



NEPS SURVEY PAPERS

Jana Kähler

NEPS TECHNICAL REPORT FOR SCIENCE: SCALING RESULTS OF STARTING COHORT 2 FOR GRADE 1

NEPS Survey Paper No. 58
Bamberg, November 2019

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).

The NEPS Survey Papers are edited by a review board consisting of the scientific management of LifBi and 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 available at <https://www.neps-data.de> (see section "Publications").

Editor-in-Chief: Corinna Kleinert, LifBi/University of Bamberg/IAB Nuremberg

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 2 for Grade 1

Jana Kähler

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

Email address of the lead author:

jkaehler@ipn.uni-kiel.de

Bibliographic data:

Kähler, J. (2019). *NEPS Technical Report for Science: Scaling Results of Starting Cohort 2 for Grade 1* (NEPS Survey Paper No. 58). Bamberg, Germany: Leibniz Institute for Educational Trajectories, National Educational Panel Study. doi:10.5157/NEPS:SP58:1.0

NEPS Technical Report for Science:

Scaling Results of Starting Cohort 2 for Grade 1

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 in grade 1 of starting cohort 2. The scientific literacy test contained 25 items with different response formats representing different contexts as well as different areas of knowledge. The test was administered to 6,734 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 a good reliability and that all items 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. A limitation of the test was the lack of very difficult items. However, 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. Additionally, the design and results of the linking study for the competence scores in kindergarten and grade 1 are presented.

Key words: scientific literacy, 1st grade, linking kindergarten and grade 1, 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.....	7
4	Analyses.....	7
4.1	Missing responses	7
4.2	Scaling model	8
4.3	Checking the quality of the test	8
4.4	Software	10
5	Results	10
5.1	Descriptive statistics of the responses	10
5.2	Missing Responses.....	10
5.2.1	Missing responses per person.....	10
5.2.2	Missing responses per item.....	13
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.....	20
5.4.6	Unidimensionality of the test.....	21
6	Discussion	21
7	Data in the Scientific Use file.....	22
7.1	Naming conventions.....	22
7.2	Linking of competence scores	22
7.2.1	Samples	22
7.2.2	The design of the link study.....	23
7.2.3	Results	23
7.3	Scientific literacy scores	25

1 Introduction

Within the National Educational Panel Study (NEPS) different competences are measured coherently across the life span. Tests have been developed for different competence domains. These include, among other things, reading competence, mathematical competence, scientific literacy, information and communication literacy (computer literacy), metacognition, vocabulary, and domain-general cognitive functioning. An overview of the competences measured in the NEPS is given by Weinert et al. (2011) and by Fuß, Gnambs, 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 in grade 1 of starting cohort 2. First, the main concepts of the scientific literacy test are introduced. Then, the scientific literacy data of starting cohort 2 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 one dimensional 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 matter, system, development and interaction. KAS is divided into the process-related components 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.

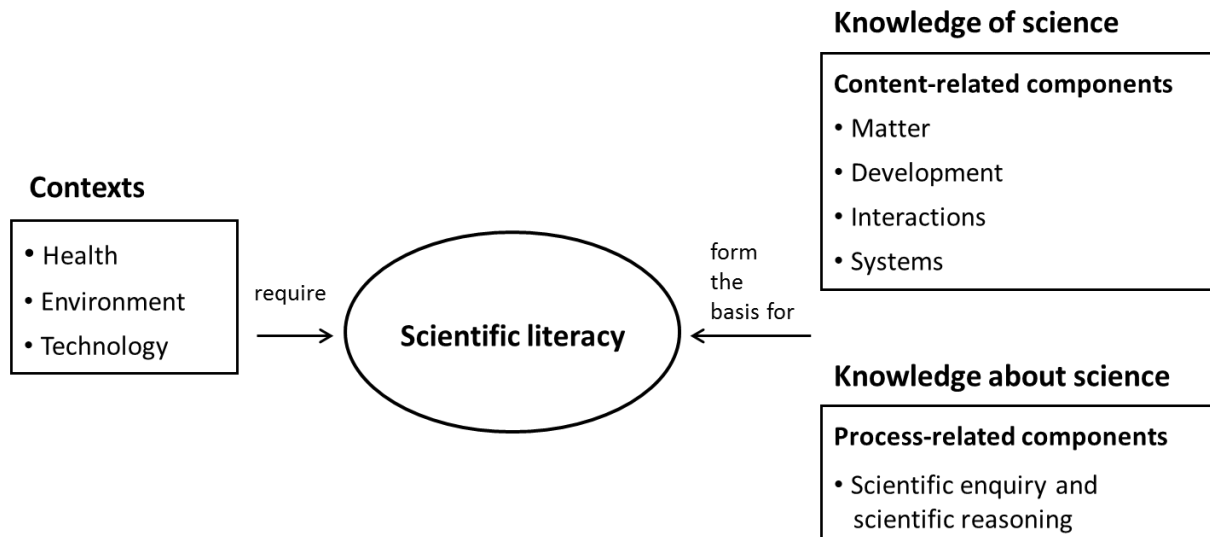


Figure 1. Assessment framework for scientific literacy (Hahn et al., 2013).

In the scientific literacy test for grade 1 of starting cohort 2 (Kindergarten) 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 including scientific literacy and, mathematical competences as well as procedural metacognition, receptive vocabulary, receptive grammatical competencies, and metacognition.

The competence tests for these domains were administered on two consecutive days. On one testing day, the children's mathematical competence and scientific literacy were assessed; the other competence domains were assessed on the other testing day. In order to control for possible effects of test position within testing days and order effects of the testing days, the two domains as well as the testing days were rotated. For this purpose, the sample was split into four groups receiving the tests in different orders and days. Assignment to the test booklets was random. Therefore, one testing group first completed the science test, followed by the mathematics test (including procedural metacognition), while the other group completed the two tests in the opposite order. Moreover, one group started with these two tests on the first testing day, the other group started with receptive vocabulary, receptive grammatical competencies, and metacognition (or in the order: metacognition, receptive vocabulary, receptive grammatical competencies) followed by either the science test and the mathematics test (including the procedural metacognition) or, in the opposite order, the mathematics test and the science test on the second testing day (see Table 1). 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 30 minutes.

Table 1

Design of the study

Test day	Rotation 1	Rotation 2	Rotation 3	Rotation 4
1 st	Math Science	Science Math	VOC GRA MC	MC VOC GRA
2 nd	VOC GRA MC	MC VOC GRA	Science Math	Math Science

Note. Math = mathematical competence, Science = Scientific literacy, VOC = vocabulary, GRA = grammatical competencies, MC = declarative metacognition.

A special challenge of this test was to take into account that the reading competences of this age group are very heterogeneous. Regarding the status of the early readers, all items – including the response options – were read out to the children by a test instructor. There were up to 14 children bundled in one test session. As a consequence, it was up to the test instructors to keep the time limits for the whole group in mind.

