

This is an electronic reprint of the original article. This reprint may differ from the original in pagination and typographic detail.

Validation of an Early Numeracy Screener for First Graders

Lopez-Pedersen, Anita; Mononen, Riikka; Korhonen, Johan; Aunio, Pirjo; Melby-Lervåg, Monica

Published in:
Scandinavian Journal of Educational Research

DOI:
[10.1080/00313831.2019.1705901](https://doi.org/10.1080/00313831.2019.1705901)

Published: 01/01/2020

Document Version
Accepted author manuscript

Document License
CC BY

[Link to publication](#)

Please cite the original version:

Lopez-Pedersen, A., Mononen, R., Korhonen, J., Aunio, P., & Melby-Lervåg, M. (2020). Validation of an Early Numeracy Screener for First Graders. *Scandinavian Journal of Educational Research*, 1–22.
<https://doi.org/10.1080/00313831.2019.1705901>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

This paper is not the copy of record and may not exactly replicate the final, authoritative version of the article. Please do not copy or cite without authors permission. The final article is available, via <https://doi.org/10.1080/00313831.2019.1705901>

Anita Lopez-Pedersen,
Department of special needs education
University of Oslo
anita.lopez-pedersen@isp.uio.no

Riikka Mononen,
Department of special needs education
University of Oslo
r.m.mononen@isp.uio.no

Johan Korhonen
Åbo Akademi University, Turku
johan.korhonen@abo.fi

Pirjo Aunio,
Helsinki University
Finland and Department of special needs education
University of Oslo
pirjo.aunio@helsinki.fi

Monica Melby-Lervåg (Corresponding author)
Department of special needs education
University of Oslo
monica.melby-lervag@isp.uio.no

Word count: 9631

Abstract

This study investigated the psychometric properties of the Early Numeracy Screener. The Early Numeracy Screener is a teacher administered, paper-and-pencil test measuring counting skills, numerical relational skills and basic arithmetic skills. Three hundred and sixty-six first graders took the Early Numeracy Screener in the beginning of the school year. Confirmatory factor analysis was conducted in order to examine whether the screening tool was identified as a one-factor model or a three-factor model. The confirmatory factor analysis found evidence for the three-factor model, establishing construct validity. Evidence for criterion-related validity was found in crosstabulation and correlation with the national test measuring overall mathematics performance taken towards the end of the school year. The Early Numeracy Screener may serve as an indicator of young children's performance in early numeracy. The brevity and ease of use of the Early Numeracy Screener makes it suitable for classroom instructional settings.

Keywords: early numeracy, reliability, screener, test, validity,

Validation of an Early Numeracy Screener for First Graders

Developing well-functioning early numeracy skills is a foundation for further mathematical skills and for qualification for employment in society (Duncan et al., 2007; Geary, Bailey, & Hoard, 2009; Geary, Hoard, Nugent, & Bailey, 2013; Ritchie & Bates, 2013; Trilling, Bernie, Fadel, & Charles, 2009). Mathematical skills develop in a cumulative fashion, with early skills forming the foundation for the acquisition of later skills (Aunio & Niemivirta, 2010; Jordan, Kaplan, Ramineni, & Locuniak, 2009; Purpura, Baroody, & Lonigan, 2013). The early childhood years serves perhaps the most important developmental years in one's life (McGuire, Kinzie, & Berch, 2012). Longitudinal studies show that early numeracy skills are important for the kind of learning trajectory the child has in primary school mathematics (Aunio & Niemivirta, 2010; Desoete et al., 2009; Dougherty, 2003; Gersten et al, 2015; Jordan, Glutting, & Ramineni, 2010; Jordan, Kaplan, Ramineni, & Locuniak, 2009; Krajewski & Schneider, 2009; Zhang et al., 2014). Studies have related children's mathematical achievement to specific aspects (e.g. counting skills) of their early numerical competencies (Hannola-Sormunen, Lehtinen, & Räsänen, 2015). Differences in early numeracy are displayed before the onsets of formal schooling (Berch, 2005). Children who perform poorest in early numeracy skills may have serious deficits in all early number skills (Salminen et al., 2018). For instance, verbal counting plays an important role as a predictor of arithmetic (Zhang et al., 2014), accordingly an important skill for identifying children at-risk for developing mathematical learning disabilities, might be counting. Therefore, to assist children in establishing these skills, it is important to identify children who struggle with numeracy skills at an early stage. Despite the need, there are surprisingly few well-validated screening tools for this purpose for non-English-speaking European countries. To ensure the validity of assessment tools to be used in different countries, even in mathematics, they preferably need to be validated for each language and country. To fill this gap, we present a study of the development and validation of an early

numeracy screening tool in Norwegian to detect first graders with challenges in early numeracy skills.

Why assessment of early numeracy skills?

Mathematical skills is a continuous variable that is normally distributed in the population, and cut-offs to establish normal versus disordered development will be arbitrary. However, it is common to assume that around 15–20% of children and adults experience difficulties in developing mathematical skills in such an extent that it hampers their school or work performance (Geary, 2011). Out of these, 5–7% have problems so severe that they are often diagnosed as having specific mathematical learning disabilities or dyscalculia (American Psychiatric Association, 2013; Butterworth, Varma, & Lurillard, 2011). Identifying and remediating the early numeracy that predict poor school-entry mathematical knowledge has the potential to substantially reduce these risks, and accordingly considerable resources have been devoted to these efforts in recent years (Clarke et al., 2016; Fuchs et al., 2013; Gersten et al., 2015; Jordan, Glutting, Dyson, Hassinger - Das & Irwin, 2012).

Core numerical skills model: What do we need to assess?

When it comes to developing appropriate targeted assessment tools for early numeracy skills, it has been suggested that the skills that need to be considered generally fall into three different domains: understanding numerical relations, counting skills, and basic arithmetic skills (Aunio & Räsänen, 2016; Jordan, Kaplan, Locuniak, & Ramineni, 2009; National Research Council [NRC], 2009; Purpura & Longian, 2013). Aunio & Räsänen, 2016 theorized a model of these crucial numerical factors for the development of mathematical skills among children aged five to eight years old. Their model was based on the results of longitudinal studies, and a series of analyses of standardized tests intended to measure the development of mathematical skills. Support for the content of these domains can also be

found in previous research describing early numeracy (Desoete, Ceuемans, De Weerdт, & Pieters, 2012; Gersten et al. 2012; Moeller, Pixner, Zuber, Kaufmann, & Nuerk, 2011).

Understanding numerical relations

If we more closely examine the three putative domains that constitute early numeracy skills, the first foundation is to understand numerical relations. Numerical relations refer to the understanding of the quantitative and non-quantitative relationships between the elements in the task (Aunio & Räsänen, 2016). Numerical relational skills serve as a prerequisite to basic arithmetic skills. The knowledge of basic arithmetic principles is often referred to as the understanding of part-whole relations in addition or subtraction tasks (Canobi, Reeve, & Pattison 2002; Wilkins, Baroody, & Tiilikainen 2001). Numerical relational skills include a set of subskills such as the ability to compare the magnitudes of numbers, to understand cardinal value, one-to-one correspondence and early mathematical-logical principles and to understand the meaning of the 10-base system (Aunio & Räsänen, 2016; Geary & vanMarle, 2016). Longitudinal studies clearly establish that numerical relational skills are a crucial part of early numeracy development (Aunio & Niemivirta, 2010; Desoete et al., 2009; Stock, Desoete, & Roeyers, 2009). Research has pointed out that children's understanding of numerical magnitudes predicts individual differences in mathematics achievement (e.g. De Smedt, Noël, Gilmore, & Ansari, 2013; De Smedt, Verschaffel, & Ghesquière, 2009; Halberda, Mazocco, & Feigenson, 2008; Holloway & Ansari, 2009; Schneider et al., 2016; Vanbinst, Ghesquière & De Smedt, 2015).

