

The Impact of Survey Mode Design and Questionnaire Length on Measurement Quality

Cernat, Alexandru; Sakshaug, Joseph; Christmann, Pablo; Gummer, Tobias

Veröffentlichungsversion / Published Version

Zeitschriftenartikel / journal article

Zur Verfügung gestellt in Kooperation mit / provided in cooperation with:

GESIS - Leibniz-Institut für Sozialwissenschaften

Empfohlene Zitierung / Suggested Citation:

Cernat, A., Sakshaug, J., Christmann, P., & Gummer, T. (2022). The Impact of Survey Mode Design and Questionnaire Length on Measurement Quality. *Sociological Methods & Research*, OnlineFirst, 1-50. <https://doi.org/10.1177/00491241221140139>

Nutzungsbedingungen:

Dieser Text wird unter einer CC BY-NC Lizenz (Namensnennung-Nicht-kommerziell) zur Verfügung gestellt. Nähere Auskünfte zu den CC-Lizenzen finden Sie hier: <https://creativecommons.org/licenses/by-nc/4.0/deed.de>

Terms of use:

This document is made available under a CC BY-NC Licence (Attribution-NonCommercial). For more information see: <https://creativecommons.org/licenses/by-nc/4.0>

The Impact of Survey Mode Design and Questionnaire Length on Measurement Quality

Sociological Methods & Research

1–50

© The Author(s) 2022



Article reuse guidelines:

sagepub.com/journals-permissions

DOI: 10.1177/00491241221140139

journals.sagepub.com/home/sm

Alexandru Cernat¹ , Joseph Sakshaug^{2,3} ,
Pablo Christmann⁴, and Tobias Gummer⁴ 

Abstract

Mixed-mode surveys are popular as they can save costs and maintain (or improve) response rates relative to single-mode surveys. Nevertheless, it is not yet clear how design decisions like survey mode or questionnaire length impact measurement quality. In this study, we compare measurement quality in an experiment of three distinct survey designs implemented in the German sample of the European Values Study: a single-mode face-to-face design, a mixed-mode mail/web design, and a shorter (matrix) questionnaire in the mixed-mode design. We compare measurement quality in different ways, including differences in distributions across several data quality indicators as well as equivalence testing over 140 items in 25 attitudinal scales. We find similar data quality across the survey designs, although the mixed-mode survey shows more item nonresponse compared to the single-mode survey.

¹ Department of Social Statistics, University of Manchester, Manchester, Lancashire, UK

² Department of Statistical Methods, Institut für Arbeitsmarkt- und Berufsforschung, Nürnberg, Bayern, Germany

³ Department of Sociology, Universität Mannheim Fakultät für Sozialwissenschaften, Mannheim, Germany

⁴ Data and Research on Society, GESIS – Leibniz Institute for the Social Sciences, Mannheim, Baden-Württemberg, Germany

Corresponding Author:

Alexandru Cernat, Department of Social Statistics, University of Manchester, Manchester, Lancashire, UK.

Email: alexandru.cernat@manchester.ac.uk

Using equivalence testing we find that most scales achieve metric equivalence and, to a lesser extent, scalar equivalence across the designs.

Keywords

data quality, matrix sampling, mixed-mode survey, measurement invariance, split questionnaire

Introduction

Face-to-face and telephone surveys are currently being challenged by decreasing response rates (Beullens et al., 2018; Brick & Williams, 2013; de Leeuw et al., 2018) and a corresponding rise in costs due to intensified fieldwork efforts (Wolf et al., 2021). These forces have prompted many social surveys to consider moving away from traditional interviewer-administered single-mode designs to mixed-mode designs involving a combination of interviewer- and self-administered (e.g., web) modes, or solely self-administered using either paper-based mail or online modes. For example, longitudinal studies such as the Health and Retirement Study, the UK Household Longitudinal Study (*Understanding Society*), and the National Longitudinal Study of Adolescent to Adult Health (Add Health) have transitioned from fully face-to-face or telephone survey designs to mixed-mode designs involving online and/or mail-based data collection (Bianchi et al., 2017; Biemer et al., 2021; Cernat et al., 2016). Similarly, repeated cross-sectional and cross-national studies, such as the European Social Survey and the European Values Study (EVS) have recently experimented with online and mail-based data collection as an alternative to face-to-face data collection in some countries (Cernat & Revilla, 2020; Luijkx et al., 2020). Most recently, the COVID-19 pandemic has forced many surveys to switch from face-to-face data collection to telephone and self-administered modes (e.g., Gummer et al., 2020; Sakshaug et al., 2020).

Yet, while adopting mixed-mode designs with greater use of self-administration have been shown to be cost-effective (Bianchi et al., 2017; Kappelhof, 2015; Wagner et al., 2014), such a design change can have non-trivial effects on data quality. This is due to the fact that different interview modes have inherently different properties (e.g., level of interviewer presence and aural versus visual stimuli) that can influence respondents' answers. These properties are particularly distinct between interviewer- and self-administered modes, and there is a broad literature showing that the same

groups of respondents can provide different answers to the same items depending on which mode type is used (de Leeuw, 2005). The length of the questionnaire also has data quality implications, especially when a long questionnaire designed for interviewer-administration is applied via self-completion, which may lead to greater respondent burden and poorer data quality for questions positioned later in the questionnaire (Galesic & Bosnjak, 2009). Dropping questions or administering subsets of the questionnaire to different sets of respondents minimizes this risk (Raghunathan & Grizzle, 1995), but may result in a lack of measurement comparability with the full-length version. Thus, there is a need to understand the implications of these design decisions on data quality and data comparability.

Against this backdrop, it is important for practitioners to understand the measurement implications of changing from an interviewer-administered mode design with a long questionnaire to a self-administered mode design with a shorter questionnaire. In this article, we shed light on these issues by investigating differences in measurement quality in the German sample of the EVS. Sampled individuals of the general population were randomized to be interviewed using the traditional face-to-face mode design, a mixed-mode paper/web design, or a mixed-mode paper/web design with a shorter questionnaire implemented as part of a matrix sampling design. Using this experimental setup, it is possible to assess the extent to which measurement quality differs between the mode designs and questionnaire lengths. These comparisons will reveal potential tradeoffs of shifting to a self-administered mode design either as a supplement or as a replacement for traditional face-to-face data collection. The following research questions (RQ) are addressed:

RQ1. What is the impact of mode design and questionnaire length on point estimates and distributions of survey items?

RQ2. To what extent does mode design and questionnaire length affect data quality indicators (e.g., item nonresponse, response style indicators)?

RQ3. Does measurement equivalence hold for attitudinal scales measured under different mode designs and questionnaire lengths?

Previous Research

The choice of data collection mode(s) is a crucial decision for any survey as it affects how respondents receive and answer questions. Different modes have inherently distinct features that affect the presentation of the questions and influences the answers that respondents provide. Consequently, different modes can elicit different answers to the same questions posed to the same respondents (de

Leeuw, 2005). These “measurement mode effects” can arise when some respondents are interviewed in one mode, while other respondents are interviewed in an alternative mode. Measurement mode effects are problematic as they can bias comparisons between respondents who are interviewed in different modes. Such effects can occur when comparing respondents in a mixed-mode survey or in cross-national settings where different countries administer the same questionnaire in different modes. Measurement mode effects tend to be more severe when mixing interviewer-administered (e.g., face-to-face and telephone) and self-administered (e.g., mail and web) modes, rather than mixing within them (Cernat et al., 2016; Klausch et al., 2013). These two mode types differ with respect to two key features which have been attributed to their potential for measurement mode effects: the communication channel (aural or visual) and the presence/absence of an interviewer. Below we describe each feature and its potential impact on measurement mode effects in turn.

The communication channel refers to whether the survey questions are presented visually or orally to respondents. Face-to-face and telephone are primarily aural modes as interviewers read the questions out loud and respondents process the information aurally. In contrast, mail and web are primarily visual modes as respondents read the questions and process the information using visual cues. There are some exceptions. For instance, face-to-face interviews are known to use showcards to supplement the oral presentation (Lynn et al., 2012) and web surveys may use audio components to supplement the visual presentation. Both communication channels are suggested to affect the cognitive process and memory capacity in different ways (Krosnick & Alwin, 1987; Schwarz et al., 1991). Because of this, different communication channels can lead to different response behaviors. For example, aural modes are often associated with recency effects, defined as the tendency to select the last spoken answer categories, which are most easily remembered compared to the earlier spoken answer categories (Smyth et al., 1987). Visual modes are usually associated with primacy effects, or the tendency to endorse the first answer categories, which is in line with the notion that respondents select the first acceptable answer category they see (Krosnick & Alwin (1987). Though the empirical evidence on primacy and recency effects in different modes is inconsistent (de Leeuw, 2005). Other response behaviors, including the level of nondifferentiation (or straightlining) to attitudinal item batteries and the frequency of selecting extreme answers to rating scales, have also been shown to vary between self- and interviewer-administered modes more so than within them (Dillman et al., 2009; Kim et al., 2019).

Regarding the presence/absence of an interviewer, this mode feature can have a strong influence on the answers that respondents provide.

Interviewers are a well-known source of variance inflation (or “interviewer effects”) in survey estimates (West & Blom, 2017), but they also have an important effect on social desirability bias. A consistent research finding is that interviewer-administered modes produce higher levels of socially desirable responding compared to self-administered modes (Cernat et al., 2016; Heerwegh, 2009; Kreuter et al., 2008; Tourangeau & Yan, 2007). However, interviewer presence may reduce the frequency of other suboptimal response behaviors, such as nondifferentiation and item nonresponse (Hope et al., 2014), by keeping respondents motivated and focused on the response task, which is more difficult to do in self-completion modes.

With respect to measurement mode effects in attitudinal items and Likert scales, Klausch et al. (2013) report higher category thresholds for attitudes about police and traffic in the Netherlands in face-to-face and telephone modes compared to mail and web modes, an indication of stronger socially desirable responding in the former modes. Heerwegh & Loosveldt (2011) report evidence of socially desirable responding among telephone responses in a crime victimization survey in Belgium, reflected in more favorable attitudes toward the police, compared to mail responses. Cernat et al. (2016) identified higher levels of depression reported online compared to telephone and face-to-face in the US Health and Retirement Study. Cernat & Revilla (2020) investigated differences in measurement quality between the face-to-face ESS round 8 and the CROSS-National Online Survey (CRONOS) panel. They find higher item nonresponse and higher levels of primacy effects in the CRONOS panel, but similar levels of nondifferentiation with the ESS. Using equivalence testing they find that metric and scalar equivalence holds for four out of the five scales tested.

In addition to measurement mode effects, the length of the questionnaire can also have data quality implications when transitioning from interviewer-administration to self-completion. While it is common to use interviewer modes to administer long questionnaires to respondents, administering long questionnaires via self-completion is likely to require greater effort from the respondent. In the case of web surveys, researchers have suggested the ideal length should not exceed 20 min (Callegaro et al., 2015; Revilla & Ochoa, 2017). Extending the questionnaire beyond this length bears the risk that respondents will engage in undesirable response behaviors toward the end of the interview as fatigue accumulates and they lose focus and motivation to invest the necessary effort to provide high-quality answers. Galesic & Bosnjak (2009) found that items positioned later in the questionnaire of a web survey were answered more quickly, with shorter answers given to open-ended questions, less variability in answers given to grid questions, and

slightly more item nonresponse, compared to the same items positioned earlier in the questionnaire. Peytchev & Peytcheva (2017) also found strong evidence of greater measurement error when items on diet and exercise appeared later in a web survey instrument. Similar findings have been reported in self-administered paper questionnaires (Herzog & Bachman, 1981; Sahlqvist et al., 2011).

One strategy to minimize respondent burden is to shorten the questionnaire by dropping items. However, this strategy may be undesirable from a research perspective or impossible for a country that is participating in a cross-national survey with a fixed questionnaire. An alternative strategy is to implement a split or matrix questionnaire design by partitioning the questionnaire into shorter modules and administering subsets of the modules to different respondents (Raghunathan & Grizzle, 1995). The items that are not answered by all respondents are then imputed to compensate for the data which are missing by design and produce a full rectangular dataset available for multivariate analyses – which is, however, difficult to achieve in general population survey datasets that include hundreds of variables. In an experimental evaluation of the matrix questionnaire design on measurement error reduction, Peytchev & Peytcheva (2017) found that the matrix design performed better than the full questionnaire design on several metrics (e.g., estimates of means, mean squared error, associations) that are susceptible to measurement error. Yet, additional research is needed to understand the effects of using a matrix questionnaire design on measurement quality in multi-item scales, where the measurement properties (e.g., context effects) can vary due to the different questionnaire structures.