The allocation of the 25 items to the content areas (KOS and KAS) is summarized in Table 2. Table 3 shows how the items cover the different contexts of the scientific literacy framework (Hahn et al., 2013), whereas Table 4 gives an overview of the response formats.

Table 2

Classification of Items into Knowledge Domains

Knowledge domains	Number of Items
Knowledge of Science (KOS)	15
Knowledge about Science (KAS)	10
Total number of items	25

Table 3

Number of Items by Different Contexts

Context	Number of Items
Health	5
Environment	11
Technology	9
Total number of items	25

Table 4

Number of Items by Response Formats

Response format	Number of Items
Simple Multiple-Choice	22
Complex Multiple-Choice (True false items)	3
Total number of items	25

3.2 Sample

A total of 6,734 individuals received the scientific literacy test. For 239 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 (see Pohl & Carstensen, 2012). Thus, the analyses presented in this paper are based on a sample of 6,495 individuals (48.9% 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.

Invalid responses occurred, for example, when two response options were selected in simple MC items where only one was required, or when numbers or letters that were not within the range of valid responses were given as a response, or when less than four answers were given

in a CMC item (which consists of four subtasks). Omitted items occurred when test takers skipped some items. Due to time limits, not all persons finished the test within the given time. All missing responses after the last valid response given were coded as not-reached. As CMC items are aggregated from several subtasks, different kinds of missing responses or a mixture of valid and missing responses may be found in these items. A CMC item was coded as missing if at least one subtask contained a missing response. When one subtask contained a missing response, the CMC item was coded as missing. If just one kind of missing response occurred, the item was coded according to the corresponding missing response. If the subtasks contained different kinds of missing responses, the item was labeled as a not-determinable missing response.

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 gave an indication of how well the persons were coping with the test. We then looked at the occurrence of missing responses per item in order 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 in order 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 all of the four 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. In order 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 (see 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. Overall judgment of the fit of an item was based on all fit indicators.

The 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 socio-economic 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 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 (1984) Q3. Because in 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 25 items (including all subtasks for the polytomous items) were included in the analyses.

5.1 Descriptive statistics of the responses

In order 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 15.6% to 86.4% for the MC items. For the CMC items, the percentage of persons who correctly answered all subtasks varied between 6.7% and 28.5%.

5.2 Missing Responses

5.2.1 Missing responses per person

Figure 2 shows the number of invalid responses per person. Overall, there were very few invalid responses. 87.1% of the respondents did not have any invalid response at all; overall, about 3.9% of the respondents had more than one invalid response.

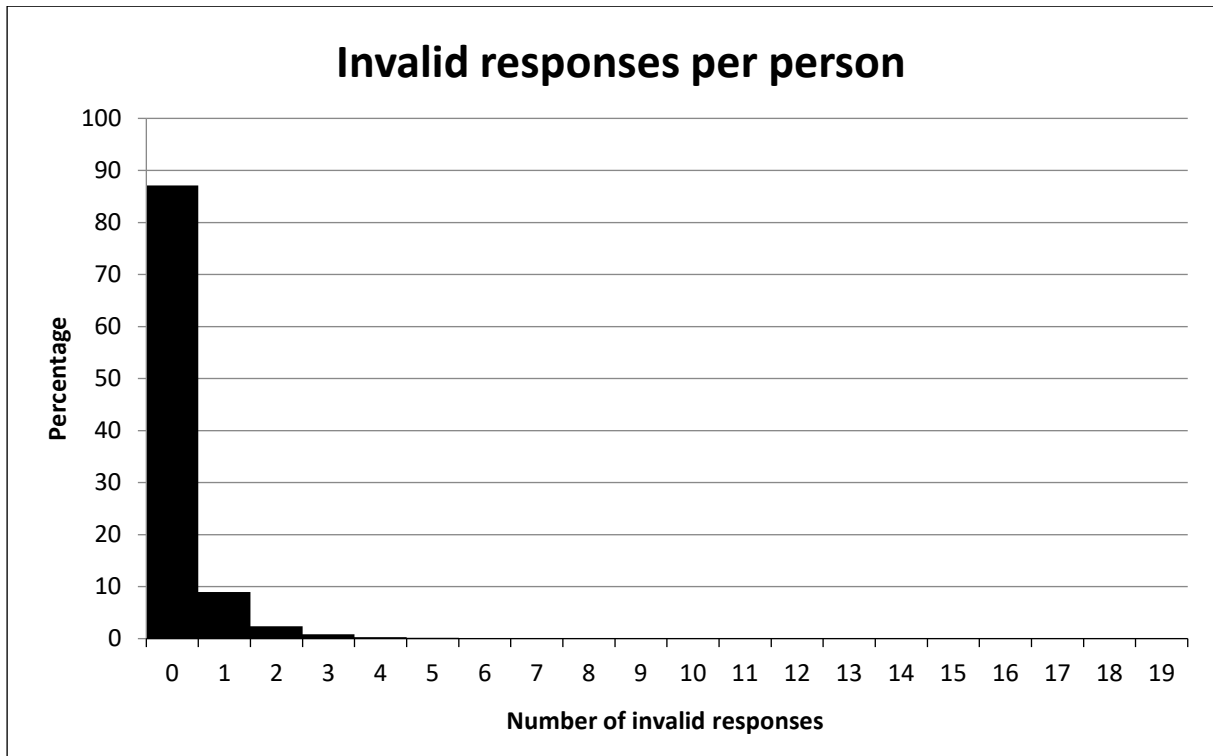


Figure 2. Number of invalid responses per person.

Missing responses may also occur when respondents omit items. As illustrated in Figure 3 most respondents, 64.9%, did not skip any item, and less than 6.0% omitted more than three items.