Counting skills

The second component in early numeracy skills is counting skills. Counting skills refers to the child's knowledge of number symbols, skills in moving within the sequence of the number words and enumeration (Aunio & Räsänen, 2016). Counting strategies are, perhaps not surprisingly, an imperative aspect of children's early numerical knowledge (

Aunio & Räsänen, 2016; Wright, Martland, & Stafford, 2006). Counting reinforces the child's understanding of the relationships between numbers (Baroody, 2006; Baroody, 2003; Baroody, Eiland, & Thompson, 2009). Counting also helps expand their quantitative knowledge to larger numbers (Baroody, 2006; Baroody, 2003; Baroody et al., 2009). Counting knowledge allows children to count on or up from addends to solve novel number combinations – a key arithmetic strategy in early elementary school (Geary, Hoard, Byrd-Craven, & DeSoto, 2004).

Basic arithmetic skills

The understanding of numerical relations and counting is a prerequisite for the core component in early mathematical skills, namely basic arithmetic. Basic arithmetic skills in 5–8 year-olds pertain to the degree to which a child masters mainly the addition and subtraction tasks with number symbols (Aunio & Räsänen, 2016). Basic arithmetic skills also depend on adequate counting skills. Frequent and successful use of counting strategies usually lead to improvements in memory representations of arithmetical facts and leads to the strategy of retrieving arithmetical facts from long-term memory (Canobi, Reeve, & Pattison, 2002; Wilkins et al., 2001). Correct and fluent number word sequences, part of counting skills are also relevant for solving basic arithmetic addition and subtraction tasks since children use counting-based strategies in the beginning when learning arithmetic, for example number word sequences advancing forward to solve addition problems and backward when solving subtraction problems.

Quality of educational assessment instruments – validation of an assessment tool

A range of different systems and criteria for judging the validity of a measurement exist (AERA et al., 1999; APA et al, 1974; Cosmin, NCME, 2018). Validity is a crucial psychometric notion since it concerns the degree to which the test scores provide information

that is related to the conclusions drawn from them. Validity is an evaluation of the degree to which empirical evidence and theoretical rationales support the adequacy and suitability of interpretations resulting from test scores of other models of assessment (Chan & Zumbo, 2009).

The European Federation of Psychologists' Associations (EFPA) (Evers, Hagemester, & Høstmælingen, 2013; Evers, Muñiz, Høstmælingen, Sjöberg, & Bartram, 2013) has provided a description and a thorough assessment of the psychological assessment tests, namely "EFPA Review Model for the Description and Evaluation of Psychological and Educational Tests". EFPA provides a manual for examining the quality of the test materials per se as well as psychometric properties of the identified documented instruments. In this study, validity and reliability are evaluated through EFPAs quality criteria.

EFPA emphasizes two types of validity – construct validity and criterion validity. Construct validity refers to whether the items represent the theoretical constructs that they are designed for. Construct validity, according to EFPA, is whether the test actually measures the intended construct or something else. Criterion-related validity is required for all kinds of tests and demonstrates the extent to which test instruments correlate with relevant valid instruments used for the same purpose to predict whether future or current performance is related to another measure of the same construct. In addition, criterion-related validity includes predictive validity, that is, whether the test is applicable in the sense that it serves its purpose, in this case to identify children at risk of developing mathematical learning difficulties (MLD). The types of validity mentioned here cannot be evaluated in isolation, they are complimentary. Overall validity refers to the appropriateness, meaningfulness, and usefulness of inferences based on scores (Cronbach & Meehl, 1955; Messick, 1995).

In addition, the extent to which a measure is prone to random measurement error is also important to assess. Reliability refers to the extent to which an instrument produces

random measurement errors and is expressed either as a standard error of measurement or as a reliability coefficient. Reliability is crucial if an assessment is to be useful, as a test that is not reliable can never be a valid instrument (Thorndike & Thorndike-Christ, 2014).

Existing assessments and screening instruments of early numeracy

Most European countries have mandatory math assessments or screening such as national tests. National assessments are designed for various purposes, not necessarily targeted to identify those who are at risk of developing MLD, but at the very least to identify children in need of extra support. Still, these tests are often introduced later in children's mathematical development and are not suitable for detecting those who are struggling at an early stage, more or less before formal math teaching starts. As for instruments other than the mandatory assessments, a review of early numeracy assessments shows that there is a scarcity of well-validated tests and screening tools (Dockrell et al., 2017). Dockrell et al. (2017) identified 23 omnibus tests and 16 tests assessing number concepts and skills, published in English. Dockrell et al. (2017) reviewed the extent to which the different instruments covered critical domains in numeracy, namely counting, transcoding (i.e. knowledge of number sequence, reading numerals, writing numerals, matching numerals to numbers), comparing numerical magnitudes and simple arithmetic.

In summary, only four of the 23 tests of numeracy featured items in all four areas that Dockrell et al. (2017) emphasized as critical domains; none of them included all types of items and none were well-validated. The tests included in Dockrell and colleagues' (2017) review were both screeners and numeracy subtests as part of a larger battery of mathematical tests (e.g. Keymath-3 (Connely, 2007)). The review concludes that choosing a test that does not cover all areas of numeracy skills can lead to misidentification of children with difficulties and as well as challenges in planning the content of support.

To ensure validity of assessments tools in mathematics to be used in different countries, they preferably need to be validated for each language and country. A challenge when validating tests in continental Europe is that many different languages are used, demanding considerable validation resources. The measures included in the review by Dockrell and colleagues (2017) were in English, but there are also validated screening instruments in other European languages. One widespread tool is the Dutch Early Numeracy Test (Van Luit, Van de Rijt & Pennings, 1994; Van de Rijt, Van Luit & Pennings, 1999) and the translation and norming of this into Finnish (Aunio, Hautamäki, Heiskari & Van Luit, 2006), German (Van Luit, Van de Rijt & Hasemann, 2001). Brankaer, Ghesquière, and DeSmedt (2017) validated the SYMP Test, a magnitude processing screening, and had a satisfactory test-retest reliability as well as construct and criterion-related validity.

In Norway, few studies exist on validation of screening instruments. The most commonly used tool for early numeracy is the summative national assessment test. Another commonly used instrument is Alle Teller! (McIntosh, 2012), but this test along with others is neither validated nor normed, hence leaving the teacher few options in assessing numeracy in early schooling. This does not necessarily imply that these assessment tools are of poor quality, but if they are not evaluated for instance relative to EFPA's criteria, they provide little or no information when it comes to the quality of the assessment tool. Thus, Norwegian schools and teachers are in an unfortunate situation where they are unable to consider the quality of the assessment tools.

The current study: Aims and research questions

The Early Numeracy Screener has a research-based theoretical foundation with three core components (Aunio & Räsänen, 2016), and can thus be used to improve the efficacy of assessment, as well as help teachers structure children's need of support more comprehensively in relation to the three core factors. The aim of the present study was to evaluate the

psychometric properties of the Early Numeracy Screener for Norwegian first graders. More specifically we investigated:

- 1) item-level difficulty and the internal consistency of the measure and its subscales (reliability)
- 2) if the hypothesized three-factor structure of the measure fits the data best (construct validity)
- 3) if the measure showed measurement invariance across gender and age (known group validity)
- 4) how the measure is related to the national tests in mathematics and if the measure identifies children who are at risk of developing mathematical learning disabilities (MLD) (compared to the national tests) (criterion-related validity)

Altogether, this provides insights into the Early Numeracy Screener as an evaluator of construct validity, criterion-related validity and reliability. Additionally, information about invariance will ensure whether the test has known groups validity, i.e. the extent to which a measurement is sensitive to differences and similarities in various groups such as gender, age groups etc. Thus, it is important to test invariance to ensure that the test works equally for boys and girls, and for children born at different times of the year (younger and older children within grade 1). If successfully validated, the Early Numeracy Screener would meet the need for a valid instrument that targets early numeracy skills for first graders that teachers would be able to use easily and efficiently. Our study will then make a contribution to educational practice in several ways: The present study enables Norwegian schools to flexibly use a validated screening tool and a screening tool that targets early numeracy skills separately and not unidimensional. This enables targeted support for children at risk of developing learning difficulties in mathematics.

