



UvA-DARE (Digital Academic Repository)

Output-driven citizenship education

A focus on learning outcomes to improve the quality of citizenship education

Hoek, L.H.M.

Publication date

2023

[Link to publication](#)

Citation for published version (APA):

Hoek, L. H. M. (2023). *Output-driven citizenship education: A focus on learning outcomes to improve the quality of citizenship education*. [Thesis, fully internal, Universiteit van Amsterdam].

General rights

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

Disclaimer/Complaints regulations

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

Chapter 2

Measuring Citizenship Competences: Assessment of Measurement Invariance*

Abstract

Standardised questionnaires are used to measure the outcomes of citizenship education. These outcomes are often compared across groups to document different outcomes, for example, between boys and girls. A prerequisite for cross-group comparisons is an assessment of measurement invariance.

This study used data from 6035 students from 87 Dutch primary schools to examine the measurement invariance of the Citizenship Competences Questionnaire (Ten Dam et al., 2011). Dutch schools use this questionnaire to gain insight into students' citizenship knowledge, attitudes, and skills. Measurement invariance was assessed across sex, socioeconomic position, and migration background.

Measurement invariance was sufficient in most cases, allowing for cross-group comparisons of associations between latent constructs and their indicators, and in some cases, for cross-group comparisons of the latent means. We conclude by emphasising that periodic assessment of measurement invariance in instruments measuring citizenship competences is important due to the dynamic nature of the construct.

Keywords: Citizenship competences; Measurement invariance; Measurement equivalence; Citizenship education; Social outcomes.

*This chapter is based on: Hoek, L., Zijlstra, B. J. H., Munniksma, A., & Dijkstra, A. B. (2023). Measuring citizenship competences: Assessment of measurement invariance. *Journal of Social Science Education*, 22(2). <https://doi.org/10.11576/jsse-5837>

Introduction

Like other educational domains, it is essential to gain insight into what students learn in citizenship education. This insight can be used to facilitate the learning process and evaluate the contents and delivery methods of the curriculum. There are several ways to obtain insight into student outcomes, including standardised measurement instruments (Daas et al., 2016). Using standardised measurement instruments is beneficial in many ways, for example, because of their practical usefulness and the opportunities for securing the quality and validity of the measurement. This is why standardised questionnaires have been widely used to measure students' citizenship competences in terms of knowledge, attitude, and skill (Ireland et al., 2006; Schulz et al., 2018).

Multi-group comparisons based on standardised questionnaires have demonstrated that citizenship competences are related to students' background characteristics, such as sex, socioeconomic position (SEP), and migration background - both internationally (Ireland et al., 2006; Kerr et al., 2007) and in the Dutch school system in which our data were gathered (Dijkstra et al., 2015; Geijsel et al., 2012). These findings are important for educational practice because they may lead to adaptations in the contents or delivery methods of citizenship education so that the learning objectives are met, and all students benefit optimally. However, a prerequisite for meaningful comparisons across different population groups (e.g., boys and girls) is that the construct to be measured is understood similarly in each group (Isac et al., 2019; Steinmetz et al., 2009). In other words: it is essential to know whether the difference in citizenship competences between, e.g., boys and girls, is a 'true' difference or a difference caused by boys and girls systematically understanding questions in a different way. To determine whether a questionnaire measures the same across groups, the measurement invariance should be assessed (Meredith, 1993).

Assessment of measurement invariance is important for the methodological quality of a measurement instrument (Meuleman et al., 2022). However, as the interpretation of a construct can change over time (Putnick & Bornstein, 2016), the establishment of measurement invariance is not a static given. This applies in particular to citizenship education, which is a dynamic construct that, like other societal phenomena, takes on meaning in and moves along with changes in society (Mattei & Broeks, 2018). This underscores the importance of periodic assessment of measurement invariance in student questionnaires, as is the case for measuring citizenship competences. Examples of such changes, amongst others, include increased cultural diversity (Eurostat, 2021; US Census Bureau, 2020) or social inequality (Organisation for

Economic Co-operation and Development, 2017). Both changes, for instance, may have influenced the public debate and, resultingly, have led to different interpretations of constructs related to citizenship education.

For example, students with and without migration backgrounds may perceive debates about how to deal with socioeconomic or socio-cultural differences differently nowadays, in times of what is sometimes referred to as a shift from diversity to ‘super-diversity’ (Vertovec, 2007), as opposed to earlier in time, also influencing the response in measurement instruments regarding these topics. Hence, it is important to periodically assess the measurement invariance in instruments that measure citizenship competences, in case changes in society may have shifted the norms and values of citizens regarding value-sensitive themes (Munck et al., 2018).

Whereas both conceptualisation and operationalisation are vital steps in developing measurement instruments, the assessment of measurement invariance relates explicitly to the operationalisation of a measurement instrument. It assumes that the conceptualisation of the instrument is carefully considered and, therefore, not associated with differences in student characteristics like sex, SEP, migration background, religiosity, age and others. The assumption of such a generic conceptualisation of citizenship is the basis of large-scale standardised measurement instruments for citizenship competences. National and international examples include instruments like the International Civic and Citizenship Education Study (ICCS) for international comparisons used in periodic cycles (Schulz et al., 2018), the Citizenship Education Longitudinal Study (CELS) for longitudinal research in the UK (Clever et al., 2006), or the Citizenship Competences Questionnaire (CCQ) for annual measurements used by schools in The Netherlands (Ten Dam et al., 2011). In the conceptualisation of such instruments, important building blocks are core democratic values, trust, participation and involvement, and the support for institutions of a democratic society. Such instruments generally move away from conceptualisations that take a position regarding an ‘ideal’ balance between values like, for example, individual freedom and solidarity towards others – also known as contested or school-specific goals in citizenship education (Eidhof et al., 2016). Rather, these conceptualisations adhere to consensus goals, referring to commonly shared values of reciprocity of ‘treating others as you want others to treat you’ – also known as *regula aurea* or the Golden Rule (cf. Etzioni, 1996; Wattles, 1996).

