

1 CALY-SWE value set: An integrated approach for a valuation study  
2 based on an online-administered TTO and DCE survey

3

4

5 *Authors:*

6 Kaspar Walter Meili (Corresponding author),  
7 Department of Epidemiology and Global Health  
8 Umeå university, Sweden  
9 [kaspar.meili@umu.se](mailto:kaspar.meili@umu.se), [kaspar.meili@yahoo.de](mailto:kaspar.meili@yahoo.de)  
10 Orcid: 0000-0002-9889-4406

11 Brendan Mulhern,  
12 Centre for Health Economics Research and Evaluation  
13 University of Technology Sydney, Australia  
14 Orcid: 0000-0003-3656-8063

15 Richard Ssegonja,  
16 Department of Public Health and Caring Sciences  
17 Uppsala University, Sweden  
18 Department of Medical Sciences,  
19 Respiratory, Allergy and Sleep medicine research unit,  
20 Uppsala University, Sweden  
21 Orcid: 0000-0002-5323-5626

22 Jan Hjelte,  
23 Department of Social work  
24 Umeå university, Sweden  
25 Orcid: 0000-0002-5269-1961

26 Kerstin Edin  
27 Department of Epidemiology and Global Health  
28 Umeå university, Sweden

29 Fredrik Norström,  
30 Department of Epidemiology and Global Health  
31 Umeå university, Sweden  
32 Orcid: 0000-0002-0457-2175

33 Inna Feldman,  
34 Department of Public Health and Caring Sciences  
35 Uppsala University, Sweden  
36 Department of Epidemiology and Global Health

37 Umeå university, Sweden  
38 Orcid: 0000-0003-3329-6066  
39 Anna Månsdotter,  
40 Department of Epidemiology and Global Health  
41 Umeå university, Sweden  
42 Lars Lindholm  
43 Department of Epidemiology and Global Health  
44 Umeå university, Sweden  
45 Orcid: 0000-0002-1633-2179  
46

## 47 **Abstract**

48

49 *Purpose:* To explore and develop methods to use TTO and DCE data collected in self-  
50 administered online surveys to elicit a CALY-SWE value set.

51

52 *Methods:* Building on existing methodological knowledge around DCE and TTO studies,  
53 we optimized the web survey in an integrated approach that consisted of a qualitative  
54 face validity study, iterative web survey development including a three-stage roll-out, a  
55 customized experimental design, a sample size simulation. Based on the inconsistencies  
56 of TTO answers per participant, we assessed TTO data quality by calculating a score,  
57 and examined the effect of excluding TTO data according to this score on the value set  
58 modelling.

59

60 *Results:* Participants in the quality study informed improvements in the survey's visual  
61 design and phrasing. Based on the sample size simulation, we judged a sample size of  
62 1500 with a balance of six DCE and five TTO questions to be appropriate for the  
63 valuation study. Change made for the second stage, for example the introduction of a  
64 *learning* state and of color-coding, improved TTO data quality. Excluding TTO answers  
65 per-participant based on the score lead to an improved TTO data foundation for the  
66 model by several metrics such as no inconsistent coefficients and reduced standard  
67 errors relative to the coefficient magnitude.

68

69 *Conclusion:* Developing value sets was feasible with online administered DCE and TTO  
70 questions if the web survey is sufficiently optimized and coherent with the experimental  
71 design. The severity of inconsistencies could be used to identify and exclude poor  
72 quality TTO data to strengthen the value set modelling.

73

74

75

## 76 **Plain English summary**

77 CALY-SWE is a new instrument for measuring quality of life broadly and for use in economic cost-  
78 effectiveness evaluations of social welfare interventions. To be used as such, a value set is needed to  
79 assign quality weights to the states of life used in the economic evaluations. However, no  
80 methodology exists yet for CALY-SWE to derive a value set, and existing similar valuation studies  
81 from the health context use costly person-to-person interviews.

82 In this study, we aimed to find appropriate methods for deriving a CALY-SWE value set using an  
83 online survey.

84 We developed an online web survey with discrete choice experiment (DCE) and time trade-off (TTO)  
85 questions. We used previous methodological knowledge from similar studies, qualitative interview, a  
86 statistical simulation for the sample size, and we rolled-out the survey in three stages to be able to  
87 implement further improvements. We also sought a way to identify and exclude participants who  
88 contributed poorer quality TTO answers to improve the underlying data for the value set.

89 The resulting survey was optimised for online administration with a shorter survey length of six DCE  
90 and five TTO questions, compensated by an increased sample size of 1500. We developed a score  
91 based on illogical TTO answers that allows to identify and exclude poorer quality TTO data. This  
92 design made it possible to perform the CALY-SWE valuation study.

93

94 *Keywords:* quality-adjusted life year, time trade-off, discrete choice experiment, capability approach,  
95 online survey, economic evaluation

96

97 *Word count:* 4000 = 4123 – 46 (f3) – 54 (f2) – 23 (f1)

98

## 99 Introduction

100 CALY-SWE (capability-adjusted life years Sweden) is a new measure for quality of life  
101 purposed for use in economic evaluations with broad social consequences [1, 2]. CALY-  
102 SWE conceptually relates to the quality-adjusted life year (QALY) concept developed  
103 within health economics [1, 2], and a value set is needed to use the instrument in  
104 economic cost-effectiveness evaluations [3]. Value sets consists of quality scores, or  
105 weights, on the [0,1] scale for all *states* that the instrument describes. CALY-SWE consists  
106 of 6 attributes (health, social relations, financial situation and housing, occupation,  
107 political and civil rights, and security) where each has 3 levels (Do not agree, Partially  
108 agree, Agree completely) , equalling 729 possible states [2]. Zero corresponds to a  
109 quality of life equivalent to death, and 1 to a quality of life sufficient for a flourishing life  
110 [4].

111 Traditionally in health economics, standard gamble and time trade-off (TTO) questions  
112 have been widely used for measuring health state values. More recently, discrete choice  
113 experiment (DCE) questions have also been widely adopted in the development of value  
114 sets [3].