Method

Participants

All children born in 2010 and attending first grade in two municipalities in Norway were invited to participate in the study. This is not a random sample, as we first contacted heads of the two municipal affairs, and from those municipalities each school's principal agreed to participate. Teachers in each class distributed and collected information about the study and a letter of consent to the parents. This resulted in a sample size of 366 participants (mean age 6.36 years, 55.7 % boys). None of the participants had been diagnosed with learning disabilities as this is often diagnosed later in school. Ethical approval was obtained from the Norwegian Social Science Data Services, and informed parental consent was obtained for each child in this study. A survey about the educational level of the parents and home language was distributed. The students were predominantly of Norwegian nationality and 13.9 % of the students had Norwegian as a second language (Table 1). The children were recruited from a district that was close to the national average on variables related to parental education (Norway: 25.8 % secondary school, 37.2 % high school, 27 % bachelor degree, and 10 % master's degree, recruitment area 26.9 % secondary school, 37.9 % high school, bachelor degree 26.1 %, and master's degree 10.1 % Statistics Norway, 2019; our sample secondary school 4.2 %, high school 31.5 %, bachelor degree 38.9 %, master's degree 20.1 %, Phd 0.8 %). Notably, 18.8 % did not reply to our survey, and prior studies indicate that those with low educational levels are less likely to answer these kinds of surveys (Goyder, Warriner, & Miller, 2002). However notably, the educational level is close to the national average with high school. Thus, since we had a cohort of 81.2 % of children it could be likely that it is approaching the national average, even if our survey do not indicate so.

Notably, in Norway children enter elementary school at the age of six and most of them have been to kindergarten from the age of one until they begin first grade. In a Norwegian context, this does not include a formalized kindergarten education program, but basically day care. There is no detailed common curriculum in mathematics, so the ways of introducing

children to early numeracy may differ from kindergarten to kindergarten. Moreover, the main focus in Norwegian kindergartens is on play and play-based learning rather than formal instruction.

Measures

Early Numeracy Screener.

The screening measure was developed based on Aunio and Räsänen's (2016) core numerical skills model. The tasks are representative of the range of skills assessed by other early numeracy measures (Clements, Sarama, & Liu, 2008; Jordan, Kaplan, Locuniak, & Ramineni, 2009). The screener consists of 52 items measuring core early numeracy skills – numerical relational skills, counting skills and basic skills in arithmetic (addition and subtraction). One point was given per correct answer and zero per wrong answer. The three components in the screener were as follows:

Numerical relational skills were assessed using 14 tasks. The children were asked to do tasks such as comparing numbers – e.g. “tick the box with the smallest number” and comparing concepts such as “one more than”, “as many as”, “one less than” (see Appendix A2 for item examples). One point was given per correct answer and zero per wrong answer.

Counting skills were measured with 28 tasks. The children were asked to do various counting tasks with ordinal numbers in number sequences e.g. “tick the third triangle” “tick the seventh star”, and tasks measuring the number-quantity correspondence (see appendix A2). One point was given per correct answer and zero per wrong answer.

Arithmetic skills were assessed with ten tasks, six of them in addition and four in subtraction. The sums and differences in the addition and subtraction tasks ranged from 0-15, that is, the answers were in this number range (see appendix A2). One point was given per correct answer and zero per wrong answer.

National test in mathematics.

The national test is a curriculum based standardized summative assessment in mathematics for first graders and is used in Norwegian schools in April of each year (Udir, 2016). It consists of 50 different early numeracy related items (Cronbach's $\alpha = .903$, from our sample); measuring number line, counting, number word knowledge and addition and subtraction fluency (Udir, 2016). The tasks in the National test also measures fluency. All the items are timed; they have for instance 1.5 minutes to solve a set of counting tasks. Item examples for counting skills and basic arithmetic are shown in (see Appendix A3) (Udir, 2016). Scoring of the national test varies somewhat throughout the test in the sense that some tasks are scored by giving one point for correct answers and zero for a wrong answer. However, some of the tasks require a correct sum score on two or three different tasks to obtain a score of 1. Consequently, in those 50 tasks there are tasks with several items all of which the children have to get correct to score a correct answer, which is given one point.

Procedure

Data for this validation study was collected October 2016. The teachers administering the screener were given a three-hour long training session in their respective schools. All the teachers administering the screener had prior experience with group-based assessment before this training, because the screener is administered in the same way as the national test in mathematics. The screener is a group-based, teacher-instructed, paper-and-pencil test that requires about 30–45 minutes to complete. The children were given instructions before each task and were told to write down the answers themselves. Teachers collected the answer sheets and sent them to the research group for correction and coding. The research group was thus involved with correcting both the Early Numeracy Screener and the national tests in an attempt to ensure that both tests were correctly scored.

Analysis

All analyses were performed in *Mplus 7.3* (Muthén & Muthén, 1998–2017). Due to categorical data, weighted least squares means and variance estimation (WLSMV) was used as estimator in all the analyses. Confirmatory factor analysis (CFA) was used to examine the factor structure of the test instrument. To evaluate the fit of the structural equation models that contained latent variables, we considered the common guidelines for model fit. These guidelines suggest an acceptable fit to the data if the Comparative Fit Index (CFI) and the Tucker Lewis Index (TLI) exceed .95, the Root Mean Square Error of Approximation is less than .08, and the Standardized Root Mean Square Residual (SRMR) is less than .10 (Hu & Bentler, 1999; Marsh, Hau, & Grayson, 2005).

Multigroup CFA with categorical data was conducted to test for measurement invariance across gender and age according to the guidelines of Muthen (Muthen & Muthen, 1998–2012, p.485). We used a model which assumes the same factor structure but allows factor loadings and item thresholds to vary freely across groups as a baseline model (configural). The configural model was then compared to a model where factor loadings and item thresholds are constrained to equality across groups (scalar invariance). When comparing nested models, a change of more than .01 in CFI, and .015 in RMSEA, indicates significant differences between the models (Chen, 2007).

To investigate whether the Early Numeracy Screener identifies the same children at risk of developing MLD as does the national test, configural frequency analysis was conducted. In this analyses, the children are grouped according to their performance in the Early Numeracy Screener into at-risk and typical performing groups and the same categorization is done with the national test scores. By means of configural frequency analysis we can then compare the extent to which the same children are identified with both tests. More specifically, the observed frequencies are compared to expected frequencies in a cross-tabulation and analyzed to ascertain whether cell frequencies are larger or smaller than could

be expected by a base model. The base model selected for frequency comparison was the first-order CFA, which assumes that all variables under study may show main effects and are independent of each other (von Eye, 1990, 1996). Thus, we focus on whether children identified as at-risk (or typical) with the Early Numeracy Screener are also identified as at-risk (or typical) with the national test (i.e., individual stability) as well as whether there are changes across the groups that cannot be explained by chance fluctuations (i.e., individual change).

Results

Table 2 shows the range, means, standard deviations, skewness, kurtosis, and reliabilities for all measures, and Table 3 shows correlations between them. Internal consistency was satisfactory. The same was the case for the distribution of the variables. Appendix A1 shows the skills measured on item level.

Item level difficulty and reliability. First, we examined the extent to which the different items discriminated between the children. We used the 95 % pass or fail as a criteria in order to secure variance, since we wanted to capture individual differences with the screener, hence the items cannot be so easy that more than 95 % solve them either correct or incorrect. Four of the items were removed because they were passed or failed by more than 95% of the children and thus were not useful in discriminating between them (Table A1). Notably, all the descriptive and results are thus based on the screener consisting of 52, and not the 56 original items. In addition, as shown in Table 2, the internal consistencies of all variables were satisfactory.