In this study, we evaluate the measurement mode effects of shifting from an interviewer-administered (face-to-face) survey to a fully self-administered mixed-mode (web/mail) survey and from a long questionnaire to a matrix questionnaire design by comparing the quality of respondent answers in the German sample of the EVS. Motivated by the above literature, we investigate whether the shift in mode design and questionnaire length affects estimates of means, distributional properties, and indicators of data quality (e.g., item nonresponse, primacy, recency, middle answers, extreme response style, and straightlining) for 140 items and 25 scales measured in the general population. Additionally, we apply confirmatory factor analysis to test the assumption of measurement equivalence for the multi-item scales across the different mode designs and questionnaire lengths.

Data and Experiment

The EVS is a cross-national repeated cross-sectional survey research program examining public opinion on topics like the environment, national identity, perception of life, politics and society, religion and morale, and work. The EVS was first conducted in 1981 with subsequent cross-sectional surveys being fielded every nine years. In the most recent data collection, the EVS allowed participating countries to test the use of self-administered data collection modes; for more details, see Luijkx et al. (2020). In this study, we use the EVS data collected in Germany (EVS, 2017a, 2017b).

Data Collection Experiment in the EVS Germany

The EVS 2017/2018 in Germany featured an interviewer-administered face-to-face (computer-assisted personal interviewing; CAPI) survey and two additional self-administered mixed-mode (web and mail) surveys (Wolf et al., 2021). A probabilistic sample of people residing in Germany aged 18 or older at the time of fieldwork was drawn from population registers. For this purpose, a two-stage sampling design was used. In the first step, municipalities were selected, and primary sampling units were generated. In the second step, addresses for randomly selected persons within the primary sampling units were drawn from population registers. Each sampled unit had the same inclusion probability (i.e., the sample was self-weighting). To allow for an experimental comparison of the data collection modes, postal addresses from the drawn sample were randomly split into three experimental groups (the group sample sizes varied due to project requirements).

The first experimental group (drawn sample size = 5,338) consisted of the traditional face-to-face survey using the full EVS-questionnaire with an average length of 59 min necessary for completion. The fieldwork of the face-to-face survey was organized in two phases, to which the drawn sample was randomly allocated before the beginning of the fieldwork (phase one: $N = 3,762$ eligible addresses, phase two: $N = 1,576$ eligible addresses). In the first phase, a 10€ postpaid incentive was announced in the initial cover letter ($N = 814$ realized interviews). For the second phase (starting in January 2018), a 5€ prepaid incentive was used to increase survey participation in this group ($N = 570$ realized interviews). After an evaluation of the first phase, it was decided to offer a 20€ postpaid incentive to the remaining nonrespondents ($N = 110$ realized interviews).

The second group (drawn sample size = 1,913) consisted of a mixed-mode (web/mail) survey using the same full-length EVS questionnaire

which – compared to the face-to-survey – yielded a comparable, but slightly shorter, average interview duration of 55 min (for the web interview). The third group (drawn sample size = 8,973) consisted of an additional mixed-mode (web/mail) survey using a matrix questionnaire design, implemented to reduce the overall response burden (Peytchev & Peytcheva, 2017; Raghunathan & Grizzle, 1995). The use of this matrix design reduced the average questionnaire length for the self-administered mixed-mode survey to 38 min (for the web interview).

The fieldwork of the face-to-face survey took place between October 2017 and April 2018. Participation in the face-to-face survey was comparatively low with an AAPOR response rate of 6 (The American Association for Public Opinion Research, 2016) of 28%, even though respondents were offered monetary incentives. The fieldwork of the mixed-mode matrix survey was conducted between November 2017 and March 2018 and was implemented in a responsive design with two phases (Gummer et al., 2021; Wolf et al., 2021). Here, the drawn sample was randomly allocated to the two phases. In a first phase, different incentive strategies (5€ prepaid and 10€ postpaid) and mode choice sequences (simultaneous and sequential) were experimentally implemented. Based on an evaluation of the phase one results, in January 2018 the second phase of the matrix survey started for which respondents were provided with a 5€ prepaid incentive in a concurrent mode choice sequence (i.e., offering both the mail questionnaire and the web questionnaire right from the beginning). The response rate of the matrix survey was 36.1%. On average, 73% of the respondents participated via mail mode (27% via web).

Encouraged by this outcome, it was decided to also field the full-length EVS-questionnaire (i.e., the second experimental group) in a self-administered mixed-mode survey with a 5€ prepaid incentive and concurrent mode choice sequence. The fieldwork period for the full mixed-mode survey was between September and November 2018 and resulted in a response rate of 35.3%. In this survey, 83% of the respondents participated via mail mode (17% via web).

Questionnaire Design

When designing the questionnaires for the data collection experiment, the general aim has been to make the questionnaires as comparable as possible while making the necessary mode-specific adaptations.¹ While the formulations of the questions and answer scales remained largely unchanged, question/interviewer instructions were sometimes reformulated to fit the self-administered context. As the EVS lacks complex filtering, question order was the same across designs. Yet, there are also several inherent differences

in the modes such as “don’t know”-options or using grid questions (de Leeuw et al., 2018, pp. 82–83). The EVS experiment optimized the usability of the survey for each mode, so question batteries were presented/read out item-by-item in the face-to-face mode and for mobile devices, while grids were used in the mail mode and on large screens in the web mode.

The second major difference relates to the usage of “don’t know” answer categories. Following the best practice (in Germany), “don’t know” options were not mentioned by the interviewers and were not included in the show cards. Yet, “don’t know” was still introduced as a response alternative for the interviewers if respondents mentioned this option spontaneously. For the web and mail modes, explicit “don’t know” options were provided for most questions.

Attitudinal Scales

We analyzed 140 questions from 21 item batteries with three or more questions, covering all major topics of the EVS. Out of these 21 item batteries, we identified 25 scales through confirmatory factor analysis (see the section on equivalence testing), which we analyze throughout this study. We define a scale here as a group of items that measure the same latent concept. A short description of the scales and the years when they were introduced in the EVS is provided in Table 1. A longer description of the scales, including number of items, scale points, and the use of showcards are given in Table A1 of the Appendix.

Weighting for Selection

We use weights to account for possible differences between the face-to-face, mixed-mode full-length, and mixed-mode matrix samples due to selection. Iterative proportional fitting (Deming & Stephan, 1940) was used to adjust the distribution of the sample to the target population’s distributions with respect to sex, age, education, household size, citizenship, and religion. Reference distributions for the target population were provided by the German Federal Statistical Office.

Methods

Data Quality Indicators

We computed several indicators to assess data quality in each sample – both at the item level and at the scale level. These indicators capture response

Table 1. Description of Scales Used in the Analysis.

Scale	Description	Est. in EVS
Action	Unconventional political participation	1981
Belong	Closeness to the world vis-à-vis their municipality	2017*
Childhood	Political sophistication of parents when they were 14 years old	2008
Concern	Compassion for foreigners	1999
Concern_grp	Compassion for vulnerable groups	1999
Democracy1	Democratic attitudes	2017*
Democracy2	Authoritarian attitudes	2017*
Elections	Assessment of the integrity of elections	2017*
Environment	Importance of environmental protection	1981–2017
European	Meaning of being European	2017
Immigration	Attitudes toward immigration	2008
Importance	Important in life	1990
Marriage	Conditions for a successful marriage	1981
National	Meaning of being German	2008
Norms1	Attitudes towards cheating and probity	1981–2017
Norms2	Liberal versus conservative values	1981–2017
Pol_system	Support for democracy	1999
Pol_watch	Frequency of news consumption	2017
Policy	Attitudes on redistribution and the welfare state	1990
Society	Importance of societal provisions	1990, 2017
Surveillance	Attitudes about state surveillance	2017
Traditional_family	Attitudes towards traditional family roles	2008–2017
Trust	Confidence in public and political institutions	1981–2017
Trust_pl	Interpersonal trust	2017*
Work	Importance of work in society	2008

*Adopted from the World Values Study (WVS) in 2017.

behaviors that could be the result of not completing every step of the cognitive response process (Tourangeau et al., 2000) or due to respondents attempting to reduce their perceived response burden (Krosnick, 1991, 1999).

To measure item nonresponse, we calculated “don’t know” and “no answer” as the proportion of all eligible items for which respondents have chosen one of these residual categories. To measure response style, we estimated indicators for extreme response style, primacy, recency, mid-point response style, and straightlining. Primacy refers to the tendency to endorse the first answer categories, while recency is defined as the tendency to select the last spoken/shown answer categories. We define an extreme

response style as the tendency to use the highest or lowest categories of a response scale – regardless of the content of the question (Van Vaerenbergh & Thomas, 2013), while a mid-point response style is defined as the tendency to select a middle alternative (Narayan & Krosnick, 1996). Straightlining refers to minimal differentiation between responses within item batteries (Roßmann et al., 2018). We employed a simple non-differentiation measure (Kim et al., 2019), which reflects the proportion of questions where only a single response category was chosen for all respective items, ranging from zero to one. Not every response style indicator could be estimated for each scale depending on the number of scale points (3, 4, 5, 10, 11; see Table A2 in the Appendix).

Turning to the outcomes of the three data collection modes, we start our investigation by comparing the distribution of the question and the scales across experimental groups. For the calculation of the scales, we estimated average indices. For comparability, we rescaled all questions and scales to a range from 0 to 1. For the analysis, we compare mean values and standard deviations. We performed unpaired *t*-tests and applied Bonferroni corrections to counteract the problem of the large number of comparisons. If a difference reaches the $p < 0.05$ level we plot this as an additional information alongside the comparisons of means. In a similar manner, we performed variance-ratio tests for the quality of variances, also applying Bonferroni corrections. To gauge the effect size of mean differences we calculated Cohen's *D* for each question and scale and plotted the results in histograms.

Equivalence Testing

In addition to comparing data quality using the indicators mentioned above we also use equivalence testing to identify potential measurement differences across mode designs and questionnaire lengths. Measurement equivalence is useful for two main reasons: first, it can tell us whether we can validly compare results collected in different groups. This is especially useful when a traditional face-to-face survey is switching to a new design and comparisons over time are important. It is also useful when data are collected from different groups using different designs and comparisons or combinations of these data are needed (such as in cross-cultural research). Second, it can indicate differences in data quality and potential causes for these differences.

In order to carry out the equivalence testing we identified items that are part of rating scales with at least three items and that have ordinal response categories. We initially identified 144 items covering 14 topics. We then

ran a series of confirmatory factor analysis (CFA) on groups of items that measure the same concept in order to evaluate their overall fit and ensure they can be used for further testing. In total eight items were excluded at this stage either because they had small loadings or correlations with the other items in the scale (belong/national identify: v168; parents: v270, v274; concern_grp: v219; national: v190; importance: v1, v5, v6). Thus, in the final CFA models, we were left with 132 items divided into 25 scales. The full list of items and the scales they belong to are presented in Table A1 of the Appendix, while the scales and their overall model fits can be seen in Table 2.

Table 2. Fit Statistics for the Scales Analyzed Based on a Simple Confirmatory Factor Analysis Using the Entire Sample.

Scale	# Response Categories	Chi ²	df	P-value	CFI	RMSEA
Policy	10	9.2	5	0.10	1.00	0.02
Trust	4	4598.4	135	0.00	0.93	0.08
Democracy1	10	150.4	5	0.00	0.94	0.09
Democracy2	10	5.1	2	0.08	1.00	0.02
Pol_system	4	10.6	2	0.01	1.00	0.03
Norms1	10	82.4	5	0.00	0.97	0.06
Norms2	10	1210.7	35	0.00	0.88	0.08
Belong	4	629.4	2	0.00	0.92	0.29
Elections	4	985.0	20	0.00	0.93	0.12
Importance*	4	0.00	0	0.00	1.00	0.00
Immigration	10	5.4	2	0.07	1.00	0.02
National	4	199.0	2	0.00	0.98	0.16
European	4	96.5	2	0.00	0.99	0.11
Environment	5	134.0	5	0.00	0.96	0.08
Surveillance*	4	0.0	0	0.00	1.00	0.00
Pol_watch	5	34.1	2	0.00	0.98	0.07
Concern	5	29.03	3	0.00	0.998	0.048
Concern_grp*	5	0.0	0	0.00	1.00	0.00
Society	4	80.7	2	0.00	0.92	0.10
Childhood	4	1089.2	9	0.00	0.89	0.18
Trust_pl	4	1021.6	9	0.00	0.99	0.17
Work	5	34.7	5	0.00	0.99	0.04
Marriage	3	226.0	9	0.00	0.97	0.08
Traditional_family*	5	0.0	0	0.00	1.00	0.00
Action	3	72.5	2	0.00	0.99	0.10

*Models just identified and fit cannot be estimated.