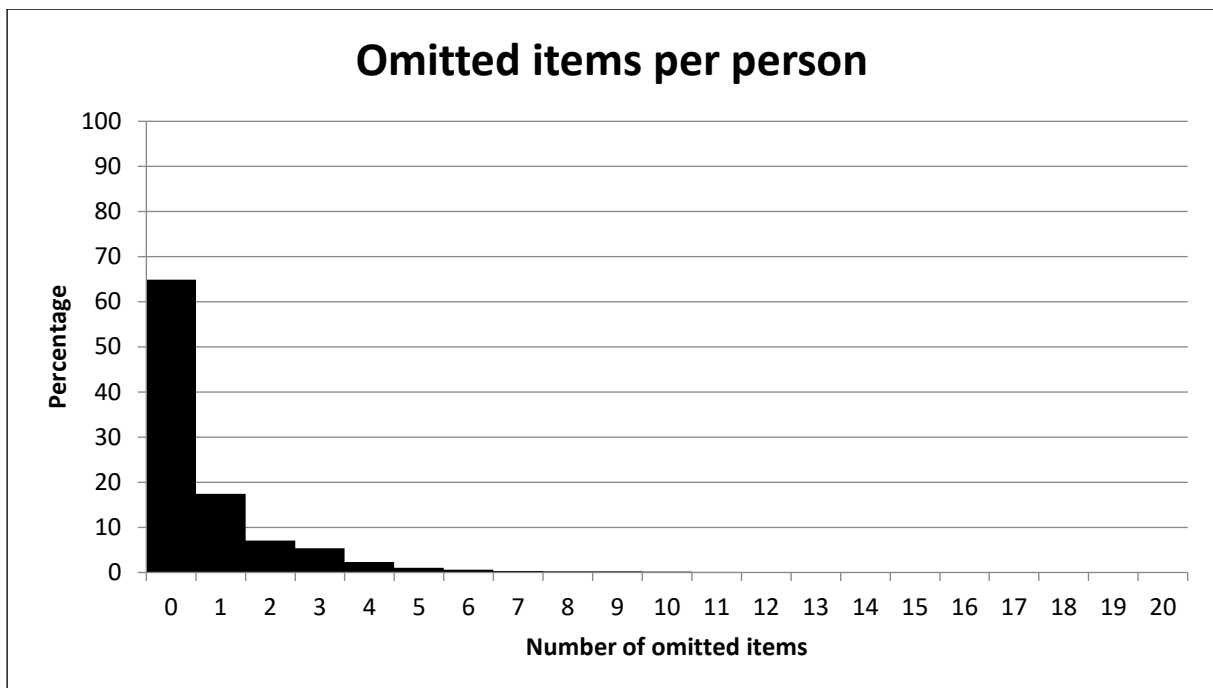


Figure 3. Number of omitted responses per person.

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 also very low, about 87.0% of the respondents were able to finish the test within the allocated time limit (Figure 4). Less than 1.0% did not finish more than half of the items.

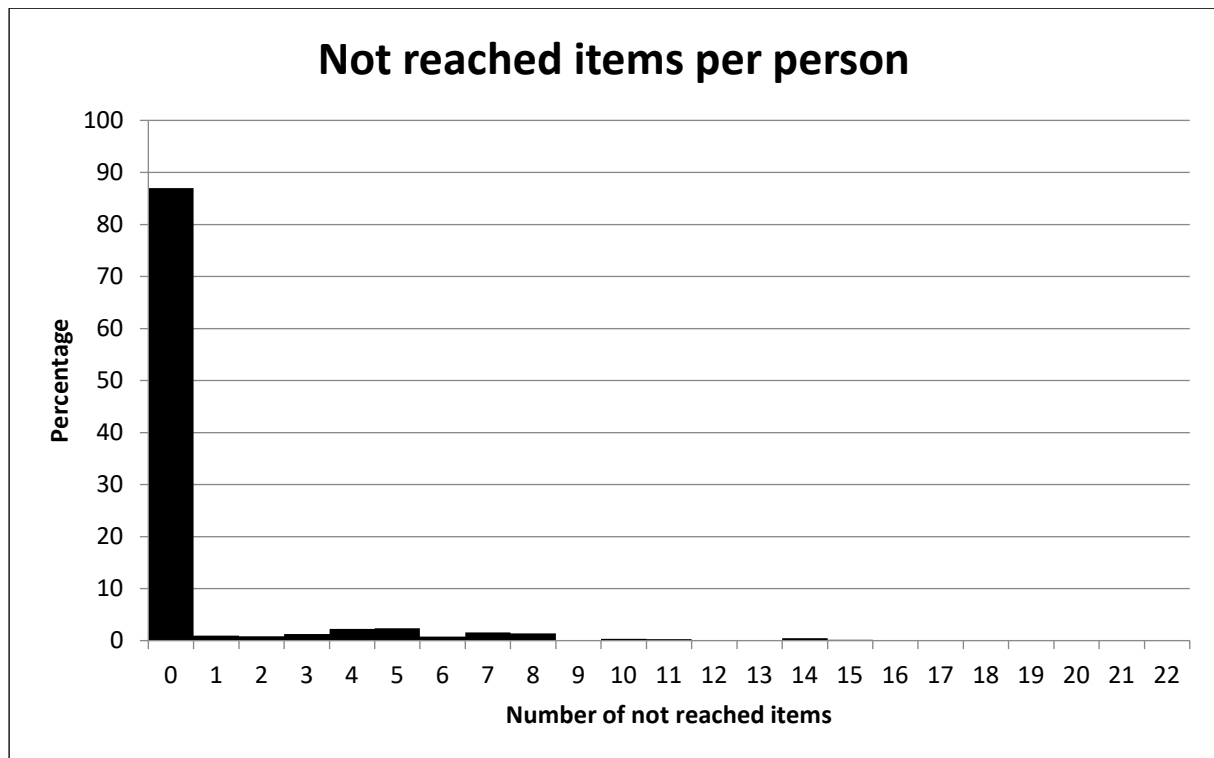


Figure 4. Number of not reached items per person.

The total number of missing responses, aggregated over invalid, omitted and not-reached missing responses, is illustrated in Figure 5. 51.1% of the students answered all questions and, consequently, had no missing responses. Only 1.3% of the students had missing responses on more than half of the items. Hence, the amount of missing responses per person can be classified as small.

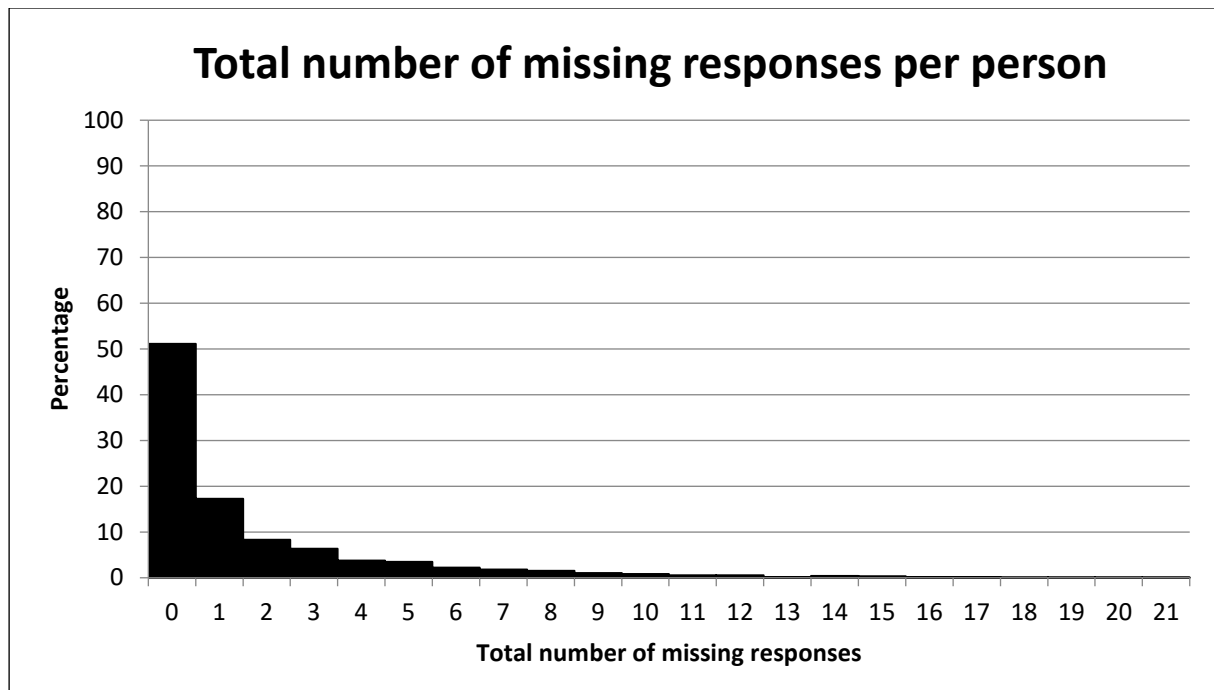


Figure 5. Total number of missing responses per person.