Construct validity and measurement invariance across gender and age for the one-factor model. Secondly, we examined the measurement model for the early numeracy construct as a one-factor model with all the items as indicators for a single overall early

numeracy factor. The fit of the measurement model was $\chi^2 (1274, N = 366) = 2956.486, p < .001$, RMSEA = .060 (90 % CI = .057 –.063), CFI = 0.911, TLI = .908; thus, this model fitted the data well. The multigroup CFAs indicated that the one-factor model showed measurement invariance across gender and age (see Table 4) and more importantly, the model fit did not worsen when constraining the factor loadings and item threshold to equality across gender, $\Delta CFI = 0.002$, $\Delta RMSEA = 0.000$

Construct validity and measurement invariance across gender and age for the three-factor model Next, we compared a three-factor model consisting of counting skills, numeric relational skills, arithmetic skills factors with the one-factor model. The three-factor model fitted the data better than the one-factor model, $\Delta CFI = 0.027$, $\Delta RMSEA = 0.01$, $\chi^2 (1271, N = 366) = 2439.967, p = .00$, RMSEA = .050 (90 % CI = .047 –.053), CFI = .938, TLI = .934. The correlation between the factors ranged from .70 to .77. These results indicate that early numeracy can be treated as a multidimensional construct and that the Early Numeracy Screener differentiates between the three subskills. Table 5 shows factor loadings, thresholds and variance in the three-factor model. As for the measurement invariance, the three-factor model was found to be invariant across gender and age (Table 6).

Criterion validity

First, we examined the correlations between the three early numeracy factors and the national test in mathematics taken 6 months after the early numeracy test. The fit of this model was good, $\chi^2 (1320, N = 366) = 2474.714, p = .00$, RMSEA = .049 (90 % CI = .046 –.052), CFI = .939, TLI = .936. The correlations between the national test scores and the early numeracy factors were low, $r = .25, p < .001$ (counting skills); $r = .20, p < .001$ (relational skills); $r = .19, p < .001$ (basic arithmetic skills).

Next, we grouped the children into at-risk of MLD and typical-performing groups based on the Early Numeracy and the national test scores using the 20th percentile established

as a critical limit in the national test as a cut-off. This resulted in four different configurations (typical/typical; typical/at-risk; at-risk/typical; at-risk/at-risk), for which the observed frequencies were compared to the expected frequencies (Table 7). To account for the increased risk of Type-I error, Bonferroni correction was applied to the significance testing ($0.05/4 = .013$) when comparing the observed and expected frequencies.

The configural frequency analysis identified 1 stable configuration and 2 antitypes of change. Of the 67 children who were identified as at-risk with the Early Numeracy Screener, 34 (51 %) were also identified as at-risk with the national test. This configuration was occurring more than expected by chance, indicating stability in classification of at-risk for MLD. This result was further supported by the fact, that moving from at-risk status to the typical group occurred less frequently than expected by chance (antitype of change) ($N = 33$, 9.2 %). However, the results for the typical-performing group was mixed as movement from typical-group (Early Numeracy) to the at-risk group (national test) was occurring less frequently than expected by chance ($N = 39$, 10.9 %). While at the same time the stability of belonging to the typical-group in both tests was not significant ($N = 253$, 70.4 %). Although, it is important to note that the p-value ($p = .024$) for the stable configuration was very close the Bonferroni-adjusted critical value of .013 indicating that it is rather likely to be identified as typical-performing in both tests. Overall, these results indicate that despite the rather low correlations between the Early Numeracy Screener and the National tests, children that are identified as at-risk for MLD with the Early Numeracy Screener are very likely to perform below the 20th percentile in the National test (Table 7).

Discussion

The main purpose of this study was to investigate the psychometric properties and to validate the Early Numeracy Screener. In line with EFPA (Evers, Hagemester, & Høstmælingen, 2013; Evers, Muñiz, Høstmælingen, Sjöberg, & Bartram, 2013) the study

evaluated reliability, construct validity and criterion-related validity. First, the items showed acceptable discriminant abilities. Regarding construct validity, the analysis suggested a three-factor model over a one-factor model. When it comes to criterion validity, the correlations between the three factors are below .85, which is often used as a cut-off (Brown, 2006). Additionally, high correlations between the core numerical skills are to be expected since the skills are related. The results indicate that the Early Numeracy Screener reliably measures three distinct subskills of early numeracy: counting skills, numerical relational skills and basic arithmetic skills. Moreover, these subskills were related to the national test scores and children identified as at risk of developing mathematical learning disabilities with the Early Numeracy test were likely to perform below the 20th percentile in the national test in mathematics six months later. The scores on the Early Numeracy Screener were on average lower than the National test. This is probably due to the different times of the school year in which the two tests are being administered. The Early Numeracy Screener was used in the very beginning in first grade, after only five weeks of schooling, and the children had short experience with mathematics instructions since this is not obligatory in Norwegian kindergartens. The National test is also designed with a ceiling effect, in which the sole aim is to identify those performing under the 20th percentile (Udir, 2016).

Construct Validity and Measurement Invariance

Concerning construct validity, firstly, we examined whether the Early Numeracy Screener worked best as a one-factor model or three-factor model. The CFAs showed better model fit for the three-factor model based on counting skills, numerical relational skills and basic arithmetic skills. This supports Aunio and Räsänen's (2016) theoretical model which suggests that the early numeracy skills are based on these three components. Support for assessing these skills is not found Aunio and Räsänen's (2016) study alone; support is also found in previous research (Desoete et al., 2012; Gersten et al. 2012; Moeller et al., 2011).

Here we see that although the three-factor model had a slightly better fit than the one-factor model, these are indeed highly related skills. Still, the three-factor model here had a better fit to the data and therefore this is chosen over the more parsimonious one-factor model.

Criterion-Related Validity

The analyses of strong concurrent and predictive relationships between the Early Numeracy Screener and the national test showed respectable criterion validity for the Early Numeracy Screener. However, the correlation between the Early Numeracy Screener and the national test proved to be quite low. There could be at least two reasons for this. Firstly, item types in the Early Numeracy Screener and National test are not very overlapping and this could imply that these two measures do not quite measure the same thing on item level, as all the items in the National test are timed, the children are given a certain amount of time, hence the National test also measures fluency, as opposed to the Early Numeracy Screener that measures accuracy and is without time limits. Since many of the tasks in the National test required children using more than one operation to solve (e.g. using both counting skills and relational skills within one task) which is not the case with the Early Numeracy Screener. In addition, in some of the tasks in the National test children were given one point if they had solved tasks consisting of multiple items (i.e. one task that was omitted one point for correct answer required the children to solve for example two items within the task correctly, if they managed only one out of two, the task was corrected with zero points). Secondly, the national test has a large ceiling effect, and this attenuates correlations. Also, the national screening was taken six months after the Early Numeracy Screener.

However, when we look at the relationship between the Early Numeracy Screener and those beyond the critical limit on the national test, there is a high correspondence between the two measures. Given that the Early Numeracy Screener measures core skills, it is not surprising that the children identified as at risk of developing mathematical learning

disabilities with the Early Numeracy Screener also have problems in national tests that are curriculum-based (the significant stability of the at-risk/at-risk configuration). At the same time, it is reasonable that some children that do not have problems in core skills might still struggle with the curriculum-based national test. However, they are probably not at risk of developing MLD; their difficulties might be related to other factors. Regarding how the screener predicts skills later on, in a general level, not solemnly children at risk of learning disabilities, the ceiling effect of the National test makes this impossible to say something about. Internal consistency of the Early Numeracy Screener and its subscales were also demonstrated to be adequate, and thus met the criteria evaluated through EFPA (Evers, Hagemester, & Høstmælingen, 2013; Evers, Muñiz, Høstmælingen, Sjöberg, & Bartram, 2013).

Implications and future studies

In this study, the Early Numeracy Screener demonstrated a three-factor structure enabling a brief screener with broader content, albeit not every aspect of mathematics, but a screener that arguably focuses on the core competencies that define early numeracy in six-year-old's (Gersten et al., 2005). To ensure the validity of assessment tools to be used in different countries, they preferably need to be validated for each language and country. A challenge when validating tests in continental Europe is that there is a large number of languages used, and therefore validation is resource demanding. However, there are validated screening instruments in other European countries (Aunio et al., 2006; Brankaer et al., 2017; Van Luit et al., 1994; Van de Rijt et al., 1999; Van Luit et al., 2001). For future studies it would be interesting to increase the validity even further with a test-retest design of the screening tool and perhaps open for doing ROC-curves analysis adding standardized mathematical measures in that respect. In this sense, we did not have any other criterion-related measure with which to relate the Early Numeracy Screener (e.g. SYMP test (Brankaer

et al.,2017)). However, the Early Numeracy Screener gives the teacher an opportunity to identify in which specific aspects of early numeracy a child needs remediation. The brevity and ease of use of the Early Numeracy Screener makes it well suited for classroom instructional settings. The present study contributes to the practice field in that it offers the schools a screening tool for measuring early numeracy skills that the schools can use whenever they want, i.e. unrestrained guidelines as to when pupils can take the national screener in mathematics that is normally administered at the end of the school year.