Despite its importance, research on measurement invariance in instruments for citizenship competences is still scarce. It is predominantly assessed in between-country research (Byrne & Watkins, 2003), such as the ICCS. The technical report of the ICCS (Schulz et al., 2016), for example, demonstrated that most of the citizenship constructs show

measurement invariance up to a level where associations of latent constructs can be safely compared across countries, but not always to the level where comparison of means of these latent constructs is justified. In addition, Isac et al. (2019) assessed the measurement invariance of a specific part of the ICCS questionnaire that focuses on young people's tolerant attitude towards immigrants (i.e., support for equal rights). The authors found that most items were measurement invariant to the level where average scale scores can be validly compared across European countries. However, the assessment of measurement invariance is also important in within-country comparisons (Steinmetz et al., 2009; Vandenberg & Lance, 2000). In this respect, some researchers mentioned that homogeneity of the population in terms of the measures being compared across groups is often implicitly assumed, regardless of whether it was assessed (Muthén, 1989; Steinmetz et al., 2009).

This study aims to examine the measurement invariance of a Dutch questionnaire for assessing citizenship competences across groups. By doing so, we assess whether this questionnaire, which was designed over fifteen years ago, is still² valid across different groups in society today, as societal changes may have affected how different groups interpret the questionnaire items. In this way, we contribute to the insight into the assessment of measurement invariance in within-country comparative research on citizenship competences, as this empirical knowledge base on measurement invariance in these questionnaires is still scarce. Indirectly, this study sheds light on whether the measurement of citizenship competences is robust amid changes in society and population. For educational practice, this study serves as a validity check of a measurement instrument that schools use to improve their citizenship education.

Following previous research showing relevant and robust differences in citizenship competences based on sex, SEP, and migration background (Geijsel et al., 2012; Ireland et al., 2006; Kerr et al., 2007), we focus on the assessment of measurement invariance across these three student-background characteristics. The findings of this study may underscore the results of these previous studies, which are largely built on the same generic conceptualisation of citizenship competences, in case the measurement instrument appears largely measurement invariant or, alternately, place caution on the findings of these studies in case measurement invariance could not be established. Moreover, various scholars have advocated investigating measurement invariance over sex, SEP, and migration background (Kline, 2015; Wray-Lake

² Still, as was the case during the initial construction phase when the measurement invariance was tested as part of a broad set of psychometric tests, meeting all necessary requirements (Ten Dam et al., 2011; Geijsel et al., 2012; based on personal communication, because these results were not available in print).

et al., 2017) or have stressed the importance of examining similarities and differences in citizenship competences based on these groups (Cleaver et al., 2006). The instrument used for testing measurement invariance is the CCQ, a large-scale standardised measurement instrument based on a generic conceptualisation of citizenship competence. The CCQ is used primarily for annual measurements of students' citizenship competences at Dutch schools. The questionnaire has been used for over a decade – yielding a rich dataset of consecutive cohorts that is eminently suitable to assess whether contextual changes have deteriorated alignment between the instrument and its context.

Theory

Meaningful Comparison of Groups

In this section, building on He and Van de Vijver (2013) and Isac et al. (2019), we outline three examples of how the operationalisation of measurement does not measure the same across groups. We also provide a 'counter-example' of a difference between groups that is not the result of measurement non-invariance but an example of a 'true' difference. To illustrate our examples, we used existing items from the CCQ (measuring citizenship attitudes, skills, and knowledge of students). Students needed to indicate to what extent the items applied to them on a four-point Likert scale ranging from 'does not apply at all to me' to 'applies completely to me' for attitude items or ranging from 'not good at all' to 'very good' for skill-items. Students had to pick the best answer in the knowledge items by choosing one out of three. The use of these items at this place is merely illustrative.

First, we outline how measurement non-invariance is caused by the fact that some underlying items are not considered indicative of the construct to be measured for some group members. To illustrate this for SEP, we look at the item: 'If we talk about the news in class, I want to add something to the conversation too'. Compared to students with a high SEP, students with a low SEP may have limited access to news sources, such as a subscription to a daily newspaper or a personal mobile device to check news websites. Therefore, students with a low SEP may answer this item with 'does not apply at all to me', whereas they are willing to contribute to the conversation. Therefore, their answer indicates their accessibility to news, not their democratic attitude. This may cause students with a low SEP to be wrongly labelled as 'less capable' of the construct 'democratic attitude' because this specific indicator of the construct is less applicable to their context.

Second, we outline how measurement non-invariance is caused by some group members understanding items differently due to linguistic differences. To illustrate this for migration

background, we look at the item: ‘How good are you at... holding on to your opinion, if you are really right?’ Linguistic differences between students with and without a migration background may influence how students understand constructs or underlying items. Language that involves ambiguity in meaning (e.g., ‘holding on to something’, ‘if you are really right’) or metaphorical language is particularly susceptible to differences in interpretation, for example, by students who have another mother language (more likely being students with a migration background). However, vice versa, multilingual students (more likely students with a migration background) may also benefit from their knowledge and skill in language comprehension. This may be less common for monolingual students (more likely students without a migration background). Either way, such linguistic differences may distort the results and wrongly label students as more or less skilled in citizenship competences.

Likewise, group members may understand items differently due to cultural differences. To illustrate this for sex, we look at the item: ‘It is normal to help in the household (for example, by preparing the dinner table, tidying up or cleaning)’. In some cultures, helping in the household is considered predominantly a task for girls and not for boys. Boys who grow up in such a culture may be more likely to answer this question with ‘does not apply to me at all’, whereas their attitude towards ‘acting in a socially responsible manner’ may, in reality, be different than this answer reveals. Or vice versa: girls who grow up in such a culture may be more likely to answer this question with ‘applies completely to me’, whereas their attitude towards ‘acting in a socially responsible manner’ may, in reality, be different than this answer reveals.