115 Both TTO and DCE questions provide complementary information on preferences. In  
116 DCE questions, participants ordinally compare two states with each other, providing  
117 information on the relative strength of levels and attributes. Results are, however, only  
118 anchored relative to each other and not on the absolute [0,1] scale. TTO questions  
119 evaluate a single state by comparing time spent in the best state with time in an  
120 impaired state. The time in the best state is then gradually changes until the participant  
121 states that both situations are equivalent. The quality weight of the impaired state  
122 equals the time spent in the impaired state divided by the time spent in the best state.  
123 TTO questions thus yield direct information on the absolute anchoring of states on the  
124 [0,1] scale but are cognitively challenging [5, 6]. Recently, hybrid models [7, 8] have been  
125 developed that can jointly estimate value sets from DCE and TTO data, making it  
126 possibly to integrate the complementary preference information of DCE and TTO.

127 EQ-5D is a widely used preference-based measure where country specific value are  
128 derived with the EQ-VT protocol [9] for person-to-person interviews. The EQ-VT protocol  
129 has been refined over time and provides routines for both DCE and TTO questions,  
130 defines their experimental design, the interview procedures and in parallel provides  
131 software for doing the interviews. The protocol also focuses on interview training and  
132 data quality monitoring per interviewer to increase data quality, for example it defines a  
133 minimum time to spend in the introductory TTO question. However, the EQ-VT is  
134 designed for person-to person interviews and unsupervised self-administered surveys  
135 are not intended. Person-to-person interviews increase the costs and requirements for  
136 sampling and data collection considerably compared to traditional self-administered  
137 surveys, thus decreasing the overall feasibility. This especially applies in the context of  
138 COVID-19 restrictions which happened to coincide with the CALY-SWE valuation study.

139 On the other hand, some evidence suggests that face-to-face supervised administration  
140 leads to better data quality compared to unsupervised online administration for TTO  
141 questions [10]. Data quality issues related to the challenging TTO format may also be  
142 exacerbated by using a commercial panel of participants together with online self-  
143 administration.

144 For the CALY-SWE valuation study we decided to rely on an online panel and a self-  
145 administrated survey with TTO and DCE questions as we judged it to be the most  
146 feasible way to time- and cost- effectively collect data representative of the Swedish  
147 population. However, given the evidence on issues related to online self-administration,  
148 we implemented and tested multiple elements to adapt the methods and increase their  
149 validity.

## 150 **Aim**

151 Our aim was to explore and develop methods for eliciting a CALY-SWE value set, based  
152 on online data collection via a web panel, and by leveraging methodological experience  
153 accumulated from existing studies using both TTO and DCE.

154 To overcome the challenging nature of data quality in online surveys, especially for TTO  
155 questions, we adopted an integrated approach over several areas, which translated into  
156 the specific goals of:

157 1) Optimizing face validity, usability and quality of the web survey and TTO questions to  
158 ease understanding and increase engagement of the participants, using an iterative  
159 approach consisting of a qualitative study followed by the staged roll-out of the main  
160 survey

161 2) Developing an experimental design optimized for online data collection that is suited  
162 for a shorter survey length to maintain engagement of participants and to ensure cost-  
163 efficient sampling

164 3) Performing a sample size calculation to determine a sufficient sample size large  
165 enough to exclude poorer quality TTO data or to use only TTO data for generating the  
166 final value set

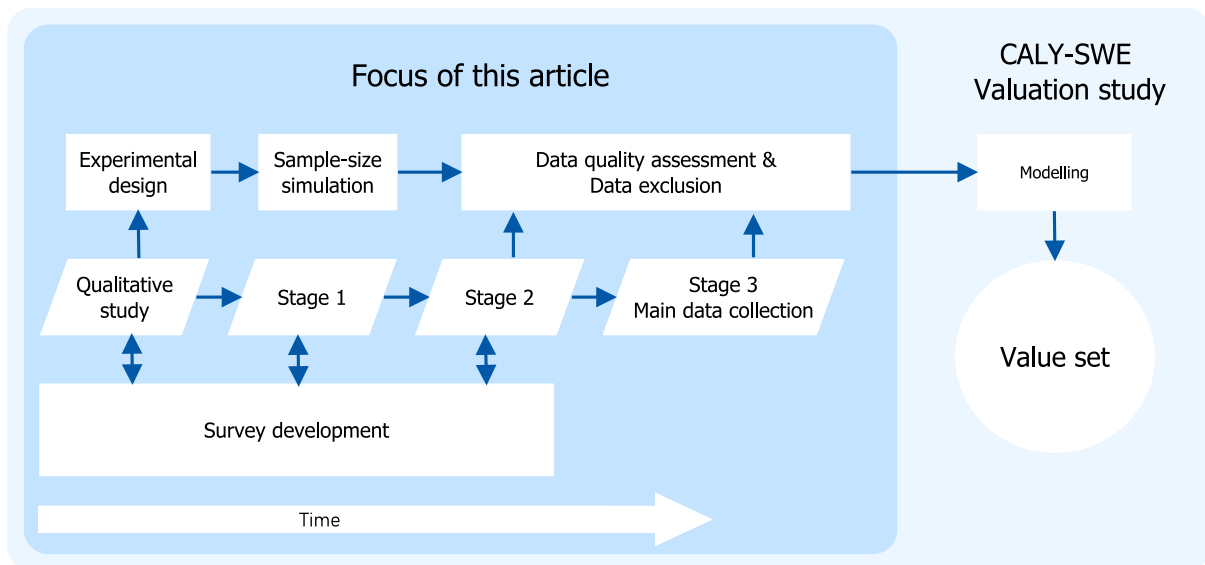
167 4) Investigating ways of assessing TTO data quality and then excluding TTO data with  
168 lower quality

169

## 170 Methods and survey development

### 171 Web survey development

172 We implemented the survey in HTML, CSS, JavaScript, PHP[11], twig [12], mariadb [13],  
173 and sqlite [14], featuring a mobile-first responsive design. This framework enabled us to  
174 customize the web survey to a larger degree than existing commercial solutions while  
175 adhering to Swedish data privacy and safety regulations. After the qualitative face-  
176 validity survey, we rolled-out the survey in three stages to iteratively analyse the  
177 collected data and in parallel improve the survey, until we deemed the data quality data  
178 to be satisfactory (Fig 1).