References

- American Educational Research Association (AERA), American Psychological Association (APA), & National Council on Measurement in Education (NCME). (1999). Standards for educational and psychological testing. Washington, DC: American Educational Research Association.
- American Psychiatric Association. (2013). *Diagnostic and statistical manual of mental disorders* (5th ed.). Arlington, VA: American Psychiatric Publishing
- American Psychological Association, American Educational Research Association, & National Council on Measurement in Education. (1974). Standards for educational and psychological tests. Washington, DC: American Psychological Association.
- Aunio, P., Hautamäki, J., Heiskari, P., & van Luit, J. (2006). The early numeracy test in Finnish: Children's norms, *Scandinavian journal of psychology*, 47(5), 369–378.
- Aunio, P. & Räsänen, P. (2016). Core numerical skills for learning mathematics in children aged five to eight years – a working model for educators. *European Early Childhood Education Research Journal*, 24:5, 684–704.
- Aunio, P. & Niemivirta, M. (2010). Predicting children's mathematical performance in grade one by early numeracy. *Learning and Individual Differences*, 20(5), 427–435.
- Baroody, A. J. (2003). The development of adaptive expertise and flexibility: The integration of conceptual and procedural knowledge. In A. J. Baroody & A. Dowker (Eds.), *The development of arithmetic concepts and skills: Constructing adaptive expertise*, 1–34. Mahwah, NJ: Erlbaum.

Baroody, J. (2006). Why children have difficulties mastering the basic number combinations and how to help them. National Council of Teachers of Mathematics,

<https://www.jstor.org/stable/41198838?origin=JSTOR->

[pdf&seq=1#metadata_info_tab_contents](https://www.jstor.org/stable/41198838?origin=JSTOR-pdf&seq=1#metadata_info_tab_contents)

Baroody, A. J., Eiland, M., & Thompson, B. (2009). Fostering at-risk preschoolers' number sense. *Early Education and Development*, 20, 80–128.

doi:10.1080/10409280802206619

Berch, D.B., (2005). Making sense of number sense: Implications for children with mathematical disabilities. *Journal of Learning Disabilities*, 38(4), 333-339.

<https://doi.org/10.1177%2F00222194050380040901>

Brankaer, C., Ghesquière, P. & De Smedt, B. (2017) Symbolic magnitude processing in elementary school children: A group administered paper-and-pencil measure (SYMP Test) *Behav Res* 49(4), 1361–1373.

Brown, T. (2006). *Confirmatory factor analysis for applied research*. New York: The Guilford Press

Butterworth, B., Varma, S., & Laurillard, D., (2011). Dyscalculia: From brain to education. *Science* 332(6033), pp. 1049–1053 DOI: 10.1126/science.1201536

Canobi, K.H., Reeve, R.A., & Pattison, P.E. (2002). *Young children's understanding of addition concepts*, *Educational Psychology*, 22(5), 513–532, DOI: [10.1080/0144341022000023608](https://doi.org/10.1080/0144341022000023608)

Chan E.K.H. (2014). Standards and Guidelines for Validation Practices: Development and Evaluation of Measurement Instruments. In: B. Zumbo & E. Chan (Eds.) *Validity and Validation in Social, Behavioral, and Health Sciences. Vol 54. Social Indicators Research Series*, p. 9–24. Switzerland: Springer

- Chen, F.F. (2007). Sensitivity of goodness of fit indexes to lack of measurement invariance, *Structural Equation Modeling: A Multidisciplinary Journal*, 14(3), pp.464–504
<https://doi.org/10.1080/10705510701301834>
- Clarke, B., Doabler, C., Smolkowski, K., Nelson, E.K., Fien, H., Baker, S.K., & Kosty, D. (2016). Testing the immediate and long-term efficacy of a tier 2 kindergarten mathematics intervention. *Journal of Research on Educational Effectiveness*, 9, 607–663. <https://doi.org/10.1080/19345747.2015.1116034>
- Clements, D.H., Sarama, J.H., & Liu, X.H. (2008) Development of a measure of early mathematics achievement using the Rasch model: the Research-Based Early Maths Assessment, *Educational Psychology*, 28(4), 457–482
- Connely, A.J. (2007). KeyMath-3 diagnostic assessment: Manual forms A and B. Minneapolis, MN: Pearson.
- Cronbach, L. J., & Meehl, P. E. (1955). Construct validity in psychological tests. *Psychological Bulletin*, 52(4), 281–302 <https://psycnet.apa.org/doi/10.1037/h0040957>
- De Smedt, B., Noël, M.-P., Gilmore, C., & Ansari, D. (2013). How do symbolic and non-symbolic numerical magnitude processing skills relate to individual differences in children’s mathematical skills? A review of evidence from brain and behavior. *Trends in Neuroscience and Education*, 2(2), pp.48–55,
<https://doi.org/10.1016/j.tine.2013.06.001>
- De Smedt, B., Verschaffel, L., & Ghesquière, P. (2009). The predictive value of numerical magnitude comparison for individual differences in mathematics achievement. 103(4), pp.469–479 <https://doi.org/10.1016/j.jecp.2009.01.010>
- Desoete, A., Ceulemans, A., De Weerd, F., & Pieters, S. (2012). Can we predict mathematical learning disabilities from symbolic and non-symbolic comparison tasks

- in kindergarten? Findings from a longitudinal study. *British Journal of Educational Psychology*, 82, 64–81. <https://doi.org/10.1348/2044-8279.002002>
- Desoete, A., Stock, P., Schepens, A., Baeyens, D., & Rieyers, H. (2009). Classification, seriation, and counting in grades 1, 2, 3 as two-year longitudinal predictors for low achieving in numerical facility and arithmetical achievement? *Journal of Psychoeducational Assessment*, 27 (3), pp. 252–264
<https://doi.org/10.1177/10734282908330588>
- Dockrell, J., Llauro, A., Hurry, J., Cowan, R., Flouri, E., & Dawson, A. (2017). Review of assessment measures in the early years. Language and literacy, numeracy, and social-emotional development and mental health, Education Endowment Foundation (EEF), Institute of Education. UCL.
- Dougherty, C. (2003). Numeracy, literacy, and earnings: Evidence from the National Longitudinal Survey of Youth. *Economics of Education Review*, 22, 511–521
- Duncan, G.J., Dowsett, C.J., Claessens, A., Magnuson, K., Huston, A.C., Klebanov, P., Pagani, L.S., Feinstein, L., Engel, M., Brooks-Gunn, J., Sexton, H., Duckworth, K., & Japel, C. (2007). School readiness and later achievement. *Developmental Psychology*, 43(6), pp. 1428–1446 DOI: 10.1037/0012-1649.43.6.1428
- Evers, A., Hagemester, C., & Hostmealingen, A. (2013). *EFPA A Review Model for the description and evaluation of psychological and educational tests*. Tech. Rep. Version 4.2.(6). Brussels: European Federation of Psychology Associations.
- Evers, A., Muñiz, C., Høstmælingen, A., Sjöberg, A., & Bartram, D. (2013). Assessing the quality of tests: Revision of the EFPA review model. *Psicothema*, 25(3), pp.283–291
- Fuchs, L.S., Geary, D.C, Compton, D.L., Fuchs, D., Schatschneider, C., Hamlett, C.L., ... Changas, P. (2013). Effects of first-grade number knowledge tutoring with contrasting

forms of practice. *Journal of Educational Psychology*, 105, 58–77.