Third, we outline how measurement non-invariance is caused by some groups characterised by a specific response style. To illustrate this for sex, we look at the item: ‘People who earn sufficient salary should together take care of people with less wealth’. Suppose that girls are more likely to answer items in a more socially desirable way and systematically answer this item with ‘applies completely to me’, whereas boys answer this item in a less socially desirable way. This may cause girls to overestimate their attitude toward the construct ‘acting in a socially responsible manner’.

At last, we outline how a difference can be not the result of measurement non-invariance but a ‘true’ difference between groups. To illustrate this, we look at the item: ‘In a sports game, the referee takes a wrong decision *against* your team. What should you do?’ And the following answer options: (a) Go to the referee and debate the decision; (b) Get the coach of your sports team; (c) Keep on playing because the decision of the referee is directive during the game. The latter is appointed the preferred answer. Student background characteristics such as sex, SEP,

or migration background may influence how group members respond to this item. However, these differences are not an example of measurement non-invariance when the differences do not adhere to differences in understanding, interpretation, or applicability to the context, but rather are an example of ‘true’ differences in what is valued about acting in a socially responsible manner. Regardless of these ‘true’ differences, the conceptualisation of acting in a socially responsible manner holds that if participants of a sports game agree on rules beforehand, they conform to these rules during the sports game – even if they do not agree during the game.

Assessment of Measurement Invariance

Multiple-group confirmatory factor analysis (MGCFA) is the most commonly used method to assess measurement invariance (Putnick & Bornstein, 2016). In MGCFA, hierarchical, subsequent models with increasing restrictions are specified and compared. These are the configural model, the metric model, and the scalar model.

The *configural model* assesses whether the instrument measures the same latent factors across groups and whether the indicators are the same across groups (Isac et al., 2019). To test for configural invariance, models with the same pattern (i.e., the same configuration) should be specified across groups (Vandenberg & Lance, 2000), meaning an equal number of latent variables, indicators, et cetera. If the configural model yields poor model fit, it indicates that in one of the groups, a different pattern fits the data. For example, for boys, one of the questions may not be an indicator of the latent construct, whereas, for girls, it is. If the configural model fits the data well, it provides ground for testing metric invariance.

The *metric model* assesses whether the factor loadings differ across groups (Horn & Mcardle, 1992). To test for metric invariance, the factor loadings are constrained to equality across groups (Reeskens & Hooghe, 2010). Thus, for example, two equal path models with the same pattern are specified with the first factor loading in the model for boys constrained to equality to the first factor loading in the model for girls, and so on for the remaining factor loadings. If the metric model yields poor model fit, it indicates that in one of the groups, a factor loading relates differently to the latent variable as compared to the other group (e.g., for boys, the third item is strongly related to the latent variable, but for girls, this is not the case). If the metric model fits the data well, it allows for subsequent analyses of testing scalar invariance. Reaching metric invariance justifies the comparison of latent variables and their associations across groups (Isac et al., 2019). In addition, achieving (partial) metric invariance is seen as a minimal prerequisite for meaningful cross-group comparisons (Little, 2013).

The *scalar model* assesses whether the intercepts (i.e., the constant or the scalar) of the

indicators are the same across groups (Steenkamp & Baumgartner, 1998). To test for scalar invariance, the intercepts of the indicators are constrained to equality across groups. Thus, in addition to the identical pattern (configural model) and equal factor loadings (metric model) across groups, the intercepts are constrained to equality. This means that the intercept of the first item for, e.g., boys is constrained to equality to the intercept of the first item for girls, and so on for the intercepts of the remaining items. If the scalar model yields poor fit, at least one intercept differs across groups (Isac et al., 2019). If the scalar model fits the data well, it justifies comparing latent means across groups (Reeskens & Hooghe, 2010).

Methodology

In this study, we examined the measurement invariance in the assessment of citizenship competences (i.e., the competence of students to function and participate in society) across sex, SEP, and migration background. We did so by looking at competences of students in terms of attitude (i.e., thoughts, desires, and willingness), knowledge (i.e., knowing, understanding, insight), and skill (i.e., an estimate by students of what they think they are able to).

Data

We used data of consecutive cohorts from the CCQ (Ten Dam et al., 2011). Schools use the CCQ to measure students' citizenship knowledge, attitude and skill in terms of four so-called 'social tasks': acting democratically, acting socially responsible, dealing with conflicts, and dealing with differences (Ten Dam et al., 2011). The CCQ is suitable for grades 5 and 6 of primary education (approximate age is 10 to 12 years old) and grades 7, 8, and 9 from secondary education (approximate age is 12 to 16 years old). We retained this study's scope to data gathered in primary education for pragmatic reasons. In this sample, we merged data from grades 5 and 6 to obtain larger group sizes and more robust results.

Sample and Procedure

The sample consists of 6035 students from 87 primary schools participating in the Dutch 'Alliance Citizenship' (Table 1). The Alliance Citizenship is a partnership of schools that organises annual measurements of citizenship competences of students. This study used data from 2015, 2016, 2017, 2018, and 2019. We used a sample that was both recent and large enough to detect possible changes in society and population composition. Each year, the CCQ is online available to participating schools during spring. The sample of schools differed each year: some schools participated more than once; some only once³.

³ The possibility that some students were measured twice (i.e., from grade 5 in one years' measurement to grade 6 in the following years' measurement), as well as the nested structure of the data, may have caused a degree of dependency that we did not account for. In practice, it means that our results may be somewhat too conservative.

Variables

This section provides information on the conceptual framework underlying the questionnaire that we examined in our analyses (Table 2). The values for Cronbach's alpha are calculated with the sample used in the present study. The values were largely consistent with the original alpha's provided by Ten Dam et al. (2011), based on 16,000 students from grade 6 and grade 9 who participated in 2005, 2006, and 2007.⁴ The variables Attitude – Acting democratically and Skill – Acting democratically are both presented in two interpretable factors. The variable Skill – Acting in a socially responsible manner is presented in a combined factor with the scale Skill – Dealing with conflicts.