179

180 **Fig 1.** Conceptual flowchart of CALY-SWE valuation methods. □ symbolizes data  
181 collections and □ general processes and analysis steps. Arrows represent the workflow.

182 To improve usability, we focused on conveying information visually instead of relying on  
183 textual information. This included attribute levels of the states and the number of years  
184 for TTO questions. For both the DCE and TTO questions we used the same basic visual  
185 layout where attribute levels were visualized by adjacent vertical bars. The bars were  
186 filled according to the level of the 3 attribute levels, either empty, half, or full. Full CALY-  
187 SWE statement phrasing, and instructions were reachable over help buttons. For TTO  
188 questions we additionally displayed a horizontal bar that was filled according to the  
189 number of years in each state (See Supplementary, section on the web survey and  
190 Supplementary Fig S1 – S4).

### 191 Qualitative face validity study

192 To assess the face validity of the web survey, we conducted a qualitative study among  
193 16 participants recruited from 2 local non-profit associations in Umeå, with varied age  
194 and gender. From September to December 2020, participants filled in the survey on a  
195 tablet at the premises of Umeå university and were afterwards interviewed by two

196 qualitative researchers (Jan Hjelte and Kerstin Edin) based on a semi-structured  
197 questionnaire.

## 198 **Experimental design and sample size simulation**

199 Our main consideration was to produce an experimental design with a relatively low  
200 number of questions per participant to keep the survey length short. To still collect  
201 enough data, we intended to compensate by increasing the number of participants.  
202 Because we aimed for a hybrid model, we also wanted to find the best balance between  
203 the number of TTO and DCE questions administered, and an adequate sample size for  
204 generating the value set with a TTO-only model.

205 To that end, we generated different experimental designs for different configurations of  
206 number of TTO and DCE questions. For each configuration, we chose the number of  
207 TTO and DCE question so that the predicted time for survey completion would be  
208 approximately 20 minutes, based on timing data from the qualitative study. Twenty  
209 minutes was the maximum we considered workable for cost-effectively sampling from  
210 the online panel, as costs per participant increase with increasing drop-out rates for  
211 longer surveys. The configurations consisted of 11 and 2, 8 and 3, 6 and 4, 4 and 5, and  
212 3 and 6 DCE and TTO questions, respectively, with each 50, 100, 200, 300, 400, 500, 700,  
213 1000, 1400, and 2000 participants.

214 To develop a D-optimal design for DCE, we used the orthogonal design approach  
215 presented in Street et al. [15]. We used this same design for all configurations and  
216 randomly picked the DCE questions for a block. For TTO, we generated the different  
217 designs specific to the number of TTO questions per configuration by using the *skpr* [16]  
218 package in R with the D-optimality criteria.

219 Thereafter, we ran simulations with different number of participants for configuration,  
220 using a hybrid and a TTO-data-only model, both with a varying intercept for the TTO  
221 part. We generated the data with the same varying intercept specification and  
222 parametrized it based on results from an earlier pilot study (unpublished). We evaluated  
223 the performance of the simulation in terms of the mean credible interval (CRI) width the  
224 mean absolute error of the weight posteriors compared to the generative model.

## 225 **Comparison stage 1 and 2**

226 Directly comparing data quality between stages was not possible due to low and differing  
227 sample sizes in each stage. We therefore applied a bootstrap approach to examine the  
228 effect of changes between the stages. We randomly sampled 500 DCE answers and 300  
229 TTO answers 10000 times with repetition for each stage and performed a logit and an  
230 ordinary least squares linear regression for the bootstrapped DCE and TTO answers per  
231 stage with a main effects model. We then compared the distribution of the log likelihood,  
232 the mean standard error of the coefficients, and the number of inconsistent coefficients  
233 (level 2 greater than level 3) between stages.



## 234 TTO data quality and exclusion

235 The EQ-VT [9] focuses on interviewer training and poor quality data (defined by criteria  
236 such as if the interviewer explained the example long enough, or if the TTO answer for  
237 the pit state is the lowest TTO answer within a margin [17]) may be discarded. By  
238 design, self-administrated online surveys are unguided, making it impossible to rely on  
239 the interviewer performance for quality assessment.

240 To address this challenge, we explored alternative methods for assessing the quality of  
241 the TTO answers and for possibility excluding them. We developed a score based on  
242 inconsistent TTO answers that considers the severity of the inconsistencies, called the  
243 *combined inconsistency severity (CIS) score*.

244 To be logically consistent, participants' answers should value states with higher levels  
245 higher than those with lower levels for the same attribute. An *inconsistency* can occur for  
246 two TTO answers from a participant for the two states in a *dominated* choice. A  
247 dominated choice occurs if all attribute levels of one state are higher or equal while at  
248 least one level is higher than in the other state. If the TTO answer for the first state is  
249 lower or equal than for the second state, the two answers are not consistent with the  
250 dominated choice. Formally, for any two states  $S_j, S_k$  from a participant's TTO block with  
251 answers  $w_j, w_k$ , and  $i$  indexing the attributes:

$$252 \quad \forall i: S_{ji} \geq S_{ki} \wedge \exists i: S_{ji} > S_{ki} \wedge w_j \leq w_k, j \neq k$$

253

254 The *severity* of an inconsistency may be expressed as the absolute difference of the TTO  
255 answers ( $W$ ) and the absolute difference in TTO level attributes for the two involved  
256 states ( $L$ ), for participant  $p$ :

257

$$258 \quad L_{pjk} = \sum_i |S_{ji} - S_{ki}|$$

$$259 \quad W_{pjk} = |w_j - w_k|$$

260

261 For example, valuing 111111 (all six attributes on level 1, 'Do not agree') with 0.8 and  
262 333333 (All six attributes on level 3, 'Agree completely') with 0.2 could be considered a  
263 quite severe inconsistency, with the level differences being 12 ( $L$ ) and the weight  
264 difference ( $W$ ) equalling 0.6. On the other hand, an inconsistency involving 232323 and  
265 232333 with answers of 0.9 and 0.8 is less severe as there is only one level difference  
266 and a 0.1 weight difference. Such an inconsistency may occur due to the difficulty of the  
267 TTO question instead lack of engagement.