<https://doi.org/10.1037/a0030127>

Geary, D.C., (2011). Cognitive predictors of achievement growth in mathematics: A 5-year longitudinal study. *Developmental Psychology*, 47(6), 1539–1552.

doi:10.1037/a0025510

Geary, D.C., Bailey, D.H., & Hoard, M.K. (2009). Predicting mathematical achievement and mathematical learning disability with a simple screening tool. *Journal of Psychoeducational Assessment*, 27(3), pp.265–278 doi:

<https://doi.org/10.1177/10734282908330592>

Geary, D.C., Hoard, M.K., Byrd-Craven, J., & DeSoto, M.C., (2004). Strategy choices in simple and complex addition: Contributions of working memory and counting knowledge for children with mathematical disability. *Journal of Experimental Child Psychology*, 88(2), pp.121–151 <https://doi.org/10.1016/j.jecp.2004.03.002>

Geary, D.C., Hoard, M.K., Nugent, L. & Bailey, D.H. (2013). Adolescents' functional numeracy is predicted by their school entry number system knowledge. *PLoS ONE* 8(1), 1–8.

Geary, D.C. & vanMarle, K., (2016) Young children's core symbolic and nonsymbolic quantitative knowledge in the prediction of later mathematics achievement.

Developmental Psychology, 52(12), 2130–2144.

<http://dx.doi.org/10.1037/dev0000214>

Gersten, R., Clarke. B., Jordan, N.C., Newman-Gonchar, R., Haymond, K., & Wilkins, C. (2012). Universal screening in mathematics for the primary grades: Beginnings of a research base. *Exceptional Children*, 78, 423–445.

<https://doi.org/10.1177/001440291207800403>

- Gersten, R., Jordan, N.C., & Flojo, J.R. (2005). Early identification and interventions for students with math difficulties. *Journal of Learning Difficulties*, 38, 293–304.
<https://doi.org/10.1177%2F00222194050380040301>
- Gersten, R., Rolffhus, E., Clarke, B., Decker, L., Wilkins, C., & Dimino, J. (2015). Intervention for first graders with limited number knowledge: Large-scale replication of a randomized controlled trial. *American Educational Research Journal*, 52, 516–546. <https://doi.org/10.3102/0002831214565787>
- Goyder, J., Warriner, K., & Miller, S. (2002). Evaluating socio-economic status (SES) bias in survey nonresponse. *Journal of Official Statistics*, 18(1), pp. 1-11.
- Halberda, J., Mazocco, M.M., & Feigenson, L., (2008). Individual differences in non-verbal number acuity correlate with maths achievement. *International Journal of Science*, 455, pp. 665–668
- Hannuka-Sormunen, M.M., Lehtinen, E., & Räsänen, P. (2015). Preschool children's spontaneous focusing on numerosity, subitizing, and counting skills as predictors of their mathematical performance seven years later at school. *Mathematical thinking and learning*, 17(2-3), pp. 155-177 <https://doi.org/10.1080/10986065.2015.1016814>
- Holloway, I.D., & Ansari, D., (2009). Mapping numerical magnitudes onto symbols: The numerical distance effect and individual differences in children's mathematics achievement. *Journal of Experimental Child Psychology*, 103(1), pp. 17–29
<https://doi.org/10.1016/j.jecp.2008.04.001>
- Hu, L.-t., & Bentler, P.M. (1991). Cutoff criteria for fit indexes in covariance structure analysis: Conventional criteria versus new alternatives, *Structural Equation Modeling*, 6(1), pp.1–55. <https://doi.org/10.1080/10705519909540118>

- Jordan, N.C., Glutting, J., Dyson, N., Hassinger-Das, B., & Irwin, C. (2012). Building kindergartners' number sense: A randomized controlled study. *Journal of Educational Psychology, 104*, 647–660. <https://doi.org/10.1037/a0029018>
- Jordan, N.C., Glutting, J. & Ramineni, C., (2010). The importance of number sense to mathematics achievement in first and third grades. *Learning and Individual Differences 20*, pp.82–88.
- Jordan, N.C., Kaplan, D., Ramineni, C., & Locuniak, M.N., (2009). Early math matters: Kindergarten number competence and later mathematics outcome. *Developmental Psychology, 45*(3), 850–867.
- Krajewski, K. & Schneider, W. (2009). Exploring the impact of phonological awareness, visual-spatial working memory, and preschool quantity-number competencies on mathematics achievement in elementary school: Findings from a 3-year longitudinal study. *Journal of Experimental Child Psychology, 103*, 516–531. <https://doi.org/10.1016/j.jecp.2009.03.009>
- Marsh, H.W., Hau, K.T., & Grayson, D. (2005). Goodness of Fit in Structural Equation Models. In A. Maydeu-Olivares & J. J. McArdle (Eds.), *Multivariate applications book series. Contemporary psychometrics: A festschrift for Roderick P. McDonald* pp. 275–340. Mahwah, NJ, US: Lawrence Erlbaum Associates Publishers
- McGuire, P., Kinzie, M.B., & Berch, D.B. (2012). Developing number sense in pre-k with five-frames. *Early Childhood Education, 40*(4), pp. 213-222
- McIntosh, A. (2012). *Alle Teller!* Trondheim: Matematikksenteret
- Messick, S. (1995). Validity of psychological assessment: Validation of inferences from persons' responses and performances as scientific inquiry into score meaning. *American Psychologist, Vol 50*(9), 741–749

- Moeller, K., Pixner, J., Zuber, J., Kaugmann, L. & Nuerk, H.-C. (2011). Early place-value understanding as a precursor for later arithmetic performance - A longitudinal study on numerical development. *Research in Developmental Disabilities, 38*, 1837–1851. <https://doi.org/10.1016/j.ridd.2011.03.012>
- Muthèn, L. K., and O. Muthèn. (1998–2004). Mplus user's guide. 3rd ed. Los Angeles, CA: Muthèn & Muthèn
- National Council on Measurement in Education. (1999)
- National Council of Teachers of Mathematics. (2006). Curriculum focal points for prekindergarten through grade 8 mathematics: A quest for coherence. Virginia: Reston
- National Mathematics Advisory Panel. Foundations for Success: The Final Report of the National Mathematics Advisory Panel. Washington, DC: U.S. Department of Education; 2008
- National Research Council. Mathematics Learning in Early Childhood: Paths Toward Excellence and Equity. Washington, DC: The National Academies Press; 2009
- Purpura, D.J., Baroody, A.J., & Lonigan, C.J., (2013). The transition from informal to formal mathematical knowledge: Mediation by numeral knowledge. *Journal of Educational Psychology, 105*(2), 453–464. <http://dx.doi.org/10.1037/a0031753>
- Purpura, D.J., & Lonigan, C.J., (2013). Informal numeracy skills: the structure and relations among numbering, relations and arithmetic operations in preschool. *American Educational Research Journal 50*(19), pp. 178–209. <https://doi.org/10.3102%2F0002831212465332>
- Ritchie, S.J. & Bates, T.C., (2013). Enduring links from childhood mathematics and reading achievement to adult socioeconomic status. *Psychological science, 24*(7), pp.1301–1308, <https://doi.org/10.1177%2F0956797612466268>

- Salminen, J.B., Koponen, T.K., & Tolvanen, A.L. (2018). Individuality in the early number skill components underlying basic arithmetic skills. *Frontiers in Psychology, 9*, pp. 1-11 <https://doi.org/10.3389/fpsyg.2018.01056>
- Schneider, M., Beeres, K., Coban, L., Merz, S., Schmidt, S., Stricker, J., & De Smedt, B. (2016). Associations of non-symbolic and symbolic numerical magnitude processing with mathematical competence: a meta-analysis. *Developmental Science*, doi:10.1111/desc.12372
- Statistics Norway. (2019). Retrieved from <https://www.ssb.no/utniv>
- Stock, P., Desoete, A., & Roeyers, H. (2009). Detecting children with arithmetic disabilities from kindergarten: Evidence from a 3-year longitudinal study on the role of preparatory arithmetic abilities. *Journal of Learning Disabilities, 43*, 250–268. <https://doi.org/10.1177/0022219409345011>
- Terwee, C. B., Prinsen, C. A. C., Chiarotto, A., De Vet, H. C. W., Westerman, M. J., Patrick, D. L. ... Mokkink, L. B. (2017). COSMIN standards and criteria for evaluating the content validity of health-related Patient-Reported Outcome Measures: a Delphi study. *Qual Life Res* in press.
- Thorndike, R.M., & Thorndike-Christ.T. (2014). *Measurement and evaluation in psychology and education* (8th ed. Ed.). Harlow: Pearson.
- Trilling, B. & Fadel, C. (2009). *21st century skills: Learning for life in our times*. San Francisco: Jossey-Bass
- Udir (2016), The Norwegian directorate of education and training
- Vanbinst, K., Ghesquière, P., & De Smedt, B., (2015). Does numerical processing uniquely predict first graders' future development of single-digit arithmetic? *Learning and Individual Differences, 37*, pp.153–160 <https://doi.org/10.1016/j.lindif.2014.12.004>