5.2.2 Missing responses per item

Table 5 shows the number of valid responses for each item as well as the percentage of missing responses. Overall, omission rates were rather low, varying across items between 0.0% and 12.5%. There were only two items with an omission rate exceeding 10.0% (scg1652s_c and scg1011s_c). The number of missing responses was uncorrelated ($r = .13$, $p = .543$) with the difficulty of the item. This result indicates that the test takers did not omit items that are more difficult. Generally, the percentage of invalid responses per item was rather low with the maximum rate being 3.6% (item scg10820_c). The relative frequency of not reached items increased towards the end of the test. Eventually, 13.0% 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 1.0% and 14.2%.

Table 5

Valid Responses and Missing Values

Item	Position in the test	Number of valid responses	Not reached items (%)	Omitted items (%)	Invalid responses (%)
scg10820_c	1	6078	0.0	2.8	3.6
scg10840_c	2	5925	0.0	7.5	1.3
scg11510_c	3	6093	0.0	5.7	0.5
scg10650_c	4	6436	0.0	0.7	0.3
scg16510_c	5	6170	0.0	4.6	0.4
scg1652s_c	6	5575	0.0	12.5	1.7
scg16110_c	7	6400	0.0	0.9	0.6
scg1091s_c	8	5817	0.0	9.3	1.1
scg10920_c	9	6124	0.0	3.5	2.2
scg1011s_c	10	5758	0.1	10.1	1.2
scg10120_c	11	6274	0.2	1.8	1.4
scg11210_c	12	6348	0.7	1.0	0.6
scg11110_c	13	6307	0.7	1.5	0.7
scg11130_c	14	6281	0.9	1.8	0.6
scg16530_c	15	6312	1.1	1.1	0.6
scg16020_c	16	6298	1.5	0.8	0.7
scg16030_c	17	6286	1.6	1.0	0.6
scg11610_c	18	6215	3.0	0.9	0.4
scg11710_c	19	5882	4.6	4.3	0.6
scg10310_c	20	6057	5.3	1.0	0.5
scg10520_c	21	5831	7.7	2.2	0.3
scg16310_c	22	5736	10.0	1.5	0.2
scg16220_c	23	5703	11.2	0.7	0.3
scg11440_c	24	5643	12.1	0.8	0.3
scg10410_c	25	5629	13.0	0.0	0.3

Table 6

Item parameters

Item	Percentage correct	Difficulty/location parameter	SE (difficulty/location parameter)	WMNSQ	t-value for WMNSQ	Pt.-bis. Corr. of correct response	Discrimination (2PL)	Yens Q3
scg10820_c	71.0	-1.000	0.030	0.98	-1.4	.40	0.89	.04
scg10840_c	53.2	-0.141	0.028	0.97	-3.8	.44	0.92	.09
scg11510_c	31.7	0.870	0.029	1.03	2.1	.33	0.60	.06
scg10650_c	86.4	-2.064	0.038	0.96	-1.5	.36	1.18	.09
scg16510_c	37.4	0.579	0.028	1.05	5.1	.31	0.53	.06
scg1652s_c	n.a.	0.230	0.021	1.10	6.5	.40	1.02	.12
scg16110_c	66.2	-0.763	0.028	1.02	1.7	.37	0.72	.09
scg1091s_c	n.a.	-0.601	0.028	0.99	-1.1	.36	1.41	.07
scg10920_c	64.7	-0.684	0.028	1.03	3.0	.34	0.64	.06
scg1011s_c	n.a.	0.017	0.028	0.98	-1.0	.22	1.78	.12
scg10120_c	54.2	-0.184	0.027	0.97	-3.0	.45	0.97	.05
scg11210_c	66.3	-0.766	0.028	0.99	-1.1	.40	0.84	.04
scg11110_c	76.7	-1.343	0.031	0.97	-1.8	.39	0.98	.06
scg11130_c	65.0	-0.697	0.028	0.91	-8.6	.52	1.44	.07
scg16530_c	22.5	1.396	0.032	1.06	3.7	.22	0.37	.06
scg16020_c	15.6	1.895	0.036	1.01	0.2	.29	0.72	.05
scg16030_c	54.6	-0.209	0.027	0.92	-9.3	.52	1.29	.09
scg11610_c	43.4	0.304	0.027	1.02	1.8	.39	0.75	.06
scg11710_c	26.3	1.164	0.031	1.05	3.2	.27	0.47	.04
scg10310_c	43.6	0.296	0.028	1.02	1.9	.36	0.66	.06
scg10520_c	69.8	-0.944	0.030	0.96	-3.0	.45	1.09	.07
scg16310_c	69.9	-0.957	0.031	0.93	-5.5	.49	1.34	.07
scg16220_c	56.5	-0.298	0.029	0.98	-2.3	.44	0.95	.06
scg11440_c	29.1	1.009	0.031	1.02	1.3	.33	0.64	.06
scg10410_c	58.1	-0.375	0.029	1.04	3.8	.34	0.57	.04

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 and MA item scores. These are denoted by n.a. For the dichotomous items, the item-total correlation corresponds to the point-biserial correlation between the correct response and the total score; for polytomous items it corresponds to the product-moment correlation between the corresponding categories and the total score (discrimination value as computed in ConQuest).

5.3 Parameter estimates

5.3.1 Item parameters

Column 2 in Table 6 shows the percentage of correct responses in relation to all valid responses for each item. Note that since there was a non-negligible amount of missing responses, this probability cannot be interpreted as an index for item difficulty. The percentage of correct responses within items varied between 6.7% and 86.4% with an average of 48.6% ($SD = 21.6$) correct responses.

The estimated item difficulties (for dichotomous items, MC items) and location parameters (for polytomous variables, CMC items) are given in Table 6. The step parameters (for polytomous variables) are depicted in Table 7. For one of the CMC items (scg1091s_c) the two lowest categories were collapsed, thus, these items were scaled using a scoring of 0, 0.5, 1, and 1.5. 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 -2.06 (scg10650_c) and 1.90 (scg16020_c). In total, the estimated item difficulties had a mean of -0.13 ($SD = 0.92$). Due to the large sample size, the standard errors of the estimated item difficulties were very small ($SE(\beta) \leq 0.038$). Overall, the item difficulties were rather low; the test did not include items with a high difficulty (above 2 logits).