⁴ In order to only share robust results, the results of the questionnaire are reported at the overarching scale-level (e.g., knowledge, attitude, or skill), and not on subscale-level. We performed the analyses of measurement invariance on both the scale- and subscale-level to gain more information on exact items that could be hindering reaching a higher level of measurement invariance. Hence, we provided the values for Cronbach's alpha also on the subscale-level (of which some values are below the advised threshold of 0.70 (Cortina, 1993)).

Table 1*Descriptive Information of the Sample*

	2015	2016	2017	2018	2019	Total
Primary schools (N)	26	18	16	11	19	87
Students (N)	1587	1349	1176	564	1359	6035
Grade						
Grade 5	809	675	605	296	678	3063
Grade 6	778	674	571	268	681	2972
Age						
10 years or younger	309	225	198	121	339	1192
11 years	676	638	570	272	643	2799
12 years	389	391	315	137	301	1533
13 years	55	33	30	9	19	146
14 years	1	3	0	0	0	4
15 years	0	0	0	0	0	0
16 years or older	1	0	1	1	1	4
Sex						
Boy	734	642	555	254	628	2813
Girl	697	648	559	285	675	2864
SEP						
Low SEP ^a	430	402	393	192	535	1952
High SEP ^b	608	526	425	178	364	2101
Migration background						
No migration background ^c	1102	989	847	401	1086	4425
Migration background ^d	334	302	267	139	217	1259

Note. SEP = socioeconomic position. Low SEP = highest educational level of mother and father is ‘no school’, ‘only primary education’ or ‘only secondary education’; high SEP = highest educational level of mother and father is ‘higher education’. No migration background = both parents are born in the Netherlands; migration background = at least one parent is born outside of the Netherlands.

Table 2

Conceptual Framework of Citizenship Competences

Components	Knowledge ($\alpha = .79$) Knowing, understanding, insight A young person with such knowledge...	Attitudes ($\alpha = .89$) Thoughts, desires, willingness A young person with such attitudes...	Skills ($\alpha = .87$) Estimate of what one can do A young person with such skills...
Social tasks Acting democratically Acceptance of and contribution to a democratic society	... knows what democratic principles are and what acting in accordance with them involves (8 items, $\alpha = .67$)	... <i>wants to hear everyone's voice, enter into a dialogue</i> (3 items; $\alpha = .65$) and <i>make an active, critical contribution</i> (3 items; $\alpha = .65$)	... <i>is able to assert own opinions</i> (3 items; $\alpha = .74$) and <i>listen to the opinions of others</i> (3 items; $\alpha = .68$)
Acting in a socially responsible manner Taking shared responsibility for the communities to which one belongs	... knows social rules (i.e. legal or unspoken rules for social interaction) (6 items, $\alpha = .54$)	... wants to uphold social justice, is prepared to provide care and assistance, does not want to harm another or the environment as a result of his or her behaviour (6 items, $\alpha = .68$)	... can adopt a socially just position (5 items, $\alpha = .76$)
Dealing with conflicts Handling of minor situations of conflict or conflicts of interest to which the child is a party	... knows methods to solve conflicts such as seeking win-win solutions, calling in help from others, admission of mistakes, prevention of escalation (7 items, $\alpha = .62$)	... is willing to explore conflicts, prepared to consider the standpoint of another, jointly searches for an acceptable solution (6 items, $\alpha = .79$)	... can listen to others, put oneself in someone else's place, seek win-win solutions (5 items, $\alpha = .76$)
Dealing with differences Handling of social, cultural, religious, and outward differences	... is familiar with cultural differences, has knowledge of rules of behaviour in different social situations, knows when one can speak of prejudice or discrimination (6 items, $\alpha = .63$)	... has a desire to learn other people's opinions and lifestyles, has a positive attitude towards differences (6 items, $\alpha = .85$)	... can adequately function in unfamiliar social situations, adjust to the desires or habits of others (4 items, $\alpha = .67$)

Note. The conceptual framework is derived from (Ten Dam et al., 2011)

For the Knowledge items, students chose the answer they thought was the best response. An example is: *All children have a right to... (a) an allowance; (b) choose who they want to live with; or (c) education.* The correct answer here is c. All knowledge items were dichotomised into correct (1) or incorrect (0). The phrasing of the Attitude items is: *How well does this statement apply to you?* A sample statement is: *I like knowing something about different religions.* The response options are: *(1) does not apply at all to me, (2) does not apply much to me, (3) applies a fair amount to me, or (4) applies completely to me.* The phrasing of the Skill items is: *How good are you at...?* A sample statement is: *... finding a solution that everyone is satisfied with for a conflict?* The response options are: *(1) not good at all; (2) not very good; (3) pretty good; or (4) very good.*

For sex, boys (49.55%) were appointed with value 1 and girls with value 2. For SEP, we converted maternal and paternal educational levels into a new variable, indicating whether the mother or father's highest obtained educational level is either no school, primary school, or secondary school (value 1) or higher education (value 2). In our sample, 48,16% of all students had a low SEP. For migration background, we converted maternal and paternal country of birth into a new variable, indicating whether both parents of the students were born in the Netherlands (value 1) or whether at least one parent of a student has a migration background, regardless of the country of origin (value 2). In our sample, 77,85% of all students had both parents born in the Netherlands.

Analytic Plan

We conducted MGCFA and compared groups based on sex (boys versus girls), SEP (low versus high), and migration background (having no migration background versus having a migration background). Analyses were conducted with R version 3.5.0 (R Core Team, 2021), using the lavaan package (Rosseel, 2012) and the semTools package (Jorgensen, 2021), with particular focus on the guidelines in the measEq.syntax. While doing so, analyses were performed according to a the following principles.