- Van de Rijt, B. A. M., Van Luit, J. E. H. & Pennings, A. H. (1999). The construction of the Utrecht Early Mathematical Competence Scale. *Educational and Psychological Measurement*, 59, 289–309. <https://doi.org/10.1177/0013164499592006>
- Van Luit, J.E.H., Van de Rijt, B.A.M, & Hasemann, K. (2001). OTZ Osnabrücker test für zahlbegriffsentwicklung [German version of the Utrecht test of number sense]. Göttingen, Germany : Hogrefe-Verlag.
- Van Luit, J.E.H., Van de Rijt, B.A.m., & Pennings, A.H. (1994). *Utrechtse Getalbegrip Toets* [Utrecht Test of Number Sense] Doetinchem, The Netherlands: Graviant
- Von Eye, A. (1990). Introduction to configural frequency analysis. The search for types and anti-types in cross classification. Cambridge: University Press
- Von Eye, A., Spiel, C., & Wood, P.H. (1996). Configural Frequency Analysis. *Applied Psychological Research*, 45(4), 310–327
- Wilkins, J.L.M., Baroody, A.J., & Tiilikainen, S. (2001). Kindergartner's understanding of additive commutativity within the context of word problems. *Journal of Experimental Child Psychology*, 79(1), 23–26.
- Wright, R. J., Martland, J., Stafford, A. K. (2006). Early numeracy. Assessment for teaching and intervention. 2. edition. London: Sage.
- Udir, 2016. The Norwegian Directorate of Education and Training.
- Zhang, X., Koponen, T., Räsänen, P, Aunola, K., Lerkkanen, M-K, & Nurmi, J-E. (2014). Linguistic and spatial skills predict early arithmetic development via counting sequence knowledge. *Child Development*, 85(3), pp. 1091-1107
<https://doi.org/10.1111/cdev.12173>
- Zumbo, B.D. & Chan, E.K.H. (2009). Setting the stage for validity and validation in social, behavioral, and health sciences: Trends in validation practices. In: B. Zumbo & E.

RUNNING HEAD: Validation of a numeracy screener for first graders

Chan (Eds.), *Validity and Validation in Social, Behavioral, and Health Sciences*.

Social Indicators Research Series: vol 54. (3–8) Switzerland: Springer.

RUNNING HEAD: Validation of a numeracy screener for first graders

Table 1

Background information: Home language, parental educational background

Home language	%	Educational background	%
Norwegian	86.10		
English	2.70		
Lithuanian	1.30	Secondary school	4.20
Urdu	1.30	High school	31.50
Somali	1.00	Bachelor	39.20
Tamil	1.00	Master	20.10
Bosnian	.06	PhD	4.20
Dutch	.06	<i>Missing</i>	.08
Kurdic	.06		
Polish	.06		
Sign language	.06		
Arabic	.03		
Burmese	.03		
Finnish	.03		
Icelandic	.03		
Portugese	.03		
Rumanian	.03		
Russian	.03		
Swedish	.03		
<i>Missing</i>	.03		

RUNNING HEAD: Validation of a numeracy screener for first graders

Table 2

Mimima, maxima, means, SD, skewness, kurtosis, and reliability for one-factor model, three-factor model and national test

	Min	Max	Mean	SD	Skewness	Kurtosis	Cronbach's α
Age	5.87	6.88	6.38	.30	.087	-1.239	
One-factor model, full scale	2	52	30.85	11.58	-.108	-.764	.943
Three-factor model							
Counting skills	2	28	18.31	6.60	-.406	-.704	.904
Numerical relational skills	0	14	7.95	3.21	.385	-.692	.792
Arithmetic skills	0	10	4.58	3.19	.803	-.935	.901
National test	6	50	43.55	7.29	-1.467	4.459	.903
	Min	Max	Mean	SD	Skewness	Kurtosis	Cronbach's α

Table 3

Correlations of the three-factor model and the national test

	Counting skills	Numerical relations	Arithmetic	National test
Counting skills	1			
Numerical relations	.744**	1		
Arithmetic	.785**	.712**	1	
National test	.415**	.395**	.300**	1

Note. **.Correlation is significant at the .01 level

Table 4

Model fit statistics for the test of invariance across gender and age– one-factor model

Model	χ^2 (<i>df</i>)	CFI	RMSEA (90 % CI)	TLI	Δ CFI	Δ RMSEA
Gender						
Configural Invariance	2548	.934	.051 (.047-.054)	.931		
Scalar invariance	2598	.932	.051 (.047-.056)	.931	.000	.000
Age						
Configural Invariance	2548	.913	.052 (.049-.056)	.909		
Metric invariance						
Scalar invariance	2598	.914	.052 (.048-.055)	.912	-.001	.000

RUNNING HEAD: Validation of a numeracy screener for first graders

Table 5
Factor loadings, thresholds and variance in -factor model

Item	Factor 1 Relational	Factor 2 Counting	Factor 3 Aritmethic	Item Threshold	Residual Variance	Explained variance (R ²)
1						
2	.703			-0.705	.506	.494
3	.596			-1.429	.645	.355
4	.701			-1.468	.509	.491
5		.552		-0.914	.695	.305
6		.596		-1.174	.645	.355
7		.528		-0.679	.722	.278
8		.441		-1.001	.805	.195
9		.539		0.048	.709	.291
10		.473		0.110	.776	.224
11	.670			-1.392	.551	.449
12	.761			-0.041	.422	.578
13	.793			0.200	.372	.628
14	.728			0.263	.470	.530
15	.543			-0.461	.705	.295
16	.551			0.796	.696	.304
17						
18	.773			-0.014	.403	.597
19	.668			0.089	.553	.447
20	.604			0.158	.635	.365
21	.627			0.228	.607	.393
22	.355			0.313	.874	.126
23		.486		-0.768	.764	.236
24		.666		-0.645	.556	.444
25		.674		-0.103	.546	.454
26		.741		-0.873	.451	.549
27		.781		-0.571	.390	.610
28		.751		0.014	.436	.564
29		.667		-0.492	.556	.444
30		.625		-0.172	.609	.391
31		.671		0.179	.550	.450
32		.751		-0.207	.435	.565
33		.862		0.357	.256	.744
34		.845		0.371	.285	.715
35		.733		-1.001	.463	.537
36		.721		-0.645	.480	.520
37		.678		-0.371	.540	.460
38		.765		-1.059	.415	.585
39		.785		-0.555	.384	.616
40		.726		0.027	.473	.527
41		.703		-1.160	.506	.494
42		.584		-0.978	.659	.341
43		.818		-0.285	.332	.668
44		.815		0.221	.336	.664
45			.895	-0.603	.198	.802
46			.891	-0.466	.205	.795
47			.845	-0.342	.285	.715
48			.974	-0.563	.050	.950
49			.870	-0.130	.244	.756
50			.835	0.379	.303	.697
51			.924	0.560	.146	.854
52			.957	0.698	.084	.916
53			.966	0.813	.067	.933
54						
55			.993	0.861	.015	.985

Note. Items number 1, 17, 54, 56 was taken out of the analysis

Table 6

Model fit statistics for the test of invariance across gender and age – three-factor model

Model	χ^2 (<i>df</i>)	CFI	RMSEA (90 % CI)	TLI	Δ CFI	Δ RMSEA
Gender						
Configural Invariance	2542	.956	.041 (.037-.056)	.954		
Scalar invariance	2588	.956	.041 (.037-.045)	.954	.000	.000
Age						
Configural Invariance	2542	.938	.044 (.040-.048)	.935		
Metric invariance						
Scalar invariance	2588	.938	.044 (.040-.047)	.937	.000	.000

RUNNING HEAD: Validation of a numeracy screener for first graders

Table 7

Change and stability in the mathematics achievement groups across the Early Numeracy test and the national test

Configuration (EN-NT)	o	e	z	p(z)	
TA – TA	253	232.62	2.25	0.0243	
TA – MLD	39	59.38	-2.89	0.0038	antitype
MLD – TA	33	53.38	-3.02	0.0025	antitype
MLD – MLD	34	13.62	5.63	0.0001	type

Note. EN = Early Numeracy test; NT = national test; o = observed frequencies; e = expected frequencies; TA = typical-achieving group; MLD = at-risk for mathematical learning difficulties group.

