming conventions (Starting Cohorts 1 to 6)*. Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study.
- Hahn, I., Schöps, K., Rönnebeck, S., Martensen, M., Hansen, S., Saß, S., . . . Prenzel, M. (2013). Assessing scientific literacy over the lifespan – A description of the NEPS science framework and the test development. *Journal for Educational Research Online*, *5*(2), 110–138.
- Masters, G. N. (1982). A Rasch model for partial credit scoring. *Psychometrika*, *47*, 149–174.
- Muraki, E. (1992). A generalized partial credit model: Application of an EM algorithm. *Applied Psychological Measurement*, *16*, 159–176.
- Pohl, S., & Carstensen, C. H. (2012). *NEPS technical report – Scaling the data of the competence tests* (NEPS Working Paper No. 14). Bamberg, Germany: Otto-Friedrich-Universität, Nationales Bildungspanel.
- Pohl, S., & Carstensen, C. H. (2013). Scaling of competence tests in the National Educational Panel Study – Many questions, some answers, and further challenges. *Journal for Educational Research Online*, *5*, 189–216.
- Rasch, G. (1980). *Probabilistic models for some intelligence and attainment tests*. Chicago, IL: University of Chicago Press.
- Scharl, A., Carstensen, C. H., & Gnamb, T. (2020). *Estimating plausible values with NEPS data: An example using reading competence in starting cohort 6* (NEPS Survey Paper No. 71). Bamberg: Leibniz Institute for Educational Trajectories, National Educational Panel Study.
doi:10.5157/NEPS:SP71:1.0
- Schwarz, G. E. (1978). Estimating the dimension of a model. *Annals of Statistics*, *6*(2), 461–464.
- Weinert, S., Artelt, C., Prenzel, M., Senkbeil, M., Ehmke, T., & Carstensen, C. H. (2011). Development of competencies across the life span. *Zeitschrift für Erziehungswissenschaft, Sonderheft 14*, 67–86.

Yen, W. M. (1984). Effects of local item dependence on the fit and equating performance of the three-parameter logistic model. *Applied Psychological Measurement*, 8(2), 125–145.

Yen, W. M. (1993). Scaling performance assessments: Strategies for managing local item dependence. *Journal of Educational Measurement*, 30(3), 187–213.

Appendix

Appendix A: ConQuest-Syntax for estimating WLE estimates in starting cohort II

Title Seven-Year-Old Children (SC1) Science analysis, Partial Credit Model;

data filename.dat;

format id 1–7 responses 8–27;

labels << filename_with_labels.txt;

recode (0,1,2,3,4) (0,0,0,1,2) !item (5,7);

codes 0,1,2;

score (0,1) (0,1) !item (1–4,6,8–20);

score (0,1,2) (0,0.5,1) !item (5,7);

set constraint=cases;

model item + item*step;

estimate; method=gauss, nodes=45;

show cases !estimates=wle >> filename.wle;

show ! estimates=latent >> filename.shw;

itanal! estimates=latent >> filename.ita;

Appendix B: Assignment of items to the content and process-related components and contexts

Variable name	Position in the test	Component	Context
scg10820_sc1n8_c	1	KOS	Health
scg10840_sc1n8_c	2	KOS	Health
scg11510_sc1n8_c	3	KOS	Technology
scg16510_sc1n8_c	4	KAS	Environment
scg1652s_sc1n8_c	5	KAS	Environment
scg10920_sc1n8_c	6	KOS	Technology
scg1011s_sc1n8_c	7	KOS	Health
scg11110_sc1n8_c	8	KOS	Environment
scg11130_sc1n8_c	9	KOS	Environment
scg16530_sc1n8_c	10	KAS	Technology
scg16020_sc1n8_c	11	KAS	Environment
scg16030_sc1n8_c	12	KAS	Environment
scg11610_sc1n8_c	13	KOS	Environment
scg10310_sc1n8_c	14	KOS	Technology
scg10520_sc1n8_c	15	KOS	Environment
scg16310_sc1n8_c	16	KAS	Technology
scg16220_sc1n8_c	17	KAS	Technology
scg33710_sc1n8_c	18	KOS	Technology
scg31010_sc1n8_c	19	KOS	Environment
scg31610_sc1n8_c	20	KOS	Health
scg30109_sc1n8_c	21	KOS	Health

Note. KOS = knowledge of science (content-related components); KAS = knowledge about science (process-related components)