268

269 Across all values, we then normalized each score to [0,1] to make their scales  
270 comparable:

271 
$$Lnorm_{pjk} = \frac{L_{pjk} - \min(L_{pjk})}{\max(L_{pjk}) - \min(L_{pjk})}$$

272

273 
$$Wnorm_{pjk} = \frac{W_{pjk} - \min(W_{pjk})}{\max(W_{pjk}) - \min(W_{pjk})}$$

274

275 Afterwards, the score was summed for each participant  $p$  so that we could approximate  
276 data quality per participant and to not merely exclude single unfitting answers:

277

278 
$$CISscore_p = \left( \sum_{jk} Lnorm_{pjk} + \sum_{jk} Wnorm_{pjk} \right)$$

279

280

281 As the TTO experimental design did not contain the same number of support points for  
282 possible inconsistencies for each participant (ranging from 5 to 7) and different blocks  
283 may have been more challenging to answer consistently, we normalized the scores per  
284 block. At last, we calculated the score percentile for each participant so that specific  
285 proportions of data can be excluded:

286

287 
$$CISpercentilescore_p = \frac{rank_{block}(score_p)}{n_{block}}$$

288

289 We assessed the impact of including only participants with better quality TTO data by  
290 comparing the results of an ordinary linear regression for including 30% to 100% of the  
291 data according to the CIS score, in 1% increments. We specifically looked at the weights  
292 for 333333 and 111111 and their difference that defines the value set's range, the  
293 adjusted  $R^2$ , the number of inconsistent coefficients (where the coefficient for level 2 is  
294 larger than for level 3), the mean standard error and mean t-score of the coefficients  
295 (coefficients divided by the standard error) and the coefficient with their 95% confidence  
296 intervals.

297 We used R [18] for all data analyses with base R regression models except for the sample  
298 simulation with Bayesian hybrid models that were estimated with stan [19] and the  
299 cmdstan R interface [20].

## 300 Results

### 301 **Qualitative face validity study**

302 Participants' statements revealed that the DCE and TTO questions were challenging but  
303 meaningful, partially engaging, and thought-provoking. Some stated that answering the  
304 DCE block before the TTO block helped them to increase their familiarity with the  
305 statement phrasing before tackling the more complex TTO questions. Most participants  
306 were able to independently finish the survey by just relying on the on-screen  
307 instructions.

308 During the qualitative interviews, we continuously developed the survey based on  
309 participant's feedback. For example, we added an explanation video for the TTO  
310 question, revised wording, and overhauled the DCE question layout to be in the same  
311 style as the TTO question layout. See Supplementary Table S3 for a detailed list of  
312 changes and differences between survey versions.

### 313 **Experimental design sample size simulation**

314 We based the DCE design on an orthogonal array with 6 columns and 45 rows, resulting  
315 in 43 pairwise comparisons after the removal of two dominated comparisons. We used  
316 random blocking where the DCE questions for each participant were randomly picked  
317 from the 43 comparisons.

318 For TTO, we deemed eight blocks with three states per block to be adequate, and we  
319 augmented each block with the pit state 111111 (all six attributes on level 1) so that  
320 111111 would be evaluated by each participant, enabling to estimate 111111 with  
321 greater precision. The DCE and TTO designs are provided in Supplementary Tables S1  
322 and S2.

323 The sample size simulation resulted in a decrease of the mean 95% CRI and mean  
324 absolute errors with increasing sample sizes. The hybrid model generally performed  
325 better than the TTO only model, especially for large proportions of TTO questions, but  
326 differences were small for both mean CRI width and mean absolute errors for  
327 configurations with at least four TTO questions. We decided for the configuration with  
328 four TTO questions and six DCE questions as we judged this configuration to offer a  
329 good balance of TTO and DCE questions. With this configuration, we found that at least  
330 500 participants for using both TTO and DCE data and 1000 participants for only TTO  
331 data would be necessary. We decided to set the target sample size at 1500 participants,  
332 to have a safety margin and to leave the option to only use TTO data for generating the  
333 value set or to exclude poor quality data (Supplementary Fig S5).