Table 7

Step parameters for the CMC items

Item	Step 1 (SE)	Step 2 (SE)	Step 2 (SE)	Step 3
scg1652s_c	-0.197 (0.027)	0.352 (0.029)	-0.448 (0.033)	0.293
scg1091s_c	-0.865 (0.027)	0.502 (0.030)	0.363	
scg1011s_c	-1.618 (0.039)	-0.359 (0.027)	0.447 (0.031)	1.530

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, 2012a). 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 (2012a).

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, 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 0.637, indicating a somewhat limited variability between subjects. The reliability of the test (EAP/PV reliability = .726; WLE reliability = .725) was acceptable. Although the items

covered a wide range of the ability distribution, there were no items covering the lower and upper peripheral ability areas. As a consequence, person ability in medium ability regions will be measured relative precisely, whereas lower and higher ability estimates will have larger standard errors of measurement.

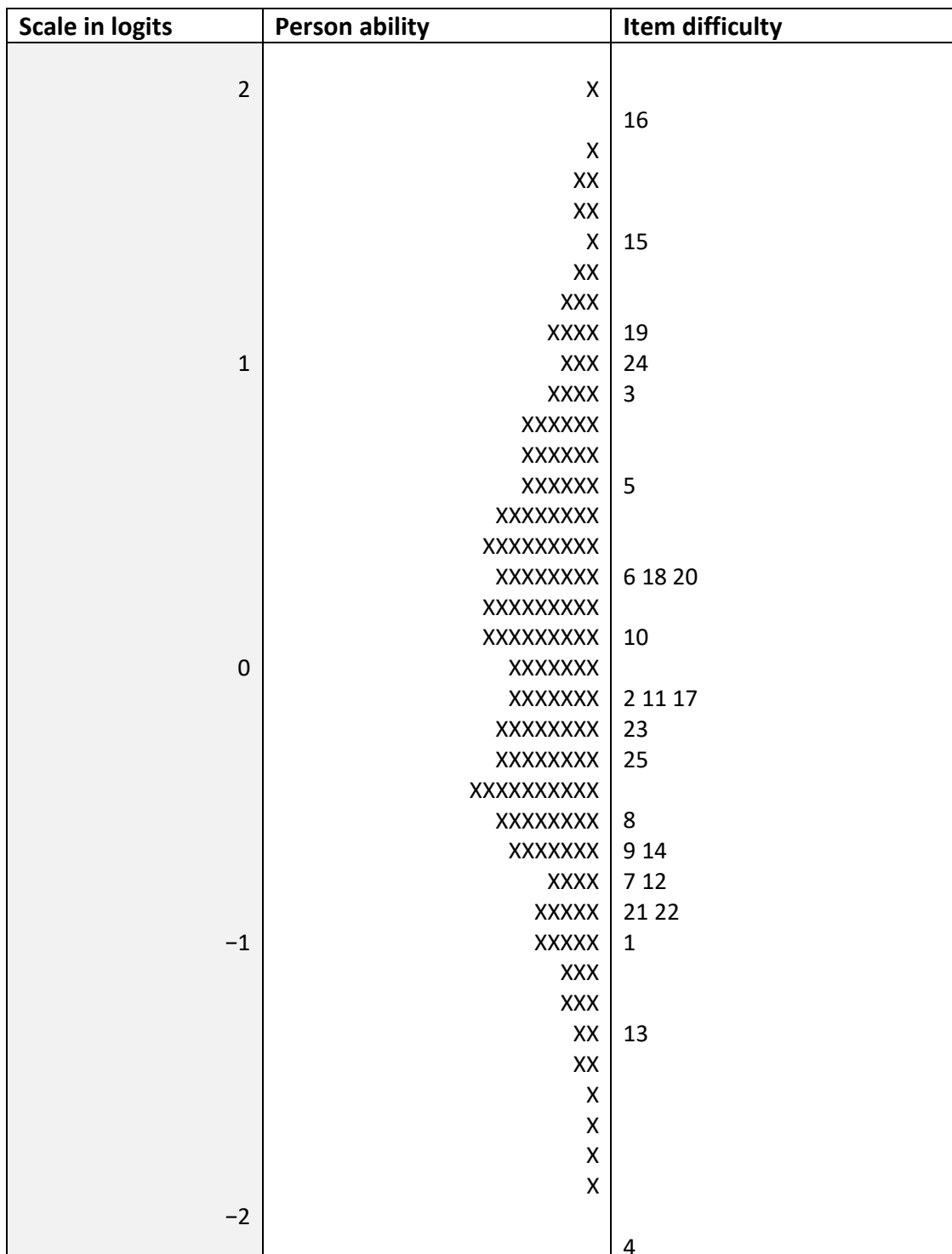


Figure 6. Test targeting. The distribution of person ability in the sample is depicted on the left side of the graph. Each 'X' represents 39.3 cases. The difficulty of the items is depicted on the right side of the graph. Each number represents an item (see Table 5).

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 34 items. The percentage of a correct response ranged from 15.6% to 87.6% across all items (*Mdn* = 53.7%). Thus, the number of correct and incorrect responses was reasonably large. All subtasks of the CMC items showed a satisfactory item fit. WMNSQ ranged from 0.92 to 1.09, the respective *t*-value from -9.0 to 8.0, and there were no noticeable deviations of the empirical 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. For three items there was one distractor with a point-biserial correlation with the total scores over zero: scg10120 (.00), scg16530 (.13) and scg11440 (.07). All of the other items only had distractors with a point-biserial correlation with the total scores below zero. Besides these three deviations, the results indicate that the distractors worked well.

5.4.3 Item fit

The evaluation of the item fit was performed on the basis of 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 6). Values of the WMNSQ ranged from 0.91 (item scg11130_c) to 1.10 (item scg1652s_c). Only one item exhibited a *t*-value of the WMNSQ greater than 6 (item scg1652s_c). Thus, there was no indication of severe item over- or underfit. Point-biserial correlations between the item scores and the total scores ranged from .22 (items scg1011s_c and scg16530_c) to .52 (items scg11130_c and scg16030_c) 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 rotation (test order and test day), 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). Table 8 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.