In order to test whether the measurement of these concepts is equivalent across groups (single-mode vs. mixed-mode and mixed-mode long vs. mixed-mode matrix), we compare a series of three multi-group CFA models with increasing restrictions (Baumgartner & Steenkamp, 2006; Davidov et al., 2014):

- **Configural model.** The structure of the CFA is the same between groups but all coefficients are allowed to be different. This model is used as the reference for subsequent models (an example is shown in Figure 1).
- **Metric model.** The loadings (λ in Figure 1) are restricted to be the same across groups. If this model is not significantly worse compared to the previous (configural) model then the variances of the latent variable can be compared across groups.
- **Scalar model.** The intercepts/thresholds (τ in Figure 1) are restricted to be the same across groups. If this model is not significantly worse than the previous one it implies that the means of the latent variables are comparable across groups.

The metric and the scalar models give information about different types of measurement error. The metric model refers to the covariance between the questions and the latent variable. In general, a higher loading (λ in

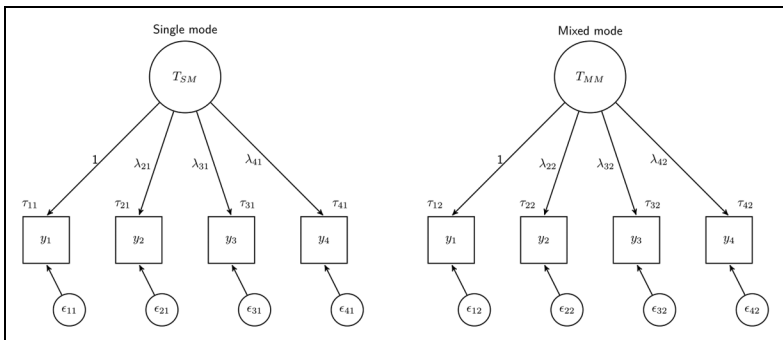


Figure 1. Visual representation of equivalence testing across two groups (single-mode vs. mixed-mode). The configural model is represented where the factor structure is the same but the coefficients are allowed to be different across groups. The large circle represents the overall concept or latent variable, the squares are the observed variables, and the small circles are the residuals. λ denotes the loadings or regression coefficients of the relationship between the latent and observed variables, while τ denotes the intercept/thresholds of these relationships.

Figure 1) highlights a stronger relationship with the latent variable and can be an indication of higher reliability (Bollen, 1989). As a result, lack of metric equivalence could indicate differences in reliability across groups. The scalar model, on the other hand, can be informative regarding systematic shifts in the responses, as it tests if the conditional average of each question is the same across groups. If scalar equivalence is not met it can indicate that the averages of the questions across groups (even after controlling for the scores on the latent variable) are different. This could be caused by systematic response styles like social desirability, acquiescence, and primacy or recency.

When equivalence is not found (metric and scalar models fit worse than the less restrictive model) it is also possible to investigate partial equivalence (Byrne et al., 1989). This implies running a series of models to find which variables can be consistently measured across groups. If at least two variables are found to be equivalent it might still be possible to compare the means and the variances of the latent variables (Byrne et al., 1989).

For all scales, we run simple CFA models for equivalence testing with one latent variable explaining the observed variables (similar to Figure 1). We treat variables with five or more categories as continuous and those with less than five as categorical. All analyses are weighted (see description of selection weights above). For the continuous variables, we use maximum likelihood robust (MLR) estimation while for the categorical ones we use the Weighted Least Squares Means and Variance adjusted estimator with Theta parametrization. We fix the means of the latent variables to 0 to aid estimation and the loading of the first variable to 1.

One model had issues during equivalence testing. For the “traditional family” scale MLR with weights leads to estimation issues, so we use ML with no weights instead. Thus, in total, we investigate 150 models (25 scales * 6 models; 3 models for single- vs. mixed-mode and 3 models for mixed-mode long vs. mixed-mode matrix) for equivalence testing. To compare model fit we use the cut-off value of 0.01 for the change in the comparative fit index (ΔCFI) to define a model as lacking equivalence (Chen, 2007). The data were cleaned in Stata and R, with the data quality indicators created in Stata and the equivalence testing performed in Mplus.

Results

Comparison of Means and Distributions

We start by investigating differences in the estimated means and standard deviations of all 140 items initially used in the CFA models (RQ1), excluding

only the scale religion for reasons of consistency with the following analysis. The first row in Figure 2 shows scatterplots of the standardized means for single versus mixed-mode and mixed-mode long versus mixed-mode matrix (short) designs (for proportion estimates, see Table A3 in the

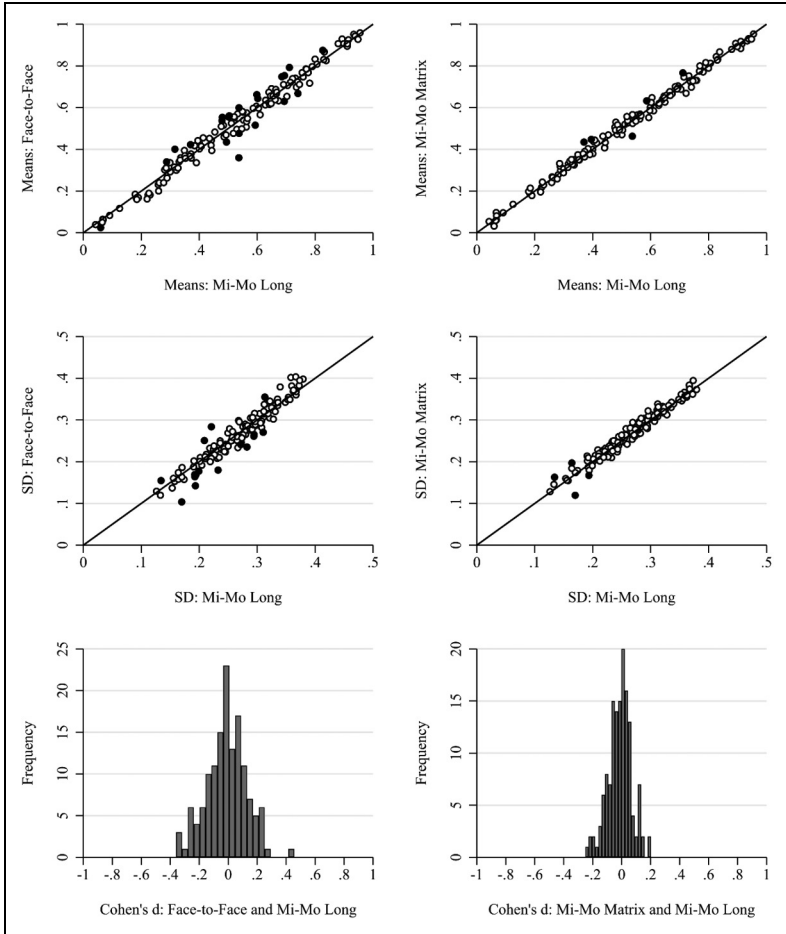


Figure 2. Item level scatterplots for means and standard deviations by survey design: face-to-face single mode, mixed-mode long, and mixed-mode matrix (short). Filled points are statistically significantly different after Bonferroni correction. Bottom row: histogram of Cohen's *d* for differences between groups. The diagonal lines are used to gauge whether values are higher or lower in one group or another.

Appendix). We observe some significant differences especially in the single-versus mixed-mode comparison. This is confirmed in the histogram of Cohen's d (third row in Figure 2). Here we see that while many of the differences are in the $+/- 0.2$ range (indicating small effects) there is a group of variables that have larger averages in face-to-face compared to the mixed-mode (mail/web) design and have moderate effect sizes between 0.2 and 0.6. This seems to be much less the case when comparing the mixed-mode data by questionnaire length.

The comparison of standard deviations (second row in Figure 2) shows fewer significant differences although a few variables seem to have significantly higher standard deviations in the mixed-mode design compared to the face-to-face one. Again, the differences are less pronounced between the mixed-mode long and the mixed-mode matrix surveys. The same patterns emerge when we replicate the analysis at the scale level for all 25 scales in Figure 3. There are some significant differences between the face-to-face mode and the mixed-mode long survey and only few differences between the mixed-mode long and the mixed-mode matrix (short) survey.

After an inspection of the six items with effect sizes larger/smaller than $+/-0.20$ (Cohen's d), we conclude that there is no discernible pattern for the differences between the mixed-mode long and the mixed-mode matrix survey (see Table A4 in the Appendix). On the other hand, a pattern emerges when qualitatively inspecting the 23 items with moderate effect sizes between the face-to-face mode and the mixed-mode long (see Table A5 in the Appendix). For example, the single item with the largest effect size of 0.46 is on how often respondents would follow news on social media, with respondents from the mixed-mode long having a higher social media news diet, on average – which appears plausible. Possibly related, respondents in the mixed-mode long are on average a little bit more skeptical about the integrity of elections in Germany. Conversely, respondents in the face-to-face mode tend to show somewhat higher levels of social and political trust, while also placing more emphasis on their national identity.

We also compare the means and standard deviations for the different survey designs by question topic and scale type (Table A4 in the Appendix), and find some indications of differences. Topics such as national identity, politics and society, and religion and morale, and four-point scales have more differences in the single- versus mixed-mode comparison, but not for the mixed-mode long versus mixed-mode matrix comparison. However, we need to be careful when interpreting these findings because of the relatively small sample size at the item level and the clustering within scales.

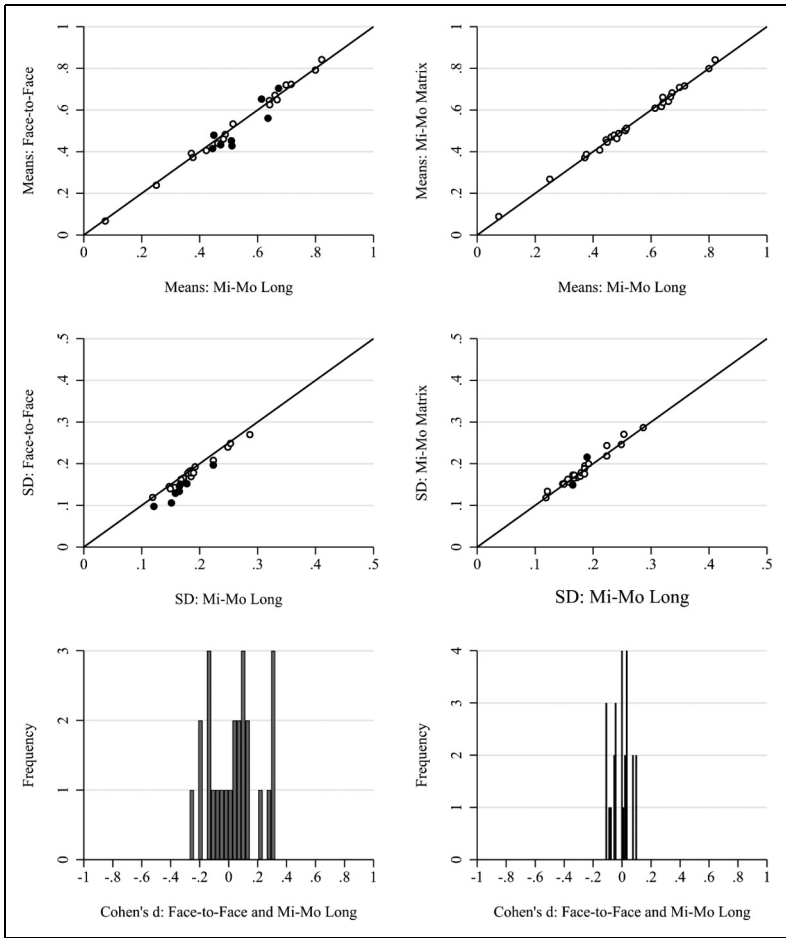


Figure 3. Scale level scatterplots for means and standard deviations by survey design: face-to-face single mode, mixed-mode long, and mixed-mode matrix (short). Filled points are statistically significantly different after Bonferroni correction. Bottom row: histogram of Cohen's *d* for differences between groups.