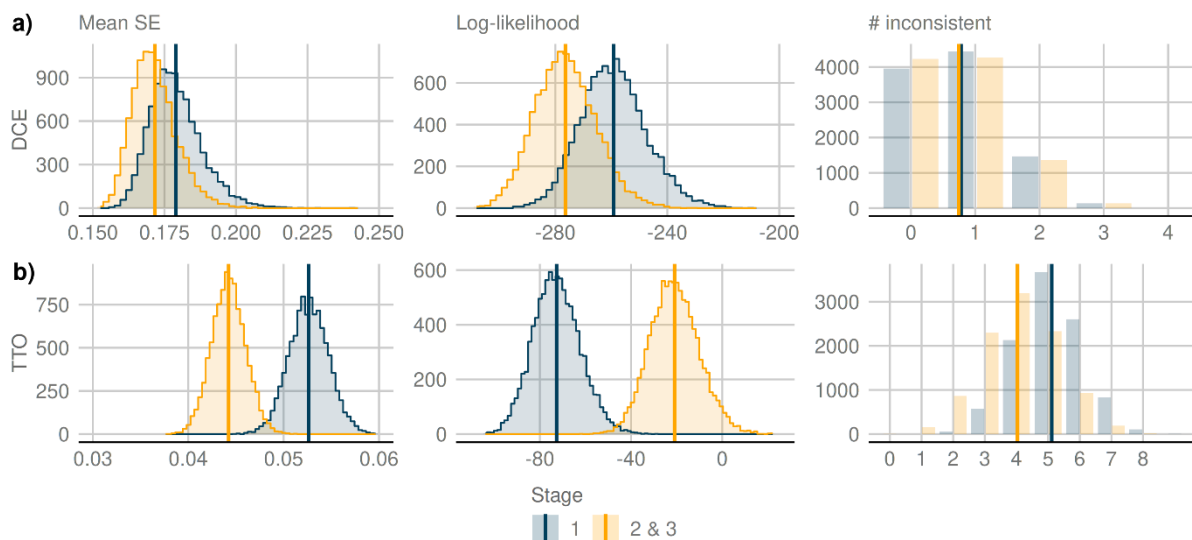
334

335 **Comparison stage 1 and 2**

336 We collected data in 3 stages (Fig 1): Stage 1 targeted 100 participants and took place  
337 from November 22 to December 2 2021. Based on preliminary analysis after stage 1,  
338 we aimed to further improve the TTO question format and piloted the changes in Stage  
339 2, targeting 200 participants (January 3 until January 12 2022). The main data collection  
340 in Stage 3 targeted 1500 participants (March 7 until April 18 2022).

341 In stage 2, we changed the iteration procedure so that participants choosing “equal” in  
342 the TTO question would not immediately proceed to the next question, but instead the  
343 interval of reachable values would shrink around the current bisection point, and  
344 renamed the option to “about equal”. The goal was to reduce the incentive to finish the  
345 question faster by choosing equal. We further decreased the TTO time frame from 20  
346 years to 10 years to reduce the expected number of needed iterations.

347 To facilitate understanding and increase engagement, we changed the graphical layout:  
348 Instead of selecting radio buttons, clicking on buttons would now directly submit the  
349 answer. We also introduced colour coding for the choices where the first and second  
350 choice were coloured in two distinct but neutral colours across DCE and TTO questions,  
351 randomized per participant (Supplementary Fig S2, S3). As a demonstrated example we  
352 also introduced a *learning* state as the first TTO state to showcase the trade-off  
353 mechanism (Details in Supplementary section on TTO learning state) and removed one  
354 DCE question per block to compensate the longer duration. Other changes included  
355 that participants could now navigate backwards. Supplementary Table S3 depicts a  
356 detailed list of changes.



357

358 **Fig 2.** Histograms of results of bootstrapped comparison between stage 1 and 2.

359 a) Logistic regression of DCE data with 500 draws with repetition.

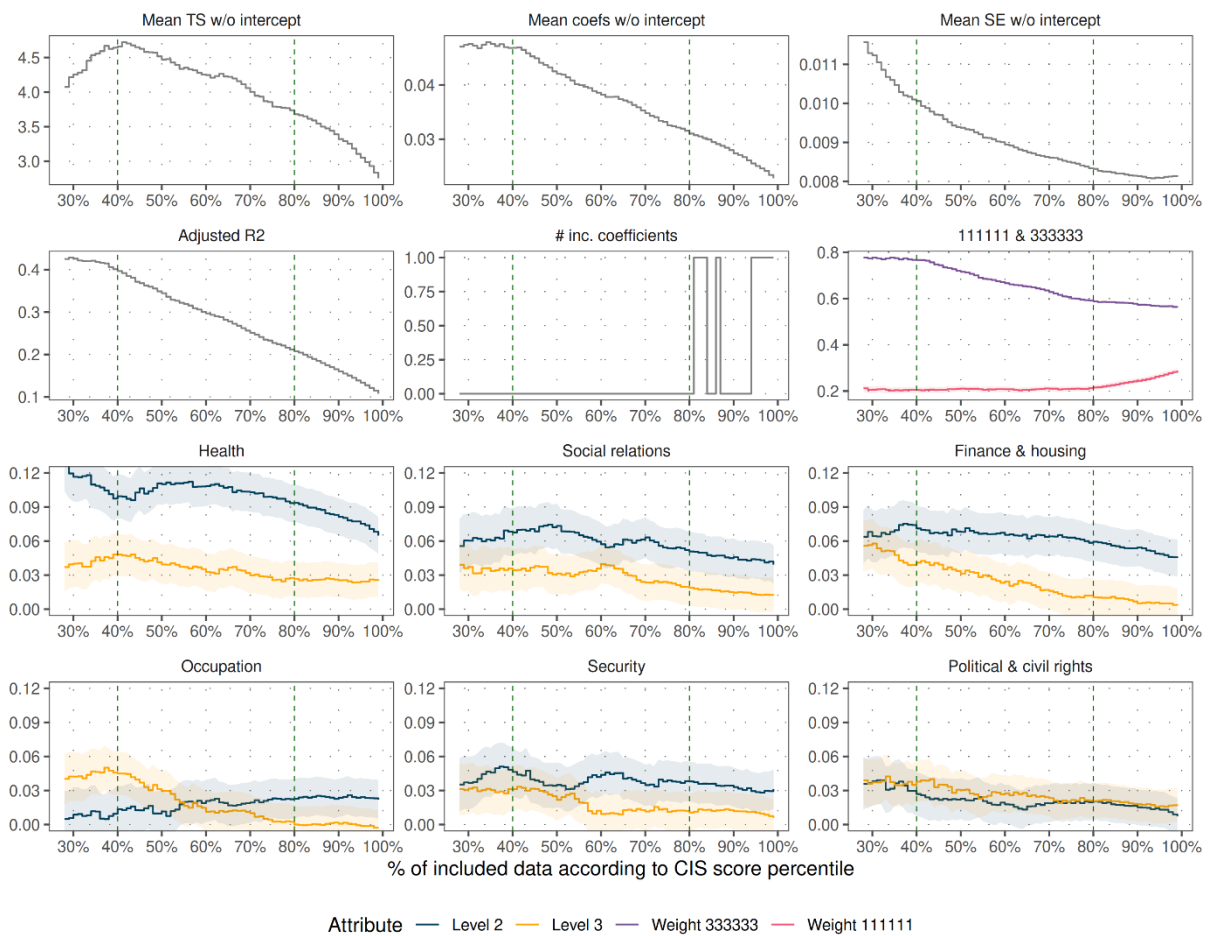
360 b) Linear regression of TTO data with 300 draws with repetition.