Table 8

Differential item functioning (absolute differences between difficulties)

Item	Rotation	Rotation	Gender	Books			Migration status		
	Testlet A vs. Testlet B	First vs. Second	Male vs. female	<100 vs. >100	<100 vs. missing	>100 vs. missing	Without vs. With	Without vs. Missing	With vs. Missing
scg10820_c	-0.204	-0.121	-0.170	-0.006	-0.052	-0.054	-0.056	-0.038	0.016
scg10840_c	0.024	-0.004	-0.006	-0.012	0.046	0.050	-0.058	0.042	0.098
scg11510_c	-0.058	0.004	0.106	-0.046	0.100	0.140	0.132	0.002	-0.134
scg10650_c	-0.114	-0.048	0.014	0.116	-0.168	-0.290	-0.348	-0.196	0.152
scg16510_c	0.134	-0.061	-0.164	-0.258	0.206	0.458	0.298	0.316	0.014
scg1652s_c	-0.170	0.002	-0.184	-0.202	0.032	0.242	0.258	0.182	-0.076
scg16110_c	-0.264	0.244	-0.158	-0.192	0.004	0.188	0.156	0.164	0.006
scg1091s_c	-0.268	-0.002	0.202	0.186	0.076	-0.106	-0.150	-0.058	0.102
scg10920_c	0.260	-0.035	0.002	-0.236	-0.190	0.040	0.178	-0.004	-0.186
scg1011s_c	-0.218	0.089	-0.130	0.018	-0.056	-0.078	-0.056	-0.114	-0.060
scg10120_c	0.128	-0.087	-0.022	-0.052	0.002	0.048	0.070	0.038	-0.036
scg11210_c	0.080	0.019	0.084	0.078	0.054	-0.030	-0.026	-0.034	-0.012
scg11110_c	-0.026	-0.042	0.530	0.182	0.022	-0.168	0.008	-0.054	-0.064
scg11130_c	0.130	-0.079	-0.030	0.368	-0.024	-0.398	-0.418	-0.210	0.208
scg16530_c	0.086	0.023	-0.222	-0.294	0.086	0.372	0.136	0.184	0.044
scg16020_c	-0.048	0.048	-0.002	-0.046	-0.106	-0.068	-0.226	0.016	0.240
scg16030_c	-0.082	0.204	-0.144	0.174	-0.034	-0.214	-0.096	-0.076	0.016
scg11610_c	0.054	-0.069	0.248	0.016	0.030	0.006	0.028	-0.038	-0.070
scg11710_c	-0.074	-0.009	0.298	-0.140	0.058	0.190	0.256	0.180	-0.080
scg10310_c	0.156	0.008	-0.212	0.054	0.078	0.016	-0.040	-0.014	0.022
scg10520_c	0.150	-0.077	-0.044	0.296	0.048	-0.254	-0.324	-0.146	0.176
scg16310_c	0.218	-0.080	0.058	0.340	0.052	-0.296	-0.078	-0.236	-0.158
scg16220_c	0.008	0.091	0.062	0.010	-0.100	-0.118	-0.164	-0.132	0.030
scg11440_c	0.006	-0.042	0.174	0.046	-0.032	-0.084	-0.174	-0.034	0.140
scg10410_c	0.234	-0.040	0.008	-0.100	-0.210	-0.116	0.164	-0.018	-0.186

Rotation – Test order

The scientific literacy test was administered in two different positions (see section 3.1 for the design of the study). A total of 3,232 (49.8%) of the test takers received the scientific literacy test first and then the mathematical test, while 3,261 (50.2%) received the mathematical literacy test before completing the scientific literacy test. The students were randomly assigned to either of the two design groups. For two test takers no information was available regarding the test position; therefore, they were excluded from this analysis. Differential item functioning of the position of the test may, for example, occur if the different certain parts or items of the test are more or less tiring for the participants. There was a small difference between the first test position and the second test position (main effect = -0.060 logits, Cohen's $d = -0.075$), indicating a lower difficulty for test takers with the first test position. Also, the highest difference in difficulties between the two groups is -0.268 logits.

Rotation – Test day

The scientific literacy test was administered on two test days (see section 3.1 for the design of the study). A total of 3,242 (49.9%) test takers were tested on the first test day, whereas 3,251 (50.1%) were tested on the second test day. The students were randomly assigned to either of the two test days. There was a no difference between these two groups (main effect = 0.014 logits, Cohen's $d = 0.018$). Also, the highest difference in difficulties between the two groups was 0.488 logits.

Gender

The sample included 3,178 (48.9%) female test takers and 3,317 (51.1%) male test takers. On average, female students had slightly higher scores in scientific literacy than male students (main effect = 0.034 logits, Cohen's $d = 0.043$). But there was no item with a considerable gender DIF. The highest difference in difficulties between the two groups was 0.530 logits.

Books

The number of books at home was used as a proxy for socio-economic status. There were 1,926 (29.7%) test takers with 0 to 100 books at home, 3,515 (54.1%) test takers with more than 100 books at home, and 1,054 (16.2%) test takers did not give a valid response. DIF was investigated using these three groups. There were considerable average differences between these three groups. Participants with 100 or less books at home on average showed lower scientific literacy scores than participants with more than 100 books (main effect = -0.534 logits, Cohen's $d = -0.727$). Participants without a valid response on the variable 'books at home' performed lower than participants with up to 100 (main effect = -0.128 logits, Cohen's $d = -0.168$) and lower than participants with more than 100 books at home, respectively (main effect = -0.664 logits, Cohen's $d = -0.899$). There was no considerable DIF comparing participants with many or fewer books (highest DIF = 0.368). Comparing the group without valid responses to the two groups with valid responses, DIF occurred up to 0.458 logits.

Migration background

There were 3,309 (50.9%) participants without a migration background and 1,137 (17.5%) participants with a migration background (for 0.8% students neither their mother, father or themselves were born in Germany, for 6.3% only the participants were born in Germany and both of their parents were born abroad, and for 10.5% of the participants only one of their

parents was born abroad). A total of 2,049 (31.5%) 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 = 0.450 logits, Cohen's $d = 0.600$) and also higher scores than students with an unknown background on migration (main effect = 0.360 logits, Cohen's $d = 0.463$). Furthermore, students with a migration background scored lower than those with an unknown background on migration (main effect = -0.092 logits, Cohen's $d = -0.115$). There was no considerable DIF comparing participants with and without a migration background (highest DIF = -0.418). Comparing the group without valid responses to the two groups with valid responses, DIF occurred up to 0.316 logits.

Besides investigating DIF for each single item, an overall test for DIF was performed by comparing models which allow for DIF with those that allow only for main effects. In Table 9, the models including only 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 for all four DIF variables. The Bayesian information criterion (BIC, Schwarz, 1978) 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 three of the four DIF variables (rotation, books and migration background).

Table 9

Comparison of models with and without DIF

DIF variable	Model	Deviance	N	Number of parameters	AIC	BIC
Rotation	main effect	205427.59	6493	35	205497.59	205734.84
Test order	DIF	205246.95	6493	60	205366.95	205773.66
Rotation	main effect	205434.17	6493	35	205504.17	205741.42
Test day	DIF	205212.36	6493	60	205332.36	205739.07
Gender	main effect	205504.95	6495	35	205574.95	205812.21
	DIF	205279.37	6495	60	205399.37	205806.09
Books	main effect	171495.35	5441	35	171565.35	171796.41
	DIF	171313.51	5441	60	171433.51	171829.62
Migration background	main effect	140189.74	4446	35	140259.74	140483.73
	DIF	140050.45	4446	60	140170.45	140554.44

5.4.5 Rasch-homogeneity

An essential assumption of the Rasch (1980) model is that all item-discrimination parameters are equal. In order to test this assumption, a generalized partial credit model (GPCM; Muraki, 1982) that estimates discrimination parameters was fitted to the data. The estimated discriminations differed moderately among items (see Table 6), ranging from 0.37 (item scg16530_c) to 1.77 (item scg1011s_c). The average discrimination parameter fell at 0.91. Model fit indices suggested a better model fit of the GPCM (AIC = 204,514.27, BIC =