Comparison of Data Quality Indicators

We next estimate multiple data quality indicators (RQ2). Figure 4 shows the average of each data quality indicator as well as the confidence interval based on all the items separately by the three groups of interest: face-to-face, mixed-

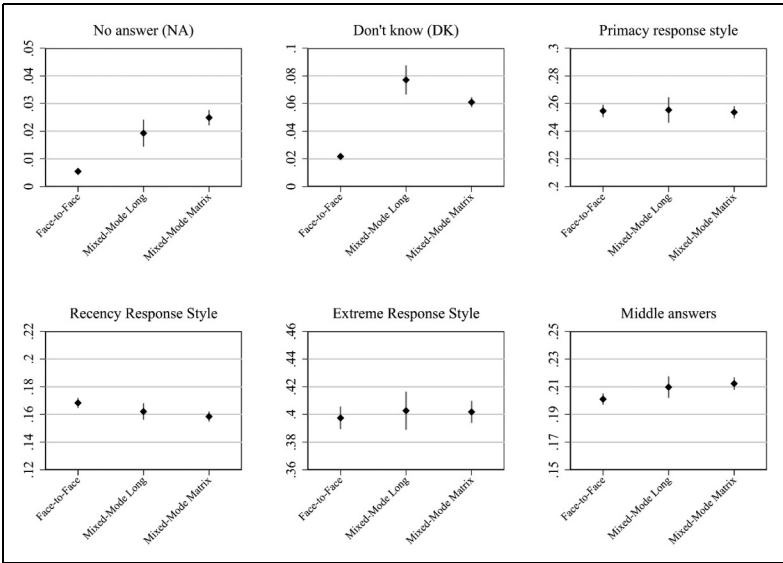


Figure 4. Item level average quality indicators with 95%-confidence intervals by group: face-to-face single mode, mixed-mode long, and mixed-mode matrix (short).

mode long, and mixed-mode matrix. Starting with item nonresponse, we see that the rate is significantly higher in the mixed-mode designs (long and short) compared to the face-to-face design (i.e., confidence intervals do not overlap). This is true both for the “no answer” and “don’t know” responses. The difference between the mixed-mode long and mixed-mode matrix designs is not significant for the “no answer” response. The “don’t know” responses appear to be less frequent in the mixed-mode matrix survey than in the mixed-mode long survey. For the other data quality indicators (primacy, recency, middle answer, and extreme response style), there are only small differences between the face-to-face, mixed-mode long, and mixed-mode matrix surveys. While there appear to be slightly higher levels of recency for the interviewer-administered survey than for the self-administered mixed-mode survey, the picture is reverted for the use of middle answers with lower levels for the face-to-face survey and somewhat higher levels for the mixed-mode surveys.

We next estimated quality indicators for the 25 scales. Figure 5 mirrors the analysis performed in Figure 4, but also adds an indicator for straightlining, which can only be computed at the scale level. Results are very similar to those presented in Figure 4: item nonresponse is significantly more frequent

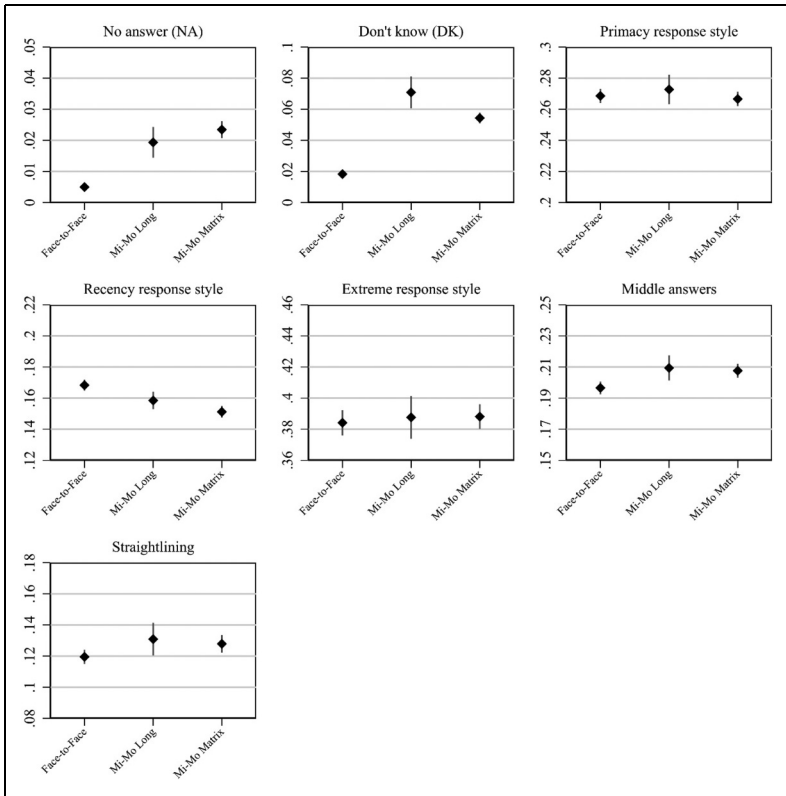


Figure 5. Scale level average quality indicators with 95%-confidence intervals by group: face-to-face single mode, mixed-mode long, and mixed-mode matrix (short).

in both mixed-mode designs than in the face-to-face mode, and the mixed-mode matrix shows a smaller extent of “don’t know” responses than the mixed-mode long. For the remaining data quality indicators differences between the surveys are small, although the face-to-face mode appears to perform slightly better than both mixed-mode designs with respect to straightlining and middle response style, while the mixed-mode surveys perform slightly better with respect to recency.

Lastly, Figure 6 shows the data quality indicators separately for each scale. When turning to this lower level of aggregation we see the results of the previous analysis reconfirmed: the mixed-mode designs consistently show greater proportions of “no answer” and “don’t know” responses across all

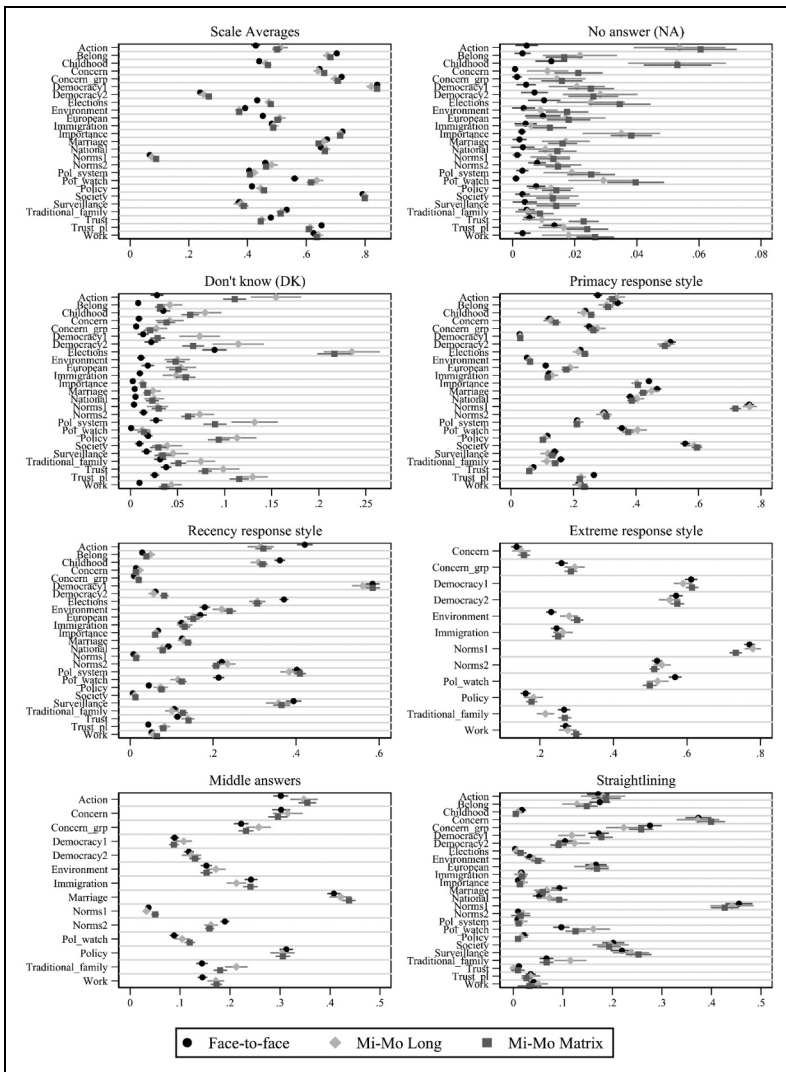


Figure 6. Scale level quality indicators with 95%-confidence intervals by group: face-to-face single mode, mixed-mode long, and mixed-mode matrix (short).

scales. These proportions also vary between the scales and the magnitude of the difference is quite large for several scales. Yet, these differences in item nonresponse seem to be only weakly related to differences in the scale averages (top-left row). Again, we find no indication of a difference in extreme response style between the surveys (bottom-left row). Finally, the mode designs seem to be only weakly and inconsistently related to the use of middle answers, primacy, recency, and straightlining; thus, the results remain consistent between the analyses presented in Figures 5 and 6.

Equivalence Testing: Mode Design Comparison

We next investigate whether the 25 scales identified have comparable measurement models (i.e., measurement equivalence) across face-to-face and mixed-mode designs and mixed-mode long and mixed-mode matrix designs (RQ3). We start by testing whether the three levels of measurement equivalence (configural, metric, and scalar) hold across the single-mode face-to-face survey and the mixed-mode mail/web survey, for each of the 25 scales. Measurement equivalence is assessed by inspecting differences in the CFI values between models. Full details of the model fit statistics are provided in the online supplementary materials. The models for one of the scales (“importance”) did not converge when analyzing it using the multi-group approach. The following results will refer to the remaining 24 scales.

Out of 24 scales analyzed the vast majority (21 out of 24) achieve metric equivalence based on our criterion (see Table 3, “metric” column). This suggests that valid comparisons of unstandardized relationships between the interviewer- and self-administered mode groups can be obtained for most scales. An alternative interpretation of equal loadings is that reliabilities are similar across the groups (Bollen, 1989).

The three scales that do not reach metric equivalence are Norms1, Pol_watch, and Traditional_family. For Norms1, the largest difference appears for the loadings of items v152 (if it is justified: “Cheating on tax if you have the chance”) where the loading is larger in the mixed-mode design and v162 (if police violence is justified) where the reverse is true. For Pol_watch, the loading for v211 (how often you follow politics on social media) has a much smaller loading in the mixed-mode compared to the single-mode design (0.40 vs. 0.07). Similarly, the loadings are less strong in the mixed-mode data compared to the face-to-face data for v83 (“It is a duty towards society to have children”) and v84 (“Adult children have the duty to provide long-term care for their parents”). This leads to a mixed result for these scales although overall the loadings are generally stronger in the

Table 3. Equivalence Testing Conclusions Based on CFI Differences – Mode Design Comparison. Problematic Items Are Denoted in Parentheses.

Scale	Metric	Scalar
Action	✓	✓
Belong	✓	✓
Childhood	✓	✓
Concern	✓	✓
Concern_grp	✓	X (v217, v218, v220)
Democracy1	✓	X (v133, v135, v136, v138, v141)
Democracy2	✓	X (v134, v137, v139, v140)
Elections	✓	✓
Environment	✓	X (v199–v203)
European	✓	✓
Immigration	✓	✓
Marriage	✓	✓
National	✓	✓
Norms1	X (v150, v152, v159, v162)	✓
Norms2	✓	X (v151, v153–v158, v160, v161, v163)
Pol_system	✓	✓
Pol_watch	X (v209–v211)	X (v208–v211)
Policy	✓	X (v103–v107)
Society	✓	✓
Surveillance	✓	✓
Traditional_family	X (v83, v84)	X (v82–v84)
Trust	✓	✓
Trust_pl	✓	✓
Work	✓	X (v46–v50)
TOTAL	21/24	15/24

face-to-face compared to the mixed-mode design, indicating slightly higher reliability in the interviewer-administered mode for these items.