361 10000 bootstrap runs, main effects model. Discrete choice experiments (DCE). Time

362 trade-off (TTO). Standard error (SE). Time trade-off (TTO).

363 The bootstrapped regressions from stage 2 compared to stage 1 indicated  
 364 improvements in TTO data quality in terms of lower mean standard deviation and  
 365 higher log likelihood with distinct distributions. The number of inconsistent coefficients  
 366 also decreased on average but the distributions were overlapping. In comparison, the  
 367 DCE distributions did not differ clearly between stages for all the indicators and  
 368 indicated on average a decrease in data quality in stage 2 compared to stage 1. We  
 369 deemed the data quality satisfactory and launched stage 3 without any further  
 370 substantial changes. The final dataset included 199 stage 2 and 1498 stage 3 answers  
 371 but no stage 1 data because of the survey differences.

372 **TTO data quality and exclusion**



373

374 **Fig 3.** Line plots of effects of including participants according to the CIS score percentile,  
 375 in 1% increments, in terms of the results of a linear regression main effects model.  
 376 Shaded error bands represent 95% confidence intervals. T-score (TS). Standard error  
 377 (SE). Combined inconsistency severity (CIS).

378 For higher percentages of included participants (100% corresponding to 1694)  
 379 according to TTO CIS percentiles, the amount of variation explained by the model  
 380 decreased. This was indicated by a decrease in the adjusted  $R^2$ , but also by a decrease  
 381 of the mean t scores and by a general decrease of the coefficients and the range of

382 333333 to 111111. Put differently, the model was less able to systematically pick up  
383 inference patterns related to trade-offs between dimensions the more data was  
384 included. Instead, the fitted values converged towards the mean of the TTO answers:  
385 The intercept increased if more than 80% of data was included, together with a  
386 decrease of the level-attribute coefficients. Beyond including 80% the coefficient for  
387 occupation level 3 became increasingly inconsistent, likely related to the general  
388 decrease in coefficient magnitude. Interestingly, while in general the proportion of level  
389 2 to level 3 coefficients remained constant, for occupation level 2 became larger than  
390 level 3 and for finance & housing the level 3 decreased compared to level 2 the more  
391 TTO data was included.

392 Therefore, the results suggest using between 40% to 80% of TTO data for generating the  
393 value set. Beyond 80%, the increase of the intercept would affect the value of 111111  
394 and thus the lower anchoring of the value set on the [0,1] scale. In addition, the effect of  
395 the attribute level coefficients increasingly loses strength beyond 80%. Model  
396 performance of jointly estimated hybrid models with DCE data as well as  
397 representativity considerations may also play a role when deciding the exact amount of  
398 TTO data for generating the value set.

## 399 Discussion

400 We created an online-administered web survey for collecting TTO and DCE data for a  
401 CALY-SWE value set including an experimental design and we estimated the necessary  
402 sample size. We leveraged qualitative interviews and a two-stage survey roll-out to  
403 improve the web survey and developed a novel and sensitive way of excluding TTO data  
404 with poorer quality This work paved the way for eliciting a CALY-SWE value set which will  
405 be reported elsewhere.

406 Strengths include that we addressed the methodological challenges connected to online  
407 administration with an integrated approach that resulted in multiple benefits: 1) The  
408 qualitative face validity study allowed us to fine tune instructions and appearance of the  
409 web survey, and we were able to develop a web survey that focused on visually  
410 conveying the task instead of textual information or oral guidance by an interviewer. 2)  
411 Timing data from the qualitative study informed the sample size simulation which in  
412 turn facilitated an informed decision on the balance between the number of TTO and  
413 DCE questions and the sample size, enabling cost-efficient online sampling due to  
414 shorter survey length that is compensated by an increased sample size to collect  
415 enough data 4) The staged roll-out allowed to further improve the TTO data quality in  
416 stage 2 as indicated by the results of the bootstrapped regression analysis. 5) We  
417 developed a per-participant score that reflects the severity of TTO inconsistencies and  
418 enabled the exclusion of TTO data of an explicitly chosen proportion of participants to  
419 partially offset limited data quality connected to online self-administration. The CIS  
420 score also enable quality control for TTO data that is independent of the administration  
421 mode as it does not depend on interviewer performance.

422 Limitations include that we focused on a main effects only experimental design, thus  
423 not considering interactions. For this valuation study we focused on method  
424 development and producing a readily interpretable value set. For generating the DCE  
425 experimental designs we relied on orthogonal arrays as it was readily implementable  
426 without relying on additional software, but it may be less efficient compared to  
427 Bayesian experimental designs [21, 22].

428 Excluding TTO data based on inconsistencies is a contentious issue. In EQ-5D valuation  
429 studies sometimes evidently invalid answers such as always responding with the same  
430 value or the same pattern, or those selected in the feedback module, are excluded [23–  
431 26]. Excluding TTO data based on inconsistencies via the CIS score indeed constitutes a  
432 risk for data curation, were data is made to corresponds to prior norms because  
433 participants with larger uncertainties are discriminated against [27]. Similarly, concern  
434 over consistency of preferences over the range of excluded data, while generally stable  
435 in our study, still constitutes an important limitation as indicated by the changing order  
436 of coefficients for level 2 and 3 of occupation. On the other hand, under the assumption  
437 that a portion of the participants did not engage in the task sufficiently to correctly state  
438 their preferences, the score made it possible to transparently examine the impact of  
439 excluding different amounts of data, on a continuous range. Hence an informed,  
440 normative decision about how much data should be included could be made, and the  
441 data quality itself is defined in relation to the total sample. In contrast, excluding for  
442 example all participants that valued all states at 0.5 [9] is an absolute either-or decision  
443 without granular control over how much data is excluded.