First, we treated the indicators that were directly observed (questionnaire items) as ordered categorical factors because they were measured using a four-point Likert scale (attitude and skill) or a three-option multiple-choice (knowledge). Therefore, they could not be treated as continuous indicators. Assessing measurement invariance using ordered categorical factors makes the procedure more complex because it also involves testing for threshold invariance. Testing for threshold invariance is a means to do justice to the ordered categorical nature of indicators. The threshold model assumes that a normally distributed latent item-response lies underneath each observed categorical indicator (Kite et al., 2018). The threshold can be seen

as a ‘tipping point’ between the different category responses (e.g., between answering ‘does not apply at all to me’ and ‘does not apply much to me’). The threshold model is estimated before the metric model, but only if polytomous indicators are involved (thus: only for attitude and skill).

In the case of dichotomous indicators (such as the knowledge items), it is not possible to distinguish between equality constraints on the thresholds (testing threshold invariance), factor loadings (testing metric invariance), and intercepts (testing scalar invariance) (Wu & Estabrook, 2016). Instead, the equality constraints need to be added simultaneously. Therefore, we tested only for configural and scalar invariance for dichotomous indicators.

Data Preparation

All empty cases were dropped via list-wise deletion by removing all rows with less than nine items answered. This removed the rows in which only certain standard school characteristics (e.g., unique identifier for school, year and class) were automatically filled in (8 items) but none of the actual questionnaire items. The resulting data had 1.28% missing observations that did not show any pattern of missingness. The percentage of missing data is considered acceptable as basis for further analyses (Bennett, 2001). Second, four variables were excluded from our subsequent analyses because they consisted of three items, which was insufficient as input to our fit measures in the configural model. These excluded variables are: ‘Attitude – Acting democratically 1’, ‘Attitude – Acting democratically 2’, ‘Skill – Acting democratically 1’, and ‘Skill – Acting democratically 2’.

Model Selection Procedure

To assess how well the configural, threshold, metric, and scalar models fit our data, we consulted the Chi-square test of overall model fit (χ^2), and the relative difference in model fit ($\Delta \chi^2$), the comparative fit index (*CFI*), and the root mean square error of approximation (*RMSEA*) – where the *CFI* and *RMSEA* are functions from the Chi-square test statistic (Shi et al., 2019). We followed five decision rules in the model selection procedure using these model fit measures.

The first decision rule is based on the *CFI*. We consulted the *CFI* to indicate how well the model fit improved compared to the null model. The *CFI* considers the complexity of the model and ranges from 0 to 1. A value above 0.95 is preferred (Hu & Bentler, 1999).

The second decision rule is based on the *RMSEA*. The *RMSEA* indicates the “badness-of-fit” (Adelson, 2012) and considers the model complexity by estimating approximation error per model degree of freedom. Larger values of the *RMSEA* indicate a worse model fit. An *RMSEA* lower than 0.05 indicates ‘close fit’, and an *RMSEA* between 0.05 and 0.08 indicates

‘satisfactory fit’ (Browne & Cudeck, 1993). In this study, we perceived a value for $RMSEA \leq 0.08$ as acceptable.

The third decision rule is based on the Chi-square test of overall model fit. This fit index tests the assumption of respectively configural, threshold (if applicable), metric (if applicable), and scalar invariance. A downside of the Chi-square test for overall model fit is that even minor deviations from the suggested models can reject the null hypotheses of model fit for large samples (Shi et al., 2019). Therefore, we proceeded to test additional invariance models as long as values for CFI and $RMSEA$ of the configural model are acceptable – even if the Chi-square value is significant. At the same time, we performed permutation tests to ensure that the rejection of the overall model fit is most likely due to the same underlying reason across groups (Jorgensen, 2017). In permutation tests, all observations are randomly reassigned (a thousand times) to groups. These permuted group compositions are expected to fit the data equally poor or well as the hypothesised group composition (in which we ‘manually’ assigned, e.g., all boys to the one group and all girls to the other group). If this is the case, it is indicated by a non-significant p-value. However, a significant p-value suggests that there might be different underlying reasons across groups causing the significant Chi-square test of overall fit.

The fourth decision rule is based on the p-value of the Chi-square test of the difference in model fit ($\Delta \chi^2$). This fit index was applied to test the null hypothesis of equal model fit for two models. If the p-value of the Chi-square test of difference is significant, it indicates a meaningful difference between the two models in comparison, and we reject the more restricted model.

The fifth decision rule entails whether or not to test for partial invariance. A partial model can be seen as an ‘in-between model’ because some, but not all, of the equality constraints may need to be rereleased. To assess what constraints this involves, we inspected standardised mean residuals and modification indices. Standardised mean residuals can be consulted to identify which factor loadings, thresholds, or intercepts differ most across groups and may need to be rereleased. In addition, modification indices estimate how much the Chi-square test statistic of a model decreases if a constrained parameter is set to free again. We used a Bonferroni-adjusted alpha level to see whether modification indices were significant and could be considered released. When values for CFI and $RMSEA$ are within the thresholds, but the p-value for the Chi-square test of difference is significant, we additionally tested for partial measurement invariance. To maintain analyses within the scope of this article, we only searched for partial invariance if it established a meaningful difference, i.e. reaching partial metric invariance (for polytomous indicators) or partial scalar invariance (for polytomous and

dichotomous indicators). We did not test for partial invariance if the meaningful minimum of partial metric invariance (Little, 2013) was still more than one step away (e.g., when the configural model fitted well, but the threshold model did not).

Results

We specified two types of first-order one-factor models: (1) models in which social tasks such as ‘acting democratically’ are specified as indicators of overarching constructs such as citizenship attitude, knowledge, or skill; and, zooming in, (2) models in which items from the questionnaire are specified as indicators of the aforementioned social tasks. We summarised the findings in Tables 3, 4 and 5. In this section, we elaborate on the results by first sharing the results of the models with social tasks as indicators of citizenship knowledge, attitude, and skill, and secondly by sharing the results of the models with questionnaire items as indicators of the social tasks.