RUNNING HEAD: Validation of a numeracy screener for first graders

Appendix A1

Item	Skill	% correct	alpha if item removed	Retained
1	Numerical relational	95.4 %	.943	No
2	Numerical relational	76 %	.942	Yes
3	Numerical relational	92.3 %	.943	Yes
4	Numerical relational	92.9 %	.942	Yes
5	Counting	82 %	.942	Yes
6	Counting	88 %	.942	Yes
7	Counting	75.1 %	.942	Yes
8	Counting	84.2 %	.943	Yes
9	Counting	48 %	.942	Yes
10	Counting	45.6 %	.942	Yes
11	Numerical relational	91.8 %	.942	Yes
12	Numerical relational	51.6 %	.942	Yes
13	Numerical relational	42.1 %	.942	Yes
14	Numerical relational	39.6 %	.942	Yes
15	Numerical relational	67.8 %	.942	Yes
16	Numerical relational	21.3 %	.942	Yes
17	Numerical relational	95.4 %	.943	No
18	Numerical relational	50.5 %	.942	Yes
19	Numerical relational	46.4 %	.942	Yes
20	Numerical relational	42.7 %	.942	Yes
21	Numerical relational	41 %	.942	Yes
22	Numerical relational	37.7 %	.943	Yes
23	Counting	77.9 %	.943	Yes
24	Counting	74 %	.942	Yes
25	Counting	54.1 %	.942	Yes
26	Counting	80.9 %	.942	Yes
27	Counting	71.6 %	.941	Yes
28	Counting	49.5 %	.941	Yes
29	Counting	68.9 %	.941	Yes
30	Counting	56.8 %	.941	Yes
31	Counting	42.9 %	.941	Yes
32	Counting	58.2 %	.941	Yes
33	Counting	36.1 %	.941	Yes
34	Counting	35.5 %	.942	Yes
35	Counting	84.2 %	.942	Yes
36	Counting	74 %	.942	Yes
37	Counting	64.5 %	.942	Yes
38	Counting	85.5 %	.941	Yes
39	Counting	71 %	.941	Yes
40	Counting	48.9 %	.942	Yes

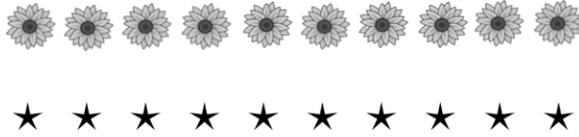
RUNNING HEAD: Validation of a numeracy screener for first graders

41	Counting	87.7 %	.942	Yes
42	Counting	83.6 %	.941	Yes
43	Counting	61.2 %	.941	Yes
44	Counting	41.3 %	.941	Yes
45	Addition	72.7 %	.941	Yes
46	Addition	67.8 %	.941	Yes
47	Addition	63.4 %	.941	Yes
48	Addition	71.3 %	.941	Yes
49	Addition	55.2 %	.941	Yes
50	Addition	35.2 %	.941	Yes
51	Subtraction	28.7 %	.941	Yes
52	Subtraction	24.6 %	.941	Yes
53	Subtraction	20.8 %	.941	Yes
54	Subtraction	20.2 %	.941	No
55	Subtraction	19.4 %	.941	Yes
56	Subtraction	18.3 %	.941	No

Appendix A2 Sample items for the Early Numeracy Screener

Note: The instructions are translated from Norwegian to English. Teacher instructions in italics.

Counting skills:



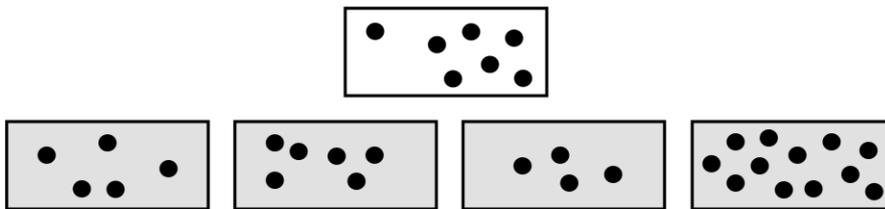
'Tick on the second flower, tick on the sixth star, tick on the fourteenth ball.'

Counting skills:



'Beside the box with black dots, there are boxes with numbers in them. First, find out how many black dots there are all together, and tick the box that says how many black dots there is in the box.'

Numerical relational skills:



'There are black dots in the white box. Tick the grey box that has one black dot less than the white does.'

Numerical relational skills:



'Here are three numbers. Tick on the smallest number.'

Basic arithmetic skills:

$$2 + 5 = \underline{\quad}$$

$$8 + 3 = \underline{\quad}$$

'Here you have addition tasks. Solve as many of them as you can.'

$$7 - 4 = \underline{\quad}$$

$$11 - 3 = \underline{\quad}$$

'Here you have subtraction tasks. Solve as many of them as you can.'

Appendix A3

Note: The instructions are translated from Norwegian to English. Teacher instructions in italics.

“How many”, time limit: 1.5 min

HVOR MANGE?

EKSEMPEL

2
4
5
6

3
6
7
12

2
10
11
20

3
8
9
10

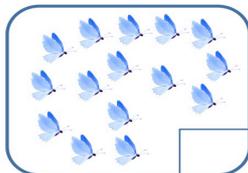
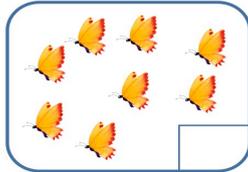
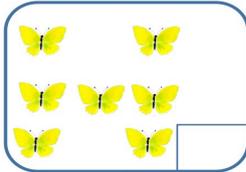
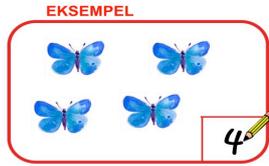
13
15
23
31

3
15
16
20

‘On this page you are going to find out how many bricks there are in each task. Look at the example (point). There are six bricks. That is why there is a circle around the number six. Now you are going to find out how many bricks there are in the other tasks’.

“How many”, time limit: 1.5 min

HVOR MANGE?



‘Here you need to find out how many butterflies there are in each task and write the amount in the little square. Look at the example (point). There it is four butterflies. That is why the number four is written in the small square. Now you are going to do the rest of the tasks on this page.’

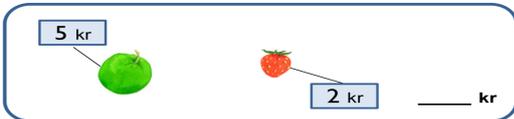
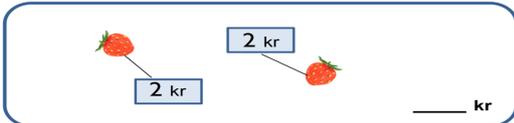
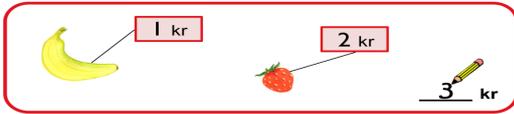
RUNNING HEAD: Validation of a numeracy screener for first graders

“How much all together”, time limit: 1.5 min

Note: krone(r) (kr =NOK)

HVOR MYE TIL SAMMEN?

EKSEMPEL



‘On this page you should imagine that you are buying things and you need to find out how much you need to pay. The price tag tells how much each thing costs. Look at the example. The banana costs 1 krone and the strawberry costs 2 kroner. Then we need to pay 3 kroner in total, which is why it is written two on the line (point). Now you are going to find out how much you have to pay for the stuff in the other tasks.’