

Appendices for *Online Volunteer Laboratories for  
Human Subjects Research*

S4 Appendix: Supplemental Materials for Study 4

1 This appendix includes additional information on the design and results of individual  
 2 response quality tests for volunteer and paid subjects.

### 3 **A Sample Composition**

4 We provide an overview of the sample sizes for each of the tests of response quality  
 5 discussed in the manuscript. Due to a combination of missingness in demographic covariates,  
 6 plus some attrition throughout the survey, the sample size for each test ranges from around  
 7 1,400 to 1,700 when combined across survey topics and volunteer status (Table A). The table  
 8 suggests even balance across paid and volunteer subjects in missingness and attrition.

Table A: Sample Sizes by Response Quality Test, Survey Topic, and Volunteer Status.

	Sample Size	By Survey Topic		By Volunteer Status	
	Combined	FEP	Secular	MTurk	Volunteer
Time answering open-ended 2	1,477	719	758	755	722
Time answering open-ended 1	1,411	664	747	732	679
Time reading prompt		750		806	741
Time answering short items		697		740	698
Straightlining	1,685	822	863	871	814
Open-ended Response Effort	1,466	708	758	755	711
Open-ended question 2 length	1,475	718	757	753	722
Open-ended question 1 length	1,484	725	759	761	723
Open-ended Response Quality	1,466	708	758	755	711
Committal question 2	1,494	730	764	771	723
Committal question 1	1,493	730	763	771	722
Consistency	1,499	734	765	773	726
Skipping	1,509	741	768	774	735
Attention Check			773	384	389

9 We compare the demographic balance across the paid and unpaid samples that took the

10 response quality surveys in Table B. We show the samples are balanced on many variables,  
11 but that the DLABSS sample tends to be older, richer, whiter, more religious, and more  
12 politically conservative. As noted in the body of the paper, many of these differences actually  
13 make DLABSS more similar to the larger US population and are the result of intentional  
14 targeting of certain populations in DLABSS recruitment efforts. Notably however, DLABSS  
15 does contain more missingness in demographic variables.

Table B: Demographics in DLABSS and MTurk.

	<i>DLABSS</i>	<i>MTurk</i>
Female	49.0% (1.8) [5.1%]	50.3% (1.8) [0.1%]
Education (mean years)	16.0 (0.1) [3.7]	14.9 (0.1) [0.0]
Age (mean years)	56.0 (0.6) [2.4]	38.7 (0.5) [0.0]
Mean income	\$60,717 (\$1,594) [9.1]	\$44,727 (\$1,192) [0.7]
Median income	\$55,000	\$37,500
Race		
White	88.9 (1.1) [2.1]	77.0 (1.5) [0.0]
Black	2.4 (0.5)	6.8 (0.9)
Hispanic	4.5 (0.7)	7.4 (0.9)
Attend Religious Service Weekly	33.6 (2.2) [2.5]	15.0 (1.8) [0.0]
Religion		
Protestant	45.7 (2.3) [0.4]	34.4 (2.4) [0.0]
Other Christian	4.9 (1.0)	2.1 (0.7)
Other	14.8 (1.6)	19.9 (2.0)
None	34.6 (2.2)	43.7 (2.5)
Party Identification		
Democrat	37.4 (1.7) [7.9]	50.2 (1.8) [3.2]
Independent	30.5 (1.6)	23.8 (1.5)
Republican	32.1 (1.7)	26.0 (1.6)
Liberal	54.4 (2.3) [4.7]	60.5 (2.5) [0.0]
Speak English at Home	99.8 (0.2) [2.0]	100.0 (0.0) [0.0]
N	476-843	387-782

*Standard errors are in parentheses. Percent missing in brackets. N varies across questions due to missingness and the religious questions only appearing on one survey.*

## 16 B Summary of Response Quality Tests

17 Table C provides a summary of each measure of response quality employed in this study.

Table C: Tests of Response Quality

Quality Dimension	Test	Measurement
Time Investment	Time answering open-ended items	Seconds spent on response page for each of two open-ended survey items (presented as an average and separately)
	Time reading prompt	Seconds viewing an article of about 400 words
	Time answering short items	Seconds spent on response page with five short-answer/multiple choice items pertaining to article
Straightlining	Straightlining in matrix-style question grids	Share of three bidirectional question matrices with entirely uniform answers
Open-Ended Investment	Subjective effort	Effort score out of five (with five being most perceived effort) given by two human coders for one open-ended survey item
	Response length	Number of characters in response to each of two open-ended survey items (presented as an average and separately)
	Subjective response quality	Summary of three dimensions coded 0 or 1 by two human coders for one open-ended survey item: whether response is long, topical, and complete
Committal Answers	“Don’t Know” answers regarding commitment to action	Whether respondent answers “Don’t know” to each of two committal questions (presented as an average and separately)
Consistency	Contradicting previous responses	Whether subject direction of subjects’ response to a multiple-choice question contradicts previous