444 Compromised TTO data quality resulting from online administration, compared to face-  
445 to-face data, or resulting from one interviewer per group compared one interviewer per  
446 person, has been previously quantified in the form of increased standard deviations  
447 [28], a lower number of trading iterations [10, 29], and a smaller range of the resulting  
448 value sets [10]. Those results align with our findings: Including poorer quality TTO data  
449 according to the CIS score resulted in a lower range and a larger standard error of the  
450 resulting weights. Similarly, changing the TTO iteration procedure in stage 2, so that  
451 choosing 'equal' still required additional iterations to arrive at an answer, improved the  
452 TTO data quality.

453 Other approaches to increase feasibility of valuation studies include videoconferencing  
454 which was found to be a viable alternative to face-to-face interviews for an Italian EQ-5D  
455 5L valuation study [30]. Lipman similarly found that tele-TTO interviews were feasible  
456 [31]. Compared to online self-administration, this still requires interviewers and training  
457 to conduct the interview, imposing significant costs.

458 Further research into refinements of self-administered online valuation surveys may  
459 additionally increase their feasibility. This particularly includes instructions and visual  
460 design properties of TTO questions, and ways to assess data quality, but also to general  
461 knowledge on the feasibility of online administration in combination with TTO data  
462 quality assessment.

463 **Conclusion**

464 TTO and DCE data collected in self-administered online surveys may be used to elicit  
465 value sets if the web survey development, the experimental design, and sample size  
466 considerations are optimized and well-coordinated. The severity of inconsistencies can  
467 be used on a per-participant basis to identify and exclude poor quality TTO data that  
468 does not contribute to the modelling of preferences.  
469



## 470 References

- 471 1. Månsdotter, A., Ekman, B., Feldman, I., Hagberg, L., Hurtig, A.-K., & Lindholm, L.  
472 (2017). We Propose a Novel Measure for Social Welfare and Public Health:  
473 Capability-Adjusted Life-Years, CALYs. *Applied Health Economics and Health Policy*,  
474 15(4), 437–440. <https://doi.org/10.1007/s40258-017-0323-0>
- 475 2. Meili, K. W., Månsdotter, A., Sundberg, L. R., Hjelte, J., & Lindholm, L. (2022). An  
476 initiative to develop capability-adjusted life years in Sweden (CALY-SWE): Selecting  
477 capabilities with a Delphi panel and developing the questionnaire. *PLOS ONE*, 17(2),  
478 e0263231. <https://doi.org/10.1371/journal.pone.0263231>
- 479 3. Brazier, J., Ratcliffe, J., Salomon, J., & Tsuchiya, A. (2016). *Measuring and Valuing*  
480 *Health Benefits for Economic Evaluation* (Second Edition.). Oxford, New York: Oxford  
481 University Press.
- 482 4. VanderWeele, T. J. (2017). On the promotion of human flourishing. *Proceedings of*  
483 *the National Academy of Sciences*, 114(31), 8148–8156.  
484 <https://doi.org/10.1073/pnas.1702996114>
- 485 5. Lugnér, A. K., & Krabbe, P. F. M. (2020). An overview of the time trade-off method:  
486 concept, foundation, and the evaluation of distorting factors in putting a value on  
487 health. *Expert Review of Pharmacoeconomics & Outcomes Research*, 20(4), 331–342.  
488 <https://doi.org/10.1080/14737167.2020.1779062>
- 489 6. Ternent, L., & Tsuchiya, A. (2013). A Note on the Expected Biases in Conventional  
490 Iterative Health State Valuation Protocols: *Medical Decision Making*.  
491 <https://doi.org/10.1177/0272989X12475093>

- 492 7. Rowen, D., Brazier, J., & Van Hout, B. (2014). A Comparison of Methods for  
493 Converting DCE Values onto the Full Health-Dead QALY Scale. *Medical Decision*  
494 *Making*, 35(3), 328–340. <https://doi.org/10.1177/0272989x14559542>
- 495 8. Ramos Goñi, J. M., Craig, B. M., Oppe, M., & Van Hout, B. (2016). *Combining*  
496 *continuous and dichotomous responses in a hybrid model*. EuroQol Working Paper  
497 Series.
- 498 9. Stolk, E., Ludwig, K., Rand, K., van Hout, B., & Ramos-Goñi, J. M. (2019). Overview,  
499 Update, and Lessons Learned From the International EQ-5D-5L Valuation Work:  
500 Version 2 of the EQ-5D-5L Valuation Protocol. *Value in Health*, 22(1), 23–30.  
501 <https://doi.org/10.1016/j.jval.2018.05.010>
- 502 10. Jiang, R., Shaw, J., Mühlbacher, A., Lee, T. A., Walton, S., Kohlmann, T., Norman, R., &  
503 Pickard, A. S. (2021). Comparison of online and face-to-face valuation of the EQ-5D-  
504 5L using composite time trade-off. *Quality of Life Research*, 30(5), 1433–1444.  
505 <https://doi.org/10.1007/s11136-020-02712-1>
- 506 11. The PHP Group. (2022, September 11). PHP: Hypertext Preprocessor. Retrieved  
507 September 11, 2022, from <http://php.net/>
- 508 12. Symphony Project. (n.d.). Home - Twig - The flexible, fast, and secure PHP template  
509 engine. Retrieved September 21, 2022, from <https://twig.symfony.com/>
- 510 13. The Maria DB foundation. (2022, November 11). Maria DB. Retrieved November 16,  
511 2022, from <https://mariadb.org/about/>
- 512 14. The SQLite project. (2022, November 11). SQLite. Retrieved November 16, 2022,  
513 from <https://www.sqlite.org>

- 514 15. Street, D. J., Burgess, L., & Louviere, J. J. (2005). Quick and easy choice sets:  
515 Constructing optimal and nearly optimal stated choice experiments. *International*  
516 *Journal of Research in Marketing*, 22(4), 459–470.  
517 <https://doi.org/10.1016/j.ijresmar.2005.09.003>
- 518 16. Morgan-Wall, T., & Khoury, G. (2021). Optimal Design Generation and Power  
519 Evaluation R: The skpr Package. *Journal of Statistical Software*, 99(1).  
520 <https://doi.org/10.18637/jss.v099.i01>
- 521 17. Ramos-Goñi, J. M., Oppe, M., Slaap, B., Busschbach, J. J. V., & Stolk, E. (2017). Quality  
522 Control Process for EQ-5D-5L Valuation Studies. *Value in Health*, 20(3), 466–473.  
523 <https://doi.org/10.1016/j.jval.2016.10.012>
- 524 18. R Core Team. (2022). R: A Language and Environment for Statistical Computing.  
525 Vienna, Austria: R Foundation for Statistical Computing. Retrieved from  
526 <https://www.R-project.org/>
- 527 19. Stan Development Team. (2022). Stan. Retrieved from <https://mc-stan.org>
- 528 20. Gabry, J., Češnovar, R., Bales, B., Morris, M., Popov, M., Lawrence, M., Landau, W.  
529 M., & Socolar, J. (2022, August 11). cmdstanr. Retrieved from [https://mc-](https://mc-stan.org/cmdstanr)  
530 [stan.org/cmdstanr](https://mc-stan.org/cmdstanr)
- 531 21. Kessels, R., Jones, B., Goos, P., & Vandebroek, M. (2011). The usefulness of Bayesian  
532 optimal designs for discrete choice experiments. *Applied Stochastic Models in*  
533 *Business and Industry*, 27(3), 173–188. <https://doi.org/10.1002/asmb.906>
- 534 22. Reed Johnson, F., Lancsar, E., Marshall, D., Kilambi, V., Mühlbacher, A., Regier, D. A.,  
535 Bresnahan, B. W., Kanninen, B., & Bridges, J. F. P. (2013). Constructing Experimental  
536 Designs for Discrete-Choice Experiments: Report of the ISPOR Conjoint Analysis

- 537 Experimental Design Good Research Practices Task Force. *Value in Health*, 16(1), 3–  
538 13. <https://doi.org/10.1016/j.jval.2012.08.2223>
- 539 23. Devlin, N. J., Shah, K. K., Feng, Y., Mulhern, B., & van Hout, B. (2018). Valuing health-  
540 related quality of life: An EQ-5D-5L value set for England. *Health Economics*, 27(1), 7–  
541 22. <https://doi.org/10.1002/hec.3564>
- 542 24. Ludwig, K., Graf von der Schulenburg, J.-M., & Greiner, W. (2018). German Value Set  
543 for the EQ-5D-5L. *Pharmacoeconomics*, 36(6), 663–674.  
544 <https://doi.org/10.1007/s40273-018-0615-8>
- 545 25. Omelyanovskiy, V. (2021). Valuation of the EQ-5D-3L in Russia. *Quality of Life*  
546 *Research*, 11.
- 547 26. Welie, A. G., Gebretekle, G. B., Stolk, E., Mukuria, C., Krahn, M. D., Enquoselassie, F.,  
548 & Fenta, T. G. (2020). Valuing Health State: An EQ-5D-5L Value Set for Ethiopians.  
549 *Value in Health Regional Issues*, 22, 7–14. <https://doi.org/10.1016/j.vhri.2019.08.475>
- 550 27. Viney, R., Mulhern, B., Norman, R., Shah, K., & Devlin, N. (n.d.). Quality control vs.  
551 'data curation': where should we draw the line in researcher judgements about.
- 552 28. Norman, R., King, M. T., Clarke, D., Viney, R., Cronin, P., & Street, D. (2010). Does  
553 mode of administration matter? Comparison of online and face-to-face  
554 administration of a time trade-off task. *Quality of Life Research*, 19(4), 499–508.  
555 <https://doi.org/10.1007/s11136-010-9609-5>
- 556 29. Shah, K. K., Lloyd, A., Oppe, M., & Devlin, N. J. (2013). One-to-one versus group  
557 setting for conducting computer-assisted TTO studies: findings from pilot studies in  
558 England and the Netherlands. *The European Journal of Health Economics*, 14(1), 65–  
559 73. <https://doi.org/10.1007/s10198-013-0509-9>

- 560 30. Finch, A. P., Meregaglia, M., Ciani, O., Roudijk, B., & Jommi, C. (2022). An EQ-5D-5L  
561 value set for Italy using videoconferencing interviews and feasibility of a new mode  
562 of administration. *Social Science & Medicine*, 292, 114519.  
563 <https://doi.org/10.1016/j.socscimed.2021.114519>
- 564 31. Lipman, S. A. (2021). Time for Tele-TTO? Lessons Learned From Digital Interviewer-  
565 Assisted Time Trade-Off Data Collection. *The Patient - Patient-Centered Outcomes*  
566 *Research*, 14(5), 459–469. <https://doi.org/10.1007/s40271-020-00490-z>
- 567
- 568

569 **Declarations and statements**

570

571 *Author contributions (CRediT):*

572 KWM: Conceptualization, data curation, formal analysis, investigation, methodology,  
573 project administration, software, visualization, data curation, writing - original draft,  
574 writing – review and editing.

575 BM: Conceptualization, methodology, supervision, writing – review and editing

576 RS: Conceptualization, methodology, formal analysis, writing – review and editing

577 JH: Conceptualization, formal analysis, investigation, methodology, supervision, writing –  
578 review and editing

579 KE: Conceptualization, formal analysis, investigation, methodology, writing- review and  
580 editing

581 FN: Supervision, writing – review and editing

582 IF: Funding acquisition, Supervision, writing – review and editing

583 AM: Funding acquisition, Supervision, writing – review and editing

584 LL: Funding acquisition, Project administration, Supervision, writing – review and editing

585

586 *Funding:* Forte – Swedish Research Council for Health, Working Life, and Welfare, Grant  
587 No 2018-00143, Principal investigator Lars Lindholm

588 *Competing interests:* The authors have no relevant financial or non-financial interests to  
589 disclose.

590 *Ethics approval:* This study was performed in line with the principles of the Declaration  
591 of Helsinki. Approval was granted by the Swedish Ethical Review Authority (Dnr 2021-  
592 04465).

593 *Consent to participate:* Informed consent was obtained from all individual participants  
594 included in the study, including the consent to publish results and anonymized data.

595