Scalar equivalence is also achieved for the majority (15 out of 24) of the scales (see Table 3, “scalar” column), indicating that it is valid to compare the latent means between mode designs. Out of the nine scalar non-equivalent scales, two of them overlap with the metric non-equivalent scales: *Pol_watch* and *Traditional_family*. These two scales are particularly problematic as they have different loadings and different intercepts and therefore are not measured in the same way across the mode designs. The other scalar non-equivalent scales include *Concern_grp*, both democracy scales,

environment, Norms2, Policy, and Work. We dissect some of the causes of scalar non-equivalence in these scales below.

Scalar equivalence is essential for ensuring that averages of the latent variables are comparable across groups. Additionally, differences in intercepts/thresholds can be informative as they may be caused by systematic biases such as social desirability or acquiescence. The differences found in the nine scales that are not scalar equivalent show mixed patterns, with larger intercept/threshold values in both mode design groups for particular items. For example, for items v136 (“it’s essential for democracy for people to receive state aid” where larger numbers mean more essential) and v211 (“how often you follow politics on social media” where larger numbers mean less often) the intercepts are larger in the mixed-mode design. In contrast, for items v156 (“Euthanasia can be justified” where larger numbers mean can be justified) and v154 (“Abortion can be justified” where larger numbers mean can be justified) the conditional mean is larger in the face-to-face single-mode design. At this point we can only speculate regarding possible causes for such differences, such as social desirability, leading to systematic shifts in the observed averages of these variables. What is clear is that systematic differences in the intercepts can happen across mode designs.

Equivalence Testing: Questionnaire Length Comparison

Next, we test for measurement equivalence across the long- and short-questionnaire groups used in the mixed-mode (web/mail) surveys. The same model evaluation criterion (i.e., change in CFI values) is used. Full model details can be found in the online supplementary materials.

Metric equivalence is achieved for all but two scales (policy and traditional_family) (see Table 4, “metric” column). Again, this means that one can be relatively assured that unstandardized relationships can be validly compared between both long- and short-questionnaire versions. The full-length questionnaire version shows stronger loadings for items v83 (“It is a duty towards society to have children”), v84 (“Adult children have the duty to provide long-term care for their parents”), v107 (“Private ownership of business and industry should be increased”), and v105 (“Competition is good”), while the shorter questionnaire version has stronger loadings for item v106 (“Incomes should be made more equal”). As mentioned above, this indicates that for some items the full-length questionnaire has higher reliability compared to the shorter length version, although this difference appears only for a small proportion of items.

Scalar equivalence is established for all but five scales (see Table 4, “Scalar” column), supporting the possibility of making valid comparisons of latent

Table 4. Equivalence Testing Conclusions Based on CFI Differences – Questionnaire Length Comparison. Problematic Items Are Denoted in Parentheses.

Scale	Metric	Scalar
Action	✓	✓
Belong	✓	✓
Childhood	✓	✓
Concern	✓	✓
Concern_grp	✓	✓
Democracy1	✓	X (v133, v135, v136, v138, v141)
Democracy2	✓	✓
Elections	✓	✓
Environment	✓	✓
European	✓	✓
Immigration	✓	✓
Marriage	✓	✓
National	✓	✓
Norms1	✓	X (v149, v150, v152, v159, v162)
Norms2	✓	✓
Pol_system	✓	✓
Pol_watch	✓	X (v208–v211)
Policy	X (v104–v107)	X (v103–v107)
Society	✓	✓
Surveillance	✓	✓
Traditional_family	X (v83, v84)	✓
Trust	✓	✓
Trust_pl	✓	✓
Work	✓	X (v46–v50)
TOTAL	22/24	19/24

means between different questionnaire lengths. Among the non-equivalent scales, policy, also lacks metric equivalence, suggesting that this scale cannot be measured in the same way in both long- and short-questionnaire forms. The other scalar non-equivalent scales include Democracy1, Norms1, Pol_watch, and work.

Looking at the largest differences in the intercepts for these non-equivalent scales we find that items v136 (“it’s essential for democracy for people to receive state aid” where larger numbers mean essential), v104 (“People who are unemployed should have the right to refuse a job they do not want” where larger numbers mean more support for statement), v211 (“how often you follow politics on social media” where larger numbers mean less often), v133 (“Governments tax the rich and subsidize the poor”

where higher numbers mean essential for democracy), v149 (“Claiming state benefits which you are not entitled to” where higher numbers mean always justified), and v138 (“Civil rights protect people from state oppression” where higher numbers mean essential for democracy) have higher intercepts in the shorter questionnaire compared to the longer one. Once again, a number of different types of systematic biases could lead to these patterns such as social desirability or recency effects.

Discussion

In this study, we investigated two important aspects of modern survey data collection: differences in measurement quality between a single-mode face-to-face design and a fully self-administered mixed-mode (mail/web) design, and differences in measurement quality between a long and short questionnaire implemented in a self-administered mixed-mode (mail/web) design. Overall, we found small differences in the distributions of variables across these different designs (RQ1), although we observed larger differences in estimated means and lower variation for items measured in the single-mode face-to-face design compared to the mixed-mode (mail/web) designs. We also observed more item missing data in the mixed-mode designs compared to the single-mode design but few systematic differences in data quality at the item level (RQ2). Moreover, there were only few differences at the item level when comparing the mixed-mode surveys with the different questionnaire lengths.

We next investigated measurement equivalence (RQ3), or the degree to which the different survey designs have the same measurement model, in 24 attitudinal scales (covering 129 items). When comparing mode designs, we found that 21 out of 24 scales achieve metric equivalence and 15 out of 24 achieve scalar equivalence. When comparing questionnaire lengths, we found that 22 out of 24 scales reach metric equivalence and 19 out of 24 reach scalar equivalence. Thus, the majority of scales allow for the comparison of means and variances of latent variables across the different survey designs.

We believe these findings carry a positive message that implies relatively similar data quality and comparable data between face-to-face and mail/web mode designs and between long and short self-completion questionnaire versions. In this context, we view the implementation of a different mode design, namely shifting from an interviewer-administered design to a self-administered one, to be more problematic than shifting from a longer to a shorter questionnaire length. This is largely consistent with previous literature

showing larger differences in measurement quality between interviewer and self-administered modes (Cernat et al., 2016; Klausch et al., 2013).

There were some indications that a single-mode face-to-face design has slightly higher reliabilities than a self-administered mixed-mode design, and that a full-length self-completion questionnaire has higher reliabilities compared to a matrix questionnaire for some items. We also found that systematic differences in the conditional means are often present. One can only speculate on their potential causes and implications for data analysis. This is a topic for future research. That being said, the differences between the two mode designs are relatively small, which may also be due to the extensive use of showcards in the face-to-face mode.

We point out some limitations of this study. Namely, we examined the EVS data for a single country. Examining the measurement quality of cross-national comparisons under different modes and questionnaire designs is a worthwhile topic for future research. Furthermore, alternative approaches to testing measurement equivalence and data quality could be considered. In the present study, we found differences in response rates between the surveys with the face-to-face survey having the lowest response rate. For a more detailed description of differences in sample composition between the designs, we refer to the study by Wolf et al. (2021). Based on our research questions, we focused on analyses of differences in measurement and utilized weighting to control for selective participation. However, future research could go one step further and try to feature both the measurement and the nonresponse simultaneously. Similarly, the weighting correction could be extended to take into account complex sample designs. Further, the EVS data did not include all the possible design options, for example, a short face-to-face survey, although we do believe the utility and prevalence of such a design would be limited in practice. Finally, we have also implemented incentive experiments (5€ prepaid vs. 10€ postpaid) in the EVS-Germany and further research could explore measurement implications between different monetary incentive types (conditional vs. unconditional) as existing research usually compares the use of incentives against no incentive.

Nonetheless, the present study undertook a comprehensive analysis of more than 100 items from more than 20 scales in one of the largest international surveys used in the social sciences. Furthermore, the study investigated data quality in a number of different ways: comparing distributional properties at the item and scale level, comparing various data quality indicators at the item and scale level, and measurement equivalence testing. Moreover, comparing both the mode design and questionnaire length brings important insights for survey practice. Often, shifting from single-mode or interviewer-administered data collection to self-completion, especially web-based, necessitates further

adaptations that can be important. For example, web and mobile data collection may have questionnaire length or interview duration constraints that can impact how the survey is adapted. As such, the switch from a long face-to-face survey to a shorter web survey seems to be especially prevalent. Our findings can help survey practitioners to assess which design alterations are likely to result in more measurement differences compared to the previous design and which aspects of data quality may be affected.

Future research should aim to replicate our findings in other countries, other questionnaires, and other survey designs. Of special interest would be to understand the extent of the measurement differences between short and long questionnaires and to investigate potential causes for such differences. Also of interest for data users would be to study the kinds of biases they should expect if measurement equivalence (metric and/or scalar) does not hold across survey designs. This is especially relevant for longitudinal surveys (panel or repeated cross-sectional studies) that switch from face-to-face to self-completion or surveys where different population groups are interviewed using different mode designs (e.g., different countries in cross-cultural research).

Author's Note

The replication code for this paper can be found online at: https://github.com/alex-cernat/evs_mm_DE_equivalence. The data is publicly available through the GESIS data archive: ZA7500_v3-0-0 (<https://doi.org/10.4232/1.13899>) and ZA7502_v1-0-0 (<https://doi.org/10.4232/1.13092>) (the full citation can be found in the references). Joseph Sakshaug, Statistical Methods Research Department, Institute for Employment Research, Nuremberg, Germany, Department of Statistics, Ludwig Maximilian University of Munich, Munich, Germany.


Declaration of Conflicting Interests


The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.


Funding

The author(s) received no financial support for the research, authorship, and/or publication of this article.

ORCID iDs

Alexandru Cernat  <https://orcid.org/0000-0003-2176-1215>

Joseph Sakshaug  <https://orcid.org/0000-0001-7520-353X>

Tobias Gummer  <https://orcid.org/0000-0001-6469-7802>

Supplemental Material

Supplemental material for this article is available online.

Note

1. All variants of the German questionnaires are available for download in the country-specific documentation of the EVS.

References

- Baumgartner, H. and J.-B. E. M. Steenkamp. 2006. "An Extended Paradigm for Measurement Analysis of Marketing Constructs Applicable to Panel Data." *Journal of Marketing Research (JMR)* 43(3):431-42.
- Beullens, K., G. Loosveldt, C. Vandenplas, and I. Stoop. 2018. "Response Rates in the European Social Survey: Increasing, Decreasing, or a Matter of Fieldwork Efforts?" *Survey Methods: Insights from the Field (SMIF)*. doi:<https://doi.org/10.13094/SMIF-2018-00003>
- Bianchi, A., S. Biffignandi, and P. Lynn. 2017. "Web-Face-to-Face Mixed-Mode Design in a Longitudinal Survey: Effects on Participation Rates, Sample Composition, and Costs." *Journal of Official Statistics* 33(2):385-408.
- Biemer, P. P., K. M. Harris, D. Liao, B. J. Burke, and C. T. Halpern. 2021 "Modeling Mode Effects for a Panel Survey in Transition." Pp. 63–88. in *Measurement Error in Longitudinal Data*, edited by A. Cernat and J. W. Sakshaug. Oxford: Oxford University Press.
- Bollen, K. 1989. *Structural Equations with Latent Variables*. New York: Wiley-Interscience Publication.
- Brick, J. M. and D. Williams. 2013. "Explaining Rising Nonresponse Rates in Cross-Sectional Surveys." *The ANNALS of the American Academy of Political and Social Science* 645(1):36-59. <https://doi.org/10.1177/0002716212456834>
- Byrne, B. M., R. J. Shavelson, and B. Muthén. 1989. "Testing for the Equivalence of Factor Covariance and Mean Structures: The Issue of Partial Measurement Invariance." *Psychological Bulletin* 105(3):456.
- Callegaro, M., K. L. Manfreda, and V. Vehovar. 2015. *Web Survey Methodology*. Los Angeles: Sage.