204,907.44) as compared to the PCM model (AIC = 205,575.17, BIC = 205,805.65). 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 two-dimensional model (BIC= 205,798.17, number of parameters = 36) fitted the data slightly better than a unidimensional model (BIC= 205,805.65, number of parameters = 34). As the correlation between the two dimensions was $r = .94$ 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 administered in grade 1 of starting cohort 2 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. However, the amount of missing responses was reasonably 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 2PL model or as a correlation of the item score with total score) were acceptable. Different variables were used for testing measurement invariance across various subgroups. No considerable DIF became evident for any of these 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 a slightly better model fit than

the unidimensional partial credit model. However, 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

There are 25 items in the data set that are either scored as dichotomous variables (MC or SCR 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. Manifest scale scores are provided in form of WLE estimates (scg1_sc1) including the respective standard error (scg1_sc2). 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 two lowest categories of one 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 and 1.5 (denoting the highest). The ConQuest Syntax for estimating the WLE scores from the items is provided in Appendix A. Students who did not take part in the test or those who did not have enough valid responses to estimate a scale score have a non-determinable missing value on the WLE score for scientific literacy.

7.2 Linking of competence scores

In starting cohort 2, the scientific literacy tests which were administered in kindergarten and grade 1 included different items that were constructed in such a way as to allow for an accurate measurement of scientific literacy within each age group. As a consequence, the competence scores derived in the different grades cannot be directly compared. Differences in observed scores would reflect differences in competences as well as differences in test difficulties. To place the different measurements onto a common scale and, thus, allow for the longitudinal comparison of competences across grades, I adopted the linking procedure described in Fischer, Rohm, Gnams, and Carstensen (2016). Following an anchor-group design, all items from the kindergarten and grade 1 scientific literacy test were administered in an independent link sample – including students from grade 1 that were not part of starting cohort 2 – within a single measurement occasion. These responses were used to link the two tests administered in starting cohort 2 across the two grades.

7.2.1 Samples

In starting cohort 2, a subsample of 534 students participated at both measurement occasions, in kindergarten and also in grade 1. Consequently, these students were used to link the two tests across both grades (see Fischer et al., 2016.). Moreover, an independent link sample of $N = 480$ students (51.0% female) from grade 1 received both tests within a single measurement occasion.

7.2.2 The design of the link study

Both tests administered in kindergarten and in grade 1 included 25 items (see above). The science tests were administered in a random order. Half of the sample received the grade 1 test before working on the kindergarten test, whereas the other half received the kindergarten test before the grade 1 test. No multi-matrix design regarding the selection and order of the items within a test was established. Thus, all test takers were given the science items in the same order.

7.2.3 Results

To examine whether the two tests administered in the link sample measured the same construct, we compared a one-dimensional model that specified a single latent factor for all items to a two-dimensional model that specified separate latent factors for the two tests. The information criteria favored the one-dimensional model, AIC = 25,146.69, BIC = 25,417.99, over the two-dimensional model, AIC = 25,150.77, BIC = 25,430.42. Moreover, an examination of the residual correlations for the one-dimensional model using the corrected Q_3 statistic (Yen, 1984) indicated a largely unidimensional scale—the average absolute residual correlation was $M = .00$ ($SD = .06$). This indicates that the scientific literacy tests administered in kindergarten and grade 1 were essentially unidimensional.

Items that are supposed to link two tests must exhibit measurement invariance; otherwise, they cannot be used for the linking procedure. Therefore, we tested whether the item parameters derived in the link sample showed a non-negligible shift in item difficulties as compared to the longitudinal subsample from the starting cohort. The differences in item difficulties between the link sample and starting cohort 2 and the respective tests for measurement invariance based on the Wald statistic (see Fischer et al., 2016) are summarized in Table 10.

Table 10

Differential Item Functioning Analyses between the Starting Cohort and the Link Sample

	Kindergarten				Grade 1			
	Item	$\Delta\sigma$	$SE_{\Delta\sigma}$	F	Item	$\Delta\sigma$	$SE_{\Delta\sigma}$	F
1.	sck10420_c	-0.16	0.23	0.53	scg10820_c	0.10	0.15	0.40
2.	sck10430_c	-0.35	0.17	4.42	scg10840_c	0.20	0.14	2.09
3.	sck16120_c	0.01	0.27	0.00	scg11510_c	0.12	0.15	0.64
4.	sck16130_c	0.58	0.21	7.48	scg10650_c	-3.87	0.20	382.65
5.	sck1102s_c	-0.44	0.14	9.40	scg16510_c	0.68	0.14	24.74
6.	sck11030_c	0.37	0.17	4.91	scg1652s_c	0.87	0.11	65.48
7.	sck1033s_c	1.58	0.21	54.83	scg16110_c	-0.11	0.14	0.61
8.	sck10210_c	2.11	1.02	4.29	scg1091s_c	0.35	0.14	6.02
9.	sck1023s_c	-0.44	0.19	5.52	scg10920_c	0.82	0.14	32.28
10.	sck11110_c	0.00	0.20	0.00	scg1011s_c	0.37	0.14	6.68
11.	sck11120_c	-0.47	0.15	10.52	scg10120_c	0.43	0.14	9.75
12.	sck16010_c	-0.52	0.18	8.03	scg11210_c	0.07	0.14	0.25
13.	sck16020_c	0.66	0.15	18.64	scg11110_c	0.25	0.16	2.50
14.	sck10510_c	0.24	0.25	0.86	scg11130_c	-2.71	0.18	230.30
15.	sck10530_c	-0.26	0.18	2.07	scg16530_c	0.15	0.16	0.86
16.	sck11610_c	0.07	0.38	0.04	scg16020_c	0.63	0.18	12.90
17.	sck1162s_c	-0.07	0.16	0.21	scg16030_c	-0.09	0.14	0.43
18.	sck10710_c	-1.26	0.21	35.52	scg11610_c	0.24	0.14	3.06
19.	sck10720_c	0.27	0.20	1.83	scg11710_c	0.34	0.15	4.87
20.	sck11310_c	-1.27	0.15	75.18	scg10310_c	0.21	0.14	2.44
21.	sck11330_c	-1.06	0.17	40.15	scg10520_c	0.43	0.15	8.32
22.	sck10910_c	-0.11	0.17	0.39	scg16310_c	0.35	0.15	5.32
23.	sck10920_c	0.13	0.33	0.16	scg16220_c	0.00	0.14	0.00
24.	sck16210_c	0.98	0.23	18.59	scg11440_c	-0.24	0.16	2.15
25.	sck16220_c	-0.60	0.14	18.62	scg10410_c	0.42	0.14	8.71

Note. $\Delta\sigma$ = Difference in item difficulty parameters between the longitudinal subsample in Kindergarten and Grade 1 and the link sample (positive values indicate easier items in the link sample); $SE_{\Delta\sigma}$ = Pooled standard error; F = Test statistic for the minimum effects hypothesis test (see Fischer et al., 2016). The critical value for the minimum effects hypothesis test using an α of .05 is $F_{.0154}(1, 1,014) = 31.87$. A non-significant test indicates measurement invariance.