Attitude

Social Tasks as Indicators

Comparing citizenship attitude across sex and SEP, metric invariance was reached. This justifies comparing associations between the construct citizenship attitude and its four underlying social tasks across sex and SEP. Comparing citizenship attitude across migration backgrounds, metric invariance could not be reached. Inspecting modification indices and standardised mean residuals pointed to no possibility of establishing a partial model. Therefore, only configural invariance was reached. This means that the path models specified for students with and without migration backgrounds are likely to follow the same pattern. For an overview, see Table 3.

Table 3*Models with Social Tasks as Indicators of Citizenship Constructs*

	Configural			Metric			Scalar			Partial Scalar		
	CFI	R		CFI	R	$\Delta \chi^2$	CFI	R	$\Delta \chi^2$	CFI	R	$\Delta \chi^2$
Attitude	Sex	+		+	+	+	-	-	-	+	+	+
	SEP	+		+	+	+	+	+	-	+	+	+
	Migr.	+		+	+	-						
Knowledge	Sex	+	-									
	SEP	+	-									
	Migr.	+	-									
Skill	Sex	+	+	+	+	+	+	+	-	+	+	+
	SEP	+	+	+	+	+	+	+	-	+	+	+
	Migr.	+	+	+	+	+	+	+	+	+	+	+

Note. Migr. = migration background; R = RMSEA. $\Delta \chi^2$ = Chi-square test of difference. A '+' for CFI indicates that the value for CFI is ≥ 0.95 ; a '-' indicates that the value is ≤ 0.95 . A '+' for RMSEA indicates that the value for RMSEA is ≤ 0.08 ; a '-' indicates that the value is ≥ 0.08 . A '+' for $\Delta \chi^2$ indicates that the p-value of the Chi-square test of difference is ≥ 0.05 ; a '-' for $\Delta \chi^2$ indicates that the p-value of the Chi-square test of difference is ≤ 0.05 . The full tables can be found in Appendices 2.1, 2.2 and 2.3.

Table 4

Models with Polytomous Questionnaire Items as Indicators of Social Tasks

	Configural			Threshold			Metric			Partial Metric			Scalar			Partial Scalar		
	CFI	R	$\Delta \chi^2$	CFI	R	$\Delta \chi^2$	CFI	R	$\Delta \chi^2$	CFI	R	$\Delta \chi^2$	CFI	R	$\Delta \chi^2$	CFI	R	$\Delta \chi^2$
Attitude – Acting in a socially responsible manner																		
Sex	+	+	+	+	+	+	+	+	+	+	+	-	+	+	-	+	+	+
SEP	+	+	+	+	+	-	+	+	+	+	+	+	+	+	+	+	+	+
Migr.	+	+	+	+	+	+	+	+	+	+	+	-	+	+	-	+	+	+
Attitude – Dealing with conflicts																		
Sex	+	+	+	+	+	+	+	+	+	+	+	-	+	+	+	+	+	+
SEP	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
Migr.	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
Attitude – Dealing with differences																		
Sex	+	-	-	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
SEP	+	-	-	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
Migr.	+	-	-	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
Skill – Acting in a socially responsible manner / Dealing with conflicts																		
Sex	+	-	-	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
SEP	+	-	-	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
Migr.	+	-	-	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
Skill – Dealing with differences																		
Sex	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
SEP	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
Migr.	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+

Note. Migr. = migration background; R = RMSEA, $\Delta \chi^2$ = Chi-square test of difference. A ‘+’ for CFI indicates that the value for CFI is ≥ 0.95 ; a ‘-’ indicates that the value is ≤ 0.95 . A ‘+’ for RMSEA indicates that the value for RMSEA is ≤ 0.08 ; a ‘-’ indicates that the value is ≥ 0.08 . A ‘+’ for $\Delta \chi^2$ indicates that the p-value of the Chi-square test of difference is ≥ 0.05 ; a ‘-’ for $\Delta \chi^2$ indicates that the p-value of the Chi-square test of difference is ≤ 0.05 . The full tables can be found in Appendices 2.1, 2.2 and 2.3.

Table 5*Models with Dichotomous Questionnaire Items as Indicators of Social Tasks*

Configural		Scalar		Partial Scalar	
CFI	R	CFI	R	CFI	R
Knowledge – Acting democratically					
Sex	+	+	+		
SEP	+	+	+		
Migr.	+	+	+		
Knowledge – Acting in a socially responsible manner					
Sex	+	+	+		
SEP	+	+	+		
Migr.	+	+	+		
Knowledge – Dealing with conflicts					
Sex	+				
SEP	+				
Migr.	+				
Knowledge – Dealing with differences					
Sex	+	-	+		
SEP	+	+	+		
Migr.	+	+	+		

Note. Migr. = migration background; R = RMSEA. $\Delta \chi^2$ = Chi-square test of difference. A '+' for CFI indicates that the value for CFI is ≥ 0.95 ; a '-' indicates that the value is ≤ 0.95 . A '+' for RMSEA indicates that the value for RMSEA is ≤ 0.08 ; a '-' indicates that the value is ≥ 0.08 . A '+' for $\Delta \chi^2$ indicates that the p-value of the Chi-square test of difference is ≥ 0.05 ; a '-' for $\Delta \chi^2$ indicates that the p-value of the Chi-square test of difference is ≤ 0.05 . The full tables can be found in Appendices 2.1, 2.2 and 2.3.

Items as Indicators

Comparing the social tasks within citizenship attitude, (partial) metric invariance was reached in most cases. This allows for comparing the associations between the concerning social tasks and the underlying items across groups. Threshold invariance was achieved for *Attitude – Dealing with conflicts* across sex, but metric invariance could not be established. This means that the thresholds are the same across boys and girls, but the factor loadings between questionnaire items and the social task differ across boys and girls. For comparisons of *Attitude – Dealing with differences* across sex, SEP and migration background, the configural model indicated poor model fit. This suggests that we cannot assume that the path models specified for boys and girls, students with low and high SEP, and students with and without migration background, follow the same pattern; the same number of indicators relating to the same number of latent constructs. For an overview, see Table 4.

Knowledge

Social Tasks as Indicators

Comparing citizenship knowledge across sex, SEP and migration background, the configural model indicated poor model fit. This suggests that the data follow a different pattern in one of the groups. For example, it may be the case that one of the indicators of citizenship knowledge, e.g., dealing with differences, is an indicator of citizenship knowledge for students with a low SEP but not for students with a high SEP. For an overview, see Table 3.

Items as Indicators

Comparing the social tasks within citizenship knowledge, the minimum of (partial) scalar invariance was reached in most cases. This means that we can meaningfully compare the mean scores of the concerning social tasks across groups. For comparisons of the social task *Knowledge – Dealing with conflicts* across sex and SEP, the configural model indicated poor model fit. For an overview, see Table 5.

Skill

Social Tasks as Indicators

Comparing citizenship skill across sex and SEP, metric invariance was reached. This justifies the comparison of citizenship skill and its associations with the underlying social tasks across boys and girls and students with a low and high SEP. Comparing citizenship skill across migration backgrounds, scalar invariance was reached. This means that we found no evidence that the intercepts differ across students with and without migration backgrounds, and we can meaningfully compare the mean scores of citizenship skill. For an overview, see Table 3.

Items as Indicators

Comparing the social tasks within citizenship skill, the required minimum of (partial) metric invariance was reached for comparisons of *Skill – Dealing with differences* across sex, SEP, and migration background. This means we may compare the associations between the social task and its underlying questionnaire items across groups. However, for comparisons of *Skill – Acting in a socially responsible manner/Dealing with conflicts* across sex, SEP, and migration background, the configural model indicated poor model fit. This means that the path models will likely follow a different pattern across the groups. For an overview, see Table 4.

Permutation Tests

The permutation tests resulted predominantly in non-significant p-values, indicating that the significant Chi-square tests of overall model fit depended on the same underlying reasons across groups (also see Appendix 2.5). However, in the comparison of *Skill – Dealing with differences* across sex and *Attitude – Acting in a socially desirable way* across migration background, the p-value was respectively 0.045 and 0.012. Here, the significant Chi-square tests of overall model fit may depend on different underlying reasons across groups. This means that we need to be cautious when interpreting the results of these two comparisons.

Exploratory Analyses

We established a partial metric model for *Skill – Dealing with differences* across migration backgrounds. For *Attitude – Acting in a socially responsible way* across SEP, *Knowledge – Dealing with conflicts* across migration background, and *Knowledge – Dealing with differences* across sex, we established partial scalar models. This means that in these models, some indicators were less invariant than others and hindered reaching a higher level of measurement invariance. We performed exploratory analyses to zoom in on these indicators by looking at the content of the questionnaire items, the response frequencies, the factor loadings and the item-rest correlations (also see Appendix 2.4). In most cases, we found that the difference in factor loadings across groups was the largest for the less invariant indicators.

For example, comparing *Attitude – Acting in a socially desirable way* across SEP, we established a partial scalar model with freed equality constraints on the second and fourth indicator intercepts. Exploratory analyses of these two indicators demonstrated that the factor loading was higher for students with a high SEP for the second indicator. For the fourth indicator, the opposite was true. This means that the second indicator (*‘If a classmate is being called names in the streets, I want to stick up for him or her’*) more strongly relates to the construct *Attitude – Acting in a socially desirable way* for students with a high SEP as compared to students with a low SEP. Contrary, the fourth indicator (*‘You should say sorry if you did*

something that hurt another person') more strongly relates to the construct for students with a low SEP as compared to students with a high SEP.

We also found a difference in factor loadings across groups. For example, comparing *Efficacy – Dealing with differences* across migration background, we established a partial metric model with freed equality constraints on the factor loading of the first indicator. Exploratory analyses demonstrated that for the first indicator ('*How good are you at... adapting your behaviour to the rules and habits of others?*'), the factor loading was larger for students without a migration background, whereas for the remaining indicators ('*How good are you at... behaving normally in an unfamiliar environment?*'; '*... adapting your language to the person you are speaking with?*'; '*... considering the wishes of others when making a decision together?*'), the factor loadings were larger for students with a migration background.

The permutation tests resulted predominantly in non-significant p-values, indicating that the significant Chi-square tests of overall model fit depended on the same underlying reasons across groups. However, in the comparison of *Skill – Dealing with differences* across sex and *Attitude – Acting in a socially desirable way* across migration background, the p-value was respectively 0.045 and 0.012. Here, the significant Chi-square tests of overall model fit may depend on different underlying reasons across groups. This means that we need to be precautionary when interpreting the results of these two comparisons.

Discussion

Standardised questionnaires are widely used to measure the outcomes of citizenship education in terms of knowledge, attitude, and skill (Ireland et al., 2006; Schulz et al., 2018). The outcomes of these questionnaires are often compared across groups based on student characteristics, such as boys and girls. Insight into the outcomes of citizenship education is important to evaluate the effectiveness of the content and delivery methods of the curriculum. Moreover, insight into the extent to which these outcomes differ across groups helps to ensure that all students have equal opportunity to acquire citizenship competences. However, a prerequisite for valid cross-group comparisons based on standardised questionnaires entails addressing whether the instrument measures the same across groups and later points in time. This can be done by examining the measurement invariance (Meredith, 1993).

The establishment of measurement invariance is, however, not a static given because the interpretation of constructs can change over time. This may particularly be the case for value-sensitive concepts like citizenship and, accordingly, citizenship education, which is considered a dynamic construct that takes on meaning in context and moves along with changes in society. Hence, it is important to periodically examine the measurement invariance of

measurement instruments capturing citizenship competences. In line with this rationale, this study aimed to investigate the measurement invariance of the CCQ (Geijsel et al., 2012; Ten Dam et al., 2011), which is a standardised questionnaire based on a generic conceptualisation of citizenship competences used over a long period. Schools use the instrument to gain insight into students' citizenship knowledge, attitudes, and skills across four social tasks: acting democratically, acting in a socially responsible manner, dealing with conflicts and dealing with differences. To check whether this instrument over time still measures the same across groups in recent samples, we conducted comparisons within three student characteristics which are known to relate to robust differences in citizenship competences: sex, SEP, and migration background (Dijkstra et al., 2015; Geijsel et al., 2012).

Our study showed that in two-thirds of the comparisons across groups, the meaningful minimum of at least (partial) metric invariance (Little, 2013) was reached. This allows for comparisons of the associations between latent constructs and their indicators across groups (i.e., 'the relation between indicator 1 and latent construct X is stronger for boys than for girls'). We even reached (partial) scalar invariance across groups in some comparisons. This allows for comparing latent means across groups (i.e., 'girls, on average, obtain a higher score on latent construct X than boys'). Our findings also mark a warning for researchers conducting multi-group comparative analyses based on citizenship constructs or social tasks where the minimum required level of measurement invariance could not be established. In this study, this applied to one-third of the comparisons, being one citizenship construct (i.e., knowledge) and two social tasks (i.e., Skill – Acting in a socially responsible manner/Dealing with conflicts, and Attitude – Dealing with differences), where fit indices of the configural model indicated overall poor model fit. This means that, for these three aspects, we have no grounds to assume that the instrument measures the same across sex, SEP and migration background, and we need to be cautious when making cross-group comparisons. Nevertheless, the results on these aspects can still be used to describe all students in a class or school as a whole. Moreover, it is advised to assess measurement invariance over these aspects again when using new samples. Whereas we advocate periodic assessment of measurement invariance of *all* citizenship constructs, our findings underscore this is indeed important for future studies examining the construct and social tasks mentioned above.

In some comparisons, specific indicators hindered establishing a higher level of invariance. Therefore, we searched for partial invariance where these particular indicators remained on the lower level of invariance while the other indicators were specified at a higher level of invariance. Exploratory analyses of the partial models pointed to no unidirectional

reason why these particular indicators may have been answered differently across groups. That is, we could not identify measurement non-invariance due to either of the three explanations mentioned by Van de Vijver (2013) and Isac et al. (2019): (1) some group members not considering some underlying indicators to be indicative of the construct; (2) cultural or linguistic differences across group members; or (3) a specific response-style of some group members. To improve the questionnaire regardless, we suggest performing additional qualitative data analysis on these non-invariant indicators to see whether this does reveal unidirectional leads of how to improve the indicators. This can be done in individual interviews or panel group interviews with students originating from all comparison groups. Alternately, new indicators can be added to the scales and tested among students to see whether they do show measurement invariance across groups. Eventually, new invariant indicators can replace the non-invariant indicators.

While the number of studies that assessed the measurement invariance of citizenship competences within a country is scarce, we can, to some extent, compare our findings to the studies that performed a between-country assessment of measurement invariance, which is more common in large-scale international assessments of citizenship competences. An example is ICCS (Schulz et al., 2016), where most scales were invariant across countries. The technical report places caution in comparing only two scales (i.e., students' perceptions of the importance of citizenship behaviours; and students' attitudes toward civic institutions and their country of residence). Isac et al. (2019) also used data from ICCS and, similarly, found that most scales measuring students' attitudes towards immigrants were invariant across countries. Only the scale measuring students' attitudes towards equal rights for immigrants was more heterogeneous across countries. The indicator 'Immigrants should have the opportunity to continue their own customs and lifestyle' was the weakest in the scale. In addition, Munck et al. (2018) found that the measurement of students' attitudes toward immigrants was largely invariant over time (from 1999 to 2009) across countries. However, the indicator 'immigrants should have the opportunity to keep their own language' was invariant in only 36 out of 92 countries. Whereas the studies of Isac et al. (2019) and Munck et al. (2018) showed that between-country measurement invariance was more difficult to establish in measuring students' attitude in dealing with differences, our findings of non-invariance in Attitude – Dealing with differences indicate the same in within-country comparisons. However, these previous studies do not comply with the non-invariance that we found for the construct Knowledge and the social task Skill – Acting in a socially responsible manner/Dealing with conflicts.

Although we grounded our findings on extensive analyses, it is important to bear in mind the limitations of this study. Caution needs to be placed on the way we composed the groups. We cared to comply with scholars who underscored that the robustness and precision of measurement invariance analyses increase as sample sizes increase (Koh & Zumbo, 2008; Meade, 2005; Meade & Lautenschlager, 2004). Therefore, we combined data from students of five subsequent years sharing the same student characteristics, and consequently, we were limited in comparing measurement invariance over time. However, the comparison over time is an important lead for future research. We were also somewhat limited in the extent to which we could differentiate within each group. For example, we combined data from students whose parents were born in countries other than the Netherlands into 'having a migration background'. Whereas it is not inconceivable that the effect on measurement invariance differs within these countries, this operationalisation suffices to provide insight into the main effect of migration background in assessing measurement invariance. The same may hold for SEP, where we combined different attained parental educational levels into low SEP.

Despite the limitations, this study contributes to the empirical knowledge base of the assessment of measurement invariance in instruments for citizenship competences. Whereas previous research on the assessment of measurement invariance in instruments for citizenship competences focused on between-country comparisons (Isac et al., 2019; Schulz et al., 2016), the assessment of measurement invariance is equally important in within-country comparisons (Steinmetz et al., 2009). This study meets this importance. However, at the same time, the findings of this study make clear that the results of the assessment of measurement invariance are not static. Citizenship and, as a result, citizenship education is dynamic and subject to changes in society and schools. These changes may lead to changes in how an instrument aligns with its context and underscores the importance of periodic assessment of measurement invariance.