- Cernat, A., M. P. Couper, and M. B. Ofstedal. 2016. "Estimation of Mode Effects in the Health and Retirement Study Using Measurement Models." *Journal of Survey Statistics and Methodology* 4(4):501-24. <https://doi.org/10.1093/jssam/smw021>
- Cernat, A. and M. Revilla. 2020. "Moving from Face-to-Face to a Web Panel: Impacts on Measurement Quality." *Journal of Survey Statistics and Methodology* smaa007. doi:<https://doi.org/10.1093/jssam/smaa007>
- Chen, F. F. 2007. "Sensitivity of Goodness of Fit Indexes to Lack of Measurement Invariance." *Structural Equation Modeling: A Multidisciplinary Journal* 14(3):464-504. <https://doi.org/10.1080/10705510701301834>
- Davidov, E., B. Meuleman, J. Cieciuch, P. Schmidt, and J. Billiet. 2014. "Measurement Equivalence in Cross-National Research." *Annual Review of Sociology* 40(1):55-75. <https://doi.org/10.1146/annurev-soc-071913-043137>
- de Leeuw, E. D. 2005. "To Mix or Not to Mix Data Collection Modes in Surveys." *Journal of Official Statistics* 21(2):233-55.
- de Leeuw, E., J. Hox, and A. Luiten. 2018. "International Nonresponse Trends Across Countries and Years: An Analysis of 36 Years of Labour Force Survey Data." *Survey Methods: Insights from the Field*: 1-11. <https://doi.org/10.13094/SMIF-2018-00008>
- Deming, W. E. and F. F. Stephan. 1940. "On a Least Squares Adjustment of a Sampled Frequency Table When the Expected Marginal Totals are Known." *The Annals of Mathematical Statistics* 11(4):427-44.
- Dillman, D. A., G. Phelps, R. Tortora, K. Swift, J. Kohrell, J. Berck, and B. L. Messer. 2009. "Response Rate and Measurement Differences in Mixed-Mode Surveys Using Mail, Telephone, Interactive Voice Response (IVR) and the Internet." *Social Science Research* 38(1):1-18. <https://doi.org/10.1016/j.ssresearch.2008.03.007>
- EVS. 2017a. *European Values Study 2017: Integrated Dataset – Matrix Design Data (EVS 2017 Matrix Design)*. GESIS Data Archive, Cologne. ZA7502 Data file Version 1.0.0, doi:10.4232/1.13092.
- EVS. 2017b. *European Values Study 2017: Integrated Dataset (EVS 2017)*. GESIS Data Archive, Cologne. ZA7500 Data file Version 3.0.0, doi:10.4232/1.13511.
- Galesic, M. and M. Bosnjak. 2009. "Effects of Questionnaire Length on Participation and Indicators of Response Quality in a Web Survey." *Public Opinion Quarterly* 73(2):349-60. <https://doi.org/10.1093/poq/nfp031>
- Gummer, T., P. Christmann, S. Verhoeven, and C. Wolf. 2022. "Using a responsive survey design to innovate self-administered mixed-mode surveys." *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 185(3):916–32. <https://doi.org/10.1111/rssa.12835>.
- Gummer, T., C. Schmiedeborg, M. Bujard, P. Christmann, K. Hank, T. Kunz, D. Lück, and F. J. Neyer. 2020. "The Impact of Covid-19 on Fieldwork Efforts

- and Planning in Pairfam and FReDA-GGS.” *Survey Research Methods* 14(2):223-7. <https://doi.org/10.18148/srm/2020.v14i2.7740>
- Heerwegh, D. 2009. “Mode Differences Between Face-to-Face and Web Surveys: An Experimental Investigation of Data Quality and Social Desirability Effects.” *International Journal of Public Opinion Research* 21(1):111-21. <https://doi.org/10.1093/ijpor/edn054>
- Heerwegh, D. and G. Loosveldt. 2011. “Assessing Mode Effects in a National Crime Victimization Survey Using Structural Equation Models: Social Desirability Bias and Acquiescence.” *Journal of Official Statistics* 27(1):49.
- Herzog, A. R. and J. G. Bachman. 1981. “Effects of Questionnaire Length on Response Quality.” *Public Opinion Quarterly* 45(4):549-59. <https://doi.org/10.1086/268687>
- Hope, S., P. C. Campanelli, G. Nicolaas, P. Lynn, and J. Annette. 2014. *The role of the interviewer in producing mode effects: Results from a mixed modes experiment comparing face-to-face, telephone and web administration* (ISER Working Paper Series No. 2014–20). Institute for Social and Economic Research. <https://econpapers.repec.org/paper/eseiserwp/2014-20.htm>.
- Kappelhof, J. W. 2015. “Face-to-Face or Sequential Mixed-Mode Surveys among Non-Western Minorities in The Netherlands: the Effect of Different Survey Designs on the Possibility of Nonresponse Bias.” *Journal of Official Statistics* 31(1):1-30.
- Kim, Y., J. Dykema, J. Stevenson, P. Black, and D. P. Moberg. 2019. “Straightlining: Overview of Measurement, Comparison of Indicators, and Effects in Mail–Web Mixed-Mode Surveys.” *Social Science Computer Review* 37(2):214-33. <https://doi.org/10.1177/0894439317752406>
- Klausch, T., J. J. Hox, and B. Schouten. 2013. “Measurement Effects of Survey Mode on the Equivalence of Attitudinal Rating Scale Questions.” *Sociological Methods & Research* 42(3):227-63. <https://doi.org/10.1177/0049124113500480>
- Kreuter, F., S. Presser, and R. Tourangeau. 2008. “Social Desirability Bias in CATI, IVR, and Web Surveys: The Effects of Mode and Question Sensitivity.” *Public Opinion Quarterly* 72(5):847-65. <https://doi.org/10.1093/poq/nfn063>
- Krosnick, J. A. 1991. “Response Strategies for Coping with the Cognitive Demands of Attitude Measures in Surveys.” *Applied Cognitive Psychology* 5(3):213-36. <https://doi.org/10.1002/acp.2350050305>
- Krosnick, J. A. 1999. “Survey Research.” *Annual Review of Psychology* 50(1):537-67.
- Krosnick, J. A. and D. F. Alwin. 1987. “An Evaluation of a Cognitive Theory of Response-Order Effects in Survey Measurement.” *Public Opinion Quarterly* 51(2):201-19. <https://doi.org/10.1086/269029>
- Luijckx, R., G. A. Jónsdóttir, T. Gummer, M. Ernst Stähli, M. Frederiksen, K. Ketola, T. Reeskens, E. Brislinger, P. Christmann, SÞ Gunnarsson, ÁB Hjaltason, D. Joye, V. Lomazzi, A. M. Maineri, P. Milbert, M. Ochsner, A. Pollien, M. Sapin, and

- I. Solanes, ... C. Wolf. 2020. "The European Values Study 2017: On the Way to the Future Using Mixed-Modes." *European Sociological Review* jcaa049. doi: <https://doi.org/10.1093/esr/jcaa049>
- Lynn, P., S. Hope, A. Jäckle, P. Campanelli, and G. Nicolaas. 2012. *Effects of visual and aural communication of categorical response options on answers to survey questions* (Working Paper No. 2012–21). ISER Working Paper Series. <https://www.econstor.eu/handle/10419/65994>.
- Narayan, S. and J. A. Krosnick. 1996. "Education Moderates Some Response Effects in Attitude Measurement." *Public Opinion Quarterly* 60(1):58-88.
- Peytchev, A. and E. Peytcheva. 2017. "Reduction of Measurement Error due to Survey Length: Evaluation of the Split Questionnaire Design Approach." *Survey Research Methods* 11(4):361-8. <https://doi.org/10.18148/srm/2017.v11i4.7145>
- Raghunathan, T. E. and J. E. Grizzle. 1995. "A Split Questionnaire Survey Design." *Journal of the American Statistical Association* 90(429):54-63. <https://doi.org/10.2307/2291129>
- Revilla, M. and C. Ochoa. 2017. "Ideal and Maximum Length for a Web Survey." *International Journal of Market Research* 59(5):557-65. <https://doi.org/10.2501/IJMR-2017-039>
- Roßmann, J., T. Gummer, and H. Silber. 2018. "Mitigating Satisficing in Cognitively Demanding Grid Questions: Evidence from two Web-Based Experiments." *Journal of Survey Statistics and Methodology* 6(3):376-400.
- Sahlgvist, S., Y. Song, F. Bull, E. Adams, J. Preston, and D. Ogilvie. 2011. "Effect of Questionnaire Length, Personalisation and Reminder Type on Response Rate to a Complex Postal Survey: Randomised Controlled Trial." *BMC Medical Research Methodology* 11(1):1-8.
- Sakshaug, J. W., J. Beste, M. Coban, T. Fendel, G.-C. Haas, S. Hülle, Y. Kosyakova, C. König, F. Kreuter, B. Kürfner, B. Müller, C. Osiander, S. Schwanhäuser, G. Stephan, E. Vallizadeh, M. Volkert, C. Wenzig, C. Westermeier, C. Zabel, and S. Zins. 2020. "Impacts of the COVID-19 Pandemic on Labor Market Surveys at the German Institute for Employment Research." *Survey Research Methods* 14(2):229-33. <https://doi.org/10.18148/srm/2020.v14i2.7743>
- Schwarz, N., F. Strack, H.-J. Hippler, and G. Bishop. 1991. "The Impact of Administration Mode on Response Effects in Survey Measurement." *Applied Cognitive Psychology* 5(3):193-212. <https://doi.org/10.1002/acp.2350050304>
- Smyth, M. M., P. E. Morris, P. Levy, and A. W. Ellis. 1987. *Cognition in Action* (pp. xiii, 346). New Jersey: Lawrence Erlbaum Associates, Inc.
- The American Association for Public Opinion Research. 2016. *Standard Definitions: Final Dispositions of Case Codes and Outcome Rates for Surveys. 9th Edition*. AAPOR.

- Tourangeau, R., L. J. Rips, and K. Rasinski. 2000. *The Psychology of Survey Response*. Cambridge: Cambridge University Press.
- Tourangeau, R. and T. Yan. 2007. "Sensitive Questions in Surveys." *Psychological Bulletin* 133(5):859-83. <https://doi.org/10.1037/0033-2909.133.5.859>
- Van Vaerenbergh, Y. and T. D. Thomas. 2013. "Response Styles in Survey Research: A Literature Review of Antecedents, Consequences, and Remedies." *International Journal of Public Opinion Research* 25(2):195-217.
- Wagner, J., J. Arrieta, H. Guyer, and M. B. Ofstedal. 2014. "Does Sequence Matter in Multimode Surveys: Results from an Experiment." *Field Methods* 26(2):141-55. <https://doi.org/10.1177/1525822X13491863>
- West, B. T. and A. G. Blom. 2017. "Explaining Interviewer Effects: A Research Synthesis." *Journal of Survey Statistics and Methodology* 5(2):175-211. <https://doi.org/10.1093/jssam/smw024>
- Wolf, C., P. Christmann, T. Gummer, C. Schnaudt, and S. Verhoeven. 2021. "Conducting General Social Surveys as Self-Administered Mixed-Mode Surveys." *Public Opinion Quarterly* 85(2):623-648. <https://doi.org/10.1093/poq/nfab039>

Author Biographies

Alexandru Cernat is an associate professor in social statistics at the University of Manchester.

Joseph Sakshaug is a professor at the Institute for Employment Research, Nuremberg, Ludwig Maximilian University of Munich, and the University of Mannheim, Germany.

Pablo Christmann is a senior researcher at GESIS - Leibniz Institute for the Social Sciences, Mannheim, Germany.

Tobias Gummer is a team leader and senior researcher at GESIS - Leibniz Institute for the Social Sciences, Mannheim, Germany.

Appendix

Table A1. Overview of Scales.

Scale	Thematic Block	Items	# Items	# Scale Points	CAPI Page	Mail Page	Split Module	Showcards
Action	Politics and society	v98–v101	4	3	16	18	D	Yes
Belong	National identity	v164–v167	4	4	25	30	B	Yes
Childhood	Respondent's parents	v267–v269, v271–v273	6	4	54	60	A	Yes
Concern	Politics and society	v212–v216	5	5	32	39	C	Yes
Concern_ggp	Politics and society	v217, v218, v220	3	5	32	40	C	Yes
Democracy1	Politics and society	v133, v135, v136, v138, v141	5	11	21	26	D	Yes

(continued)

Table A1. Continued

Scale	Thematic Block	Items	# Items	# Scale Points	CAPI Page	Mail Page	Split Module	Showcards
Democracy2	Politics and society	v134, v137, v139, v140	4	11	21	26	D	Yes
Elections	Politics and society	v176–v183	8	4	27	32	D	Yes
Environment European	Environment National	v199–v203 v194–v197	5 4	5 4	30 29	36 35	D B	Yes Yes
Immigration	identity National	v185–v188	4	10	28	33	C	Yes
Importance	Perception of life	v2, v3, v4	3	4	3	4	Core	Yes
Marriage National	Family National	v65–v70 v189, v191–v193	6 4	3 4	13 29	15 34	A B	No Yes
NormsI	Religion and morale	v149, v150, v152, v159, v162	5	10	23	29	C	Yes

(continued)

Table A1. Continued

Scale	Thematic Block	Items	# Items	# Scale Points	CAPI Page	Mail Page	Split Module	Showcards
Norms2	Religion and morale	v151, v153–v158, v160, v161, v163	10	10	23	29	B-C (B: v151, v153–v155, v157, v158, v160, v161; C: v156, v163)	Yes
Pol_system	Politics and society	v145–v148	4	4	23	28	D	Yes
Pol_watch	Politics and society	v208–v211	4	5	32	39	D	Yes
Policy	Politics and society	v103–v107	5	10	16	19	C	Yes
Society	Politics and society	v221–v224	4	4	33	40	D	Yes
Surveillance	Politics and society	v205–v207	3	4	31	38	D	Yes
Traditional_family	Family	v82–v84	3	5	14	17	A	Yes
Trust	Politics and society	v115–v132	18	4	19	25	Core-B-C (Core: v115, v118, v121, v128; B: v117, v120, v122, v124, v126, v129, v131; C: v116, v119, v123, v125, v127, v130, v132)	Yes

(continued)

Table A1. Continued

Scale	Thematic Block	Items	# Items	# Scale Points	CAPI Page	Mail Page	Split Module	Showcards
Trust_pl	National identity	v32-v37	6	4	6	7	C	Yes
Work	Work	v46-v50	5	5	8	9	A-B (A: v46, v48, v50; B: v47, v49)	Yes

Notes: Showcards were used for all scales in CAPI interviews. A total of 274 items were collected. Matrix questionnaires included core module + two out of four modules (A to D). Modules were randomized as follows: Core + A/B; Core + A/C; Core + A/D; Core + B/C; Core + B/D; Core + C/D.

Table A2. Availability of Response Style Indicators by Scale.

Scale	# Scale Points	Middle			Recency	Extreme Response Style	Straightlining
		# Scale Points	Answers	Primacy			
Action	3		Included	Included		Included	
Belong	4			Included		Included	
Childhood	4			Included		Included	
Concern	5		Included	Included	Included	Included	
Concern_grp	5		Included	Included	Included	Included	
Democracy1	11		Included	Included	Included	Included	
Democracy2	11		Included	Included	Included	Included	
Elections	4			Included		Included	
Environment	5		Included	Included	Included	Included	
European	4			Included		Included	
Immigration	10		Included	Included	Included	Included	
Importance	4			Included		Included	
Marriage	3		Included	Included		Included	
National	4			Included		Included	
Norms1	10		Included	Included	Included	Included	
Norms2	10		Included	Included	Included	Included	
Pol_system	4			Included		Included	
Pol_watch	5		Included	Included	Included	Included	
Policy	10		Included	Included	Included	Included	
Society	4			Included		Included	
Surveillance	4			Included		Included	
Traditional_family	5		Included	Included	Included	Included	

(continued)

Table A2. Continued

Scale	# Scale Points	Middle Answers	Primacy	Recency	Extreme Response Style	Straightlining
Trust	4		Included	Included		Included
Trust_pl	4		Included	Included		Included
Work	5	Included	Included	Included	Included	Included

Table A3. Means, Standard Deviations (SDs) and %-of Lowest Numerical Categories for Unstandardized Variables (Original Variables), Excluding Missing Values.

Item Name	Scale Points	Face-to-Face			Mixed-Mode Long			Mixed-Mode Matrix		
		Means	SD	%-Lowest Category	Means	SD	%-Lowest Category	Means	SD	%-Lowest Category
v1	4	1.68	0.85	50.54	1.67	0.66	39.72	1.67	0.68	42.08
v2	4	1.12	0.39	89.48	1.13	0.38	86.32	1.14	0.38	86.25
v3	4	1.50	0.61	56.40	1.60	0.57	43.60	1.63	0.61	43.13
v4	4	1.70	0.63	40.23	1.67	0.61	39.35	1.67	0.60	39.52
v5	4	2.27	0.81	16.11	2.31	0.75	11.02	2.29	0.75	12.61
v6	4	2.70	1.00	13.16	2.80	1.01	11.48	2.80	0.99	10.72
v32	4	1.15	0.41	87.51	1.20	0.46	81.07	1.24	0.48	78.51
v33	4	2.01	0.65	18.22	2.20	0.64	6.90	2.19	0.60	7.83
v34	4	1.74	0.56	31.11	1.92	0.51	17.75	1.89	0.52	19.59
v35	4	2.80	0.69	1.17	3.05	0.71	0.51	3.01	0.68	0.14
v36	4	2.33	0.68	6.47	2.50	0.73	3.06	2.59	0.76	3.56
v37	4	2.32	0.67	7.14	2.49	0.74	3.85	2.52	0.73	3.59
v46	5	2.01	1.00	33.83	1.97	0.84	27.77	2.08	0.91	26.80
v47	5	2.74	1.22	16.59	2.56	1.20	19.65	2.66	1.25	18.98
v48	5	2.47	1.14	19.30	2.30	1.10	22.92	2.33	1.17	24.83
v49	5	2.03	0.93	26.49	2.07	0.93	25.51	1.99	0.99	32.72
v50	5	3.23	1.17	7.74	3.23	1.18	6.28	3.36	1.12	5.65
v65	3	1.15	0.38	85.11	1.17	0.41	81.71	1.19	0.43	81.61
v66	3	1.97	0.71	25.05	2.03	0.71	20.49	2.06	0.69	20.06

(continued)

Table A3. Continued

Item Name	Scale Points	Face-to-Face			Mixed-Mode Long			Mixed-Mode Matrix		
		Means	SD	%-Lowest Category	Means	SD	%-Lowest Category	Means	SD	%-Lowest Category
v67	3	1.77	0.61	31.36	1.74	0.58	30.15	1.78	0.61	30.95
v68	3	1.92	0.70	29.07	1.89	0.67	27.47	1.87	0.68	28.79
v69	3	1.52	0.72	58.99	1.56	0.73	54.06	1.66	0.76	49.24
v70	3	1.62	0.61	45.50	1.70	0.64	39.70	1.75	0.63	36.69
v82	5	2.31	1.23	31.36	2.44	1.24	24.87	2.44	1.30	28.07
v83	5	3.39	1.19	6.59	3.40	1.07	4.20	3.48	1.13	4.40
v84	5	2.85	1.18	9.39	2.92	1.08	5.25	2.90	1.18	8.46
v98	3	1.57	0.76	64.15	1.44	0.68	70.61	1.46	0.68	69.64
v99	3	2.40	0.69	13.68	2.36	0.64	10.87	2.31	0.67	13.16
v100	3	1.99	0.80	36.26	1.88	0.72	36.98	1.91	0.73	36.03
v101	3	2.67	0.55	4.77	2.64	0.56	4.52	2.62	0.58	5.40
v103	10	4.91	2.42	9.05	5.45	2.52	6.96	5.53	2.67	8.97
v104	10	4.78	2.39	9.85	4.63	2.52	12.48	4.93	2.42	9.24
v105	10	3.36	2.09	24.34	3.60	1.96	16.88	3.69	2.05	17.23
v106	10	5.46	2.48	7.46	5.77	2.58	7.29	5.94	2.63	7.37
v107	10	5.20	2.04	4.23	5.38	2.20	4.83	5.40	2.15	5.51
v115	4	2.82	0.84	4.94	2.95	0.76	2.60	2.88	0.79	3.64
v116	4	2.47	0.75	7.11	2.57	0.71	2.99	2.48	0.73	5.17
v117	4	2.42	0.72	7.52	2.47	0.73	5.55	2.44	0.67	4.37
v118	4	2.75	0.75	3.79	2.85	0.77	2.82	2.79	0.74	3.10

(continued)

Table A3. Continued

Item Name	Scale Points	Face-to-Face			Mixed-Mode Long			Mixed-Mode Matrix		
		Means	SD	%-Lowest Category	Means	SD	%-Lowest Category	Means	SD	%-Lowest Category
v119	4	2.55	0.72	3.94	2.64	0.74	4.03	2.59	0.76	4.79
v120	4	1.94	0.69	23.38	2.01	0.67	17.79	1.98	0.66	19.72
v121	4	2.67	0.76	4.84	2.81	0.78	3.36	2.66	0.76	4.76
v122	4	2.39	0.65	4.74	2.56	0.67	2.45	2.52	0.70	3.97
v123	4	2.34	0.71	8.32	2.44	0.74	6.69	2.42	0.74	7.44
v124	4	2.63	0.76	4.67	2.76	0.81	4.61	2.71	0.79	4.25
v125	4	2.64	0.74	4.37	2.69	0.80	5.46	2.61	0.76	5.50
v126	4	2.27	0.72	10.37	2.35	0.70	7.76	2.35	0.75	8.88
v127	4	2.29	0.76	12.75	2.40	0.79	9.67	2.40	0.81	11.12
v128	4	2.92	0.68	1.61	2.99	0.69	1.72	3.02	0.71	1.36
v129	4	2.33	0.71	8.84	2.44	0.73	6.87	2.39	0.78	9.65
v130	4	2.98	0.63	0.84	3.14	0.64	0.66	3.04	0.67	0.91
v131	4	2.73	0.75	3.52	2.89	0.77	2.42	2.70	0.77	3.74
v132	4	3.12	0.71	1.15	3.10	0.68	1.31	3.08	0.69	0.56
v133	10	7.49	2.61	5.11	7.43	2.51	3.08	7.72	2.40	3.03
v134	10	1.75	1.69	75.65	1.82	1.71	76.44	1.86	1.92	76.51
v135	10	9.36	1.62	1.15	9.04	2.09	1.71	9.17	1.89	1.67
v136	10	8.14	2.17	1.35	7.40	2.45	2.20	7.90	2.36	2.51
v137	10	2.51	2.50	62.20	2.71	2.54	59.53	2.60	2.54	60.55
v138	10	8.43	2.40	3.18	8.54	2.33	2.62	8.81	2.18	3.19

(continued)

Table A3. Continued

Item Name	Scale Points	Face-to-Face			Mixed-Mode Long			Mixed-Mode Matrix		
		Means	SD	%-Lowest Category	Means	SD	%-Lowest Category	Means	SD	%-Lowest Category
v139	10	5.29	3.08	17.80	5.66	3.14	18.56	6.02	3.14	16.11
v140	10	2.66	2.45	57.26	2.61	2.32	55.08	2.78	2.51	56.32
v141	10	9.49	1.47	0.74	9.44	1.73	2.16	9.37	1.89	2.55
v145	4	3.31	0.85	2.71	3.22	0.92	4.40	3.28	0.94	5.19
v146	4	2.74	0.86	5.41	2.68	0.89	7.45	2.78	0.91	6.84
v147	4	3.82	0.46	0.47	3.80	0.47	0.64	3.81	0.46	0.20
v148	4	1.28	0.51	75.93	1.36	0.58	70.03	1.36	0.56	70.10
v149	10	1.59	1.42	75.44	1.60	1.48	76.39	1.87	1.77	69.55
v150	10	1.57	1.42	76.09	1.61	1.57	74.01	1.73	1.61	70.64
v151	10	2.44	2.50	62.06	2.66	2.60	59.12	2.93	2.70	51.88
v152	10	1.35	1.08	82.95	1.38	1.20	82.67	1.49	1.31	81.12
v153	10	7.51	3.14	8.84	7.66	3.18	9.34	7.48	3.25	8.81
v154	10	5.64	3.08	14.42	6.34	3.02	9.85	5.97	3.08	11.18
v155	10	7.28	2.83	4.60	7.62	2.81	4.55	7.56	2.77	4.33
v156	10	7.01	3.04	8.97	7.67	2.82	6.36	7.28	2.94	7.88
v157	10	3.84	3.14	37.89	3.83	3.21	37.31	3.82	3.16	36.90
v158	10	4.28	3.32	33.38	4.18	3.26	33.17	4.14	3.19	33.77
v159	10	2.05	1.87	61.05	2.13	2.08	62.50	2.23	2.10	60.17
v160	10	4.07	2.97	30.95	3.98	2.92	29.22	3.96	2.90	31.38
v161	10	6.89	3.09	9.47	7.32	2.91	7.26	7.21	2.97	7.88

(continued)

Table A3. Continued

Item Name	Scale Points	Face-to-Face			Mixed-Mode Long			Mixed-Mode Matrix		
		Means	SD	%-Lowest Category	Means	SD	%-Lowest Category	Means	SD	%-Lowest Category
v162	10	1.45	1.28	80.05	1.59	1.74	81.66	1.53	1.50	81.78
v163	10	2.63	2.72	65.18	3.01	2.95	58.97	2.96	2.81	58.62
v164	4	1.70	0.73	43.83	1.73	0.71	38.91	1.81	0.75	36.49
v165	4	1.76	0.74	39.88	1.94	0.82	29.82	1.90	0.78	32.39
v166	4	1.64	0.64	44.06	1.70	0.69	40.99	1.68	0.66	43.12
v167	4	2.03	0.75	25.03	2.06	0.80	24.18	2.03	0.80	25.84
v168	4	2.34	0.83	16.80	2.56	0.87	13.76	2.41	0.85	13.87
v176	4	1.38	0.63	70.80	1.52	0.70	62.20	1.47	0.68	65.51
v177	4	3.46	0.70	0.99	3.31	0.80	1.94	3.31	0.80	1.92
v178	4	2.62	0.95	10.48	2.54	0.96	12.96	2.45	0.96	15.74
v179	4	3.51	0.71	1.27	3.34	0.85	2.74	3.41	0.80	2.54
v180	4	2.09	0.79	20.34	2.11	0.80	19.68	2.14	0.81	20.26
v181	4	1.40	0.62	67.62	1.50	0.63	58.06	1.51	0.66	61.71
v182	4	3.40	0.81	2.70	3.22	0.93	5.40	3.20	0.95	6.62
v183	4	3.93	0.31	0.14	3.82	0.51	1.07	3.90	0.36	0.31
v185	10	7.36	2.60	2.85	7.34	2.64	2.98	7.19	2.60	2.84
v186	10	4.40	2.64	15.95	4.35	2.76	19.05	4.27	2.67	17.64
v187	10	4.14	2.60	19.89	4.01	2.67	22.58	4.16	2.71	20.90
v188	10	5.49	2.46	6.88	5.92	2.51	5.83	5.88	2.48	5.10
v189	4	2.59	0.95	15.72	2.52	0.91	15.77	2.52	0.92	15.91

(continued)

Table A3. Continued

Item Name	Scale Points	Face-to-Face			Mixed-Mode Long			Mixed-Mode Matrix		
		Means	SD	%-Lowest Category	Means	SD	%-Lowest Category	Means	SD	%-Lowest Category
v190	4	1.32	0.52	71.28	1.27	0.49	76.26	1.35	0.55	70.12
v191	4	2.93	0.86	7.00	2.88	0.84	7.29	2.83	0.88	9.49
v192	4	1.28	0.49	72.89	1.30	0.51	70.53	1.32	0.52	70.33
v193	4	2.15	0.79	19.46	2.07	0.83	24.15	2.07	0.83	25.79
v194	4	2.57	0.89	12.95	2.39	0.92	18.11	2.43	0.92	16.32
v195	4	2.82	0.78	5.75	2.67	0.88	10.11	2.70	0.86	8.86
v196	4	3.07	0.80	4.46	3.02	0.88	6.79	2.95	0.90	8.31
v197	4	2.11	0.78	21.55	1.92	0.81	32.76	1.92	0.79	32.51
v199	5	2.67	1.18	13.65	2.80	1.19	11.62	2.86	1.24	12.08
v200	5	3.53	1.04	2.43	3.53	1.06	2.81	3.48	1.12	4.02
v201	5	3.66	1.03	2.09	3.80	1.03	2.59	3.85	1.06	2.07
v202	5	3.52	1.17	5.46	3.60	1.22	6.32	3.60	1.24	5.96
v203	5	3.79	1.03	1.58	3.88	1.09	3.80	3.89	1.11	2.95
v205	4	2.18	1.01	28.93	2.21	0.94	23.60	2.24	1.02	27.49
v206	4	3.07	0.95	6.87	3.14	0.89	4.81	3.00	0.97	7.24
v207	4	3.43	0.86	4.36	3.32	0.85	3.91	3.27	0.91	4.99
v208	5	2.15	1.38	50.70	2.02	1.29	55.22	2.09	1.28	51.17
v209	5	2.45	1.53	45.55	2.35	1.44	45.58	2.34	1.43	45.29
v210	5	2.87	1.59	35.21	2.75	1.52	37.27	2.72	1.49	38.23
v211	5	3.56	1.53	17.60	2.85	1.49	26.79	3.15	1.45	20.07

(continued)

Table A3. Continued

Item Name	Scale Points	Face-to-Face			Mixed-Mode Long			Mixed-Mode Matrix		
		Means	SD	%-Lowest Category	Means	SD	%-Lowest Category	Means	SD	%-Lowest Category
v212	5	2.16	0.85	20.30	2.24	0.84	16.51	2.24	0.85	17.96
v213	5	2.37	0.86	11.84	2.38	0.89	13.35	2.36	0.87	12.64
v214	5	2.34	0.80	11.23	2.34	0.80	12.56	2.29	0.83	13.66
v215	5	2.61	0.87	6.92	2.65	0.91	7.22	2.47	0.84	10.25
v216	5	2.61	0.93	9.91	2.58	0.97	10.74	2.41	0.92	14.40
v217	5	1.68	0.67	42.12	1.68	0.78	47.81	1.69	0.72	43.70
v218	5	2.42	0.83	10.69	2.59	0.88	9.49	2.50	0.90	11.49
v219	5	2.60	0.90	8.10	2.85	0.97	6.78	2.73	0.98	8.36
v220	5	1.77	0.71	36.98	1.77	0.80	41.82	1.75	0.76	40.63
v221	4	1.98	0.71	23.67	2.00	0.75	23.21	2.03	0.79	25.07
v222	4	1.28	0.48	73.50	1.27	0.47	73.83	1.25	0.47	76.67
v223	4	1.61	0.60	44.45	1.63	0.66	43.98	1.55	0.63	51.98
v224	4	1.22	0.46	79.22	1.16	0.40	84.45	1.21	0.49	81.04
v267	4	2.36	1.21	37.62	2.30	1.10	30.53	2.31	1.15	33.99
v268	4	3.28	0.96	9.02	3.21	0.93	6.95	3.20	0.95	8.21
v269	4	2.14	1.07	37.06	2.06	0.94	33.11	2.06	0.99	35.99
v270	4	2.99	1.10	12.79	2.83	1.10	16.09	2.87	1.12	15.55
v271	4	2.75	1.18	24.83	2.78	1.12	19.40	2.67	1.18	25.59
v272	4	2.93	1.11	16.74	2.94	1.09	14.70	2.95	1.08	14.23
v273	4	1.65	0.92	61.00	1.69	0.88	55.57	1.60	0.85	59.83
v274	4	3.28	0.96	6.07	3.17	0.99	8.82	3.23	1.00	7.99

Table A4. Comparison of Items in Mixed-Mode Long versus Mixed-Mode Matrix with Effect Sizes Larger/Smaller than $+/-0.20$ (Cohen's d).

Item Name	Scale Points	Mixed-Mode Long			Mixed-Mode Matrix			Cohen's d	Scale	Thematic Block
		Means	SD	%-Lowest Category	Means	SD	%-Lowest Category			
v121	4	2.81	0.78	3.36	2.66	0.76	4.76	-0.20	Trust	Politics and society
v131	4	2.89	0.77	2.42	2.70	0.77	3.74	-0.25	Trust	Politics and society
v136	10	7.40	2.45	2.20	7.90	2.36	2.51	-0.21	Democracy	Politics and society
v183	4	3.82	0.51	1.07	3.90	0.36	0.31	0.20	Elections	Politics and society
v211	5	2.85	1.49	26.79	3.15	1.45	20.07	0.20	Pol_watch	Politics and society
v215	5	2.65	0.91	7.22	2.47	0.84	10.25	-0.21	Concern	Politics and society

Table A5. Comparison of Items in Face-to-Face versus Mixed-Mode Long with Effect Sizes Larger/Smaller than ± 0.20 (Cohen's d).

Item Name	Face-to-Face			Mixed-Mode Long			Cohen's d	Scale	Thematic Block
	Scale Points	Means	SD	%-Lowest Category	Means	SD			
v33	4	2.01	0.65	18.22	2.20	0.64	6.90	Trust_pl	National identity
v34	4	1.74	0.56	31.11	1.92	0.51	17.75	Trust_pl	National identity
v35	4	2.80	0.69	1.17	3.05	0.71	0.51	Trust_pl	National identity
v36	4	2.33	0.68	6.47	2.50	0.73	3.06	Trust_pl	National identity
v37	4	2.32	0.67	7.14	2.49	0.74	3.85	Trust_pl	National identity
v103	10	4.91	2.42	9.05	5.45	2.52	6.96	Policy	Politics and society
v122	4	2.39	0.65	4.74	2.56	0.67	2.45	Trust	Politics and society
v130	4	2.98	0.63	0.84	3.14	0.64	0.66	Trust	Politics and society
v131	4	2.73	0.75	3.52	2.89	0.77	2.42	Trust	Politics and society
v136	10	8.14	2.17	1.35	7.40	2.45	2.20	Trust	Politics and society

(continued)

Table A5. Continued

Item Name	Face-to-Face		Mixed-Mode Long				Cohen's <i>d</i>	Scale	Thematic Block
	Scale Points	Means	SD	%-Lowest Category	Means	SD			
v154	10	5.64	3.08	14.42	6.34	3.02	9.85	Norms2	Religion and morale
v156	10	7.01	3.04	8.97	7.67	2.82	6.36	Norms2	Religion and morale
v165	4	1.76	0.74	39.88	1.94	0.82	29.82	Belong	National identity
v168	4	2.34	0.83	16.80	2.56	0.87	13.76	Belong	National identity
v176	4	1.38	0.63	70.80	1.52	0.70	62.20	Elections	Politics and society
v177	4	3.46	0.70	0.99	3.31	0.80	1.94	Elections	Politics and society
v179	4	3.51	0.71	1.27	3.34	0.85	2.74	Elections	Politics and society
v182	4	3.40	0.81	2.70	3.22	0.93	5.40	Elections	Politics and society
v183	4	3.93	0.31	0.14	3.82	0.51	1.07	Elections	Politics and society
v194	4	2.57	0.89	12.95	2.39	0.92	18.11	European	National identity

(continued)

Table A5. Continued

Item Name	Face-to-Face		Mixed-Mode Long		Cohen's <i>d</i>	Scale	Thematic Block		
	Scale Points	Means	SD	%-Lowest Category				Means	SD
v197	4	2.11	0.78	21.55	1.92	0.81	32.76	European	National identity
v211	5	3.56	1.53	17.60	2.85	1.49	26.79	PoI_watch	Politics and society
v219	5	2.60	0.90	8.10	2.85	0.97	6.78	Concern_grp	Politics and society

Table A6. Percentage of Significant Differences in Average Item Level Means and Standard Deviations (SDs) by Survey Design and by Thematic Block and Scale Type.

	Face-to-Face versus Mixed-Mode Long		Mixed-Mode Long versus Mixed-Mode Matrix		N
	Means: %-Significant Differences*	SD: %-Significant Differences [†]	Means: %-Significant Differences*	SD: %-Significant Differences [†]	
<i>By Thematic Block</i>					
Environment	0.00%	0.00%	0.00%	0.00%	5
Family	0.00%	0.00%	0.00%	0.00%	9
National identity	37.50%	4.17%	0.00%	0.00%	24
Perception of life	0.00%	16.67%	0.00%	0.00%	6
Politics and society	14.71%	14.71%	7.35%	2.94%	68
Religion and morale	13.33%	6.67%	0.00%	13.33%	15
Respondent's parents	0.00%	12.50%	0.00%	0.00%	8
Work	0.00%	20.00%	0.00%	0.00%	5
<i>By scale type</i>					
3 point scale	0.00%	0.00%	0.00%	0.00%	10
4 point scale	19.72%	11.27%	2.82%	2.82%	71
5 point scale	11.54%	11.54%	7.69%	0.00%	26
10 point scale	12.12%	12.12%	3.03%	6.06%	33
Total	15.00%	10.71%	3.57%	2.86%	140

Notes: Proportion of significant differences between the means by survey mode and question type.

*Significance test of means based on unpaired t-tests with Bonferroni corrections.

[†]Significance test of distributions based on variance-ratio tests with Bonferroni corrections.