Thus, measurement invariance for kindergarten showed four items (sck1033s_c, sck10710_c, sck11310_c, and sck11330_c) with F -statistics exceeding the critical value of $F_{.0154}(1, 1,014) = 31.87$. In grade 1 also four items (scg10650_c, scg1652s_c, scg10920_c, and scg11130_c) exceeded the critical value. These items were therefore excluded from the estimation of the correction term.

Moreover, analyses of differential item functioning between the link sample and starting cohort 2 showed in kindergarten for 19 of the remaining 21 items of the test no DIF greater than 0.40 (difference in logits: $Min = -0.64$, $Max = 1.05$). But, two items (sck10210_c and sck16210_c) showed a DIF greater than 0.40. These items were also excluded from the estimation of the correction term. For grade 1 (difference in logits: $Min = -1.96$, $Max = 0.45$) there was no additional item with a DIF greater than 0.40, leaving 21 items for estimating the correction term.

Furthermore, there was one item in kindergarten (sck1023s_c) where categories were collapsed differently between the linking data and the main study data. This item was excluded from the estimation of the correction term.

The scientific literacy tests administered in the two grades were linked using the “mean/mean” method for the anchor-group design (see Fischer et al., 2016). The correction term was calculated as $c = 1.445$ (with a link error of 0.10). This correction term was subsequently added to each difficulty parameter estimated in grade 1 (see Table 6) to derive the linked item parameters.

7.3 Scientific literacy scores

In the SUF manifest science literacy scores are provided in the form of two different WLEs, “scg1_sc1” and “scg1_sc1u”, including their respective standard error, “scg1_sc2” and “scg1_sc2u”.

For “scg1_sc1u”, person abilities were estimated using the linked item difficulty parameters. Normally, the estimated WLE scores must be corrected for differences in the test position. Subsequently, in grade 1, the science test was either presented as the first or the second test within the test battery, whereas in kindergarten, the science test was always presented first within the test battery. To correct for differences in the test position, we added the main effect ($=0.024$) related to the test position to the WLE scores of respondents that received the science test after working on another test. As a result the WLE scores provided in “scg1_sc1u” can be used for longitudinal comparisons between kindergarten and grade 1. The resulting differences in WLE scores can be interpreted as development trajectories across measurement points. In contrast, the WLE scores in “scg1_sc1” are not linked to the underlying reference scale of kindergarten. But, it is also corrected for the position of the scientific literacy test within the booklet and can be used, if the research interest lies on cross-sectional issues. The ConQuest Syntax for estimating the WLE is provided in Appendix A. For persons 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).

References

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, *19*, 716–722.
- Carstensen, C. H. (2013). Linking PISA competencies over three cycles – Results from Germany. In M. Prenzel, M. Kobarg, K. Schöps, & S. Rönnebeck (Eds.). *Research Outcomes of the PISA Research Conference 2009* (pp. 199–214). New York, NY: Springer.
- Fuß, D., Gnams, 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., Dalehefte, I. M., & 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.
- Mislevy, R. J. (1991). Randomization-based inference about latent variables from complex samples. *Psychometrika*, *56*, 177–196.
- 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.
- Pohl, S., Haberkorn, K., Hardt, K., & Wiegand, E. (2012). Technical Report of Reading – Scaling Results of Starting Cohort 3 in Fifth Grade (NEPS Working Paper No. 15). Bamberg, Germany: Otto-Friedrich-Universität, Nationales Bildungspanel.
- Prenzel, M., Schöps, K., Rönnebeck, S., Senkbeil, M., Walter, O., Carstensen, C. H., & Hamann, M. (2007). Naturwissenschaftliche Kompetenz im internationalen Vergleich. In PISA-Konsortium Deutschland (Hrsg.), *PISA 2006 – Die Ergebnisse der dritten internationalen Vergleichsstudie* (S. 63–105). Münster, Germany: Waxmann.
- Rasch, G. (1980). *Probabilistic models for some intelligence and attainment tests*. Chicago, IL: University of Chicago Press.
- Schwarz, G. E. (1978). Estimating the dimension of a model. *Annals of Statistics*, *6*, 461–464.
- von Davier, M. (2005). *A general diagnostic model applied to language testing data* (ETS Research Rep. No. RR-05-16). Princeton, NJ: ETS.
- Warm, T. A. (1989). Weighted likelihood estimation of ability in item response theory. *Psychometrika*, *54*, 427–450.

- Weinert, S., Artelt, C., Prenzel, M., Senkbeil, M., Ehmke, T., & Carstensen C. H. (2011) Development of Competencies Across the Life Span. In H. P. Blossfeld, H. G. Roßbach & J. v. Maurice & (Eds.). *Education as a Lifelong Process: The German National Educational Panel Study (NEPS)*. (*Zeitschrift für Erziehungswissenschaft, Sonderheft 14* . Wiesbaden, Germany: VS Verlag für Sozialwissenschaften.
- Wu, M. L., Adams, R. J., & Wilson, M. R. (1997). *ACER Conquest: Generalised item response modelling software*. Melbourne, Australia: ACER Press.

Appendix

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

Title G1 Science analysis, Partial Credit Model;

data filename.dat;

format id 1–7 responses 8–32;

labels << filename_with_labels.txt;

codes 0,1,2,3,4;

score (0,1) (0,1) !item (1–5,7,9,11–25);

score (0,1,2,3,4) (0,0.5,1,1.5,2) !item (6,10);

score (0,1,2,3) (0,0.5,1,1.5) !item (8);

set constraint=cases;

model item + item*step;

estimate;

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 to contexts

Variable name	Position in the test	Component	Context
scg10820_c	1	KOS	Health
scg10840_c	2	KAS	Health
scg11510_c	3	KOS	Technology
scg10650_c	4	KAS	Health
scg16510_c	5	KAS	Environment
scg1652s_c	6	KAS	Environment
scg16110_c	7	KAS	Environment
scg1091s_c	8	KOS	Technology
scg10920_c	9	KOS	Technology
scg1011s_c	10	KOS	Health
scg10120_c	11	KOS	Health
scg11210_c	12	KOS	Environment
scg11110_c	13	KOS	Environment
scg11130_c	14	KOS	Environment
scg16530_c	15	KAS	Technology
scg16020_c	16	KAS	Environment
scg16030_c	17	KAS	Environment
scg11610_c	18	KOS	Environment
scg11710_c	19	KOS	Technology
scg10310_c	20	KOS	Technology
scg10520_c	21	KOS	Environment
scg16310_c	22	KAS	Technology
scg16220_c	23	KAS	Technology
scg11440_c	24	KOS	Technology
scg10410_c	25	KOS	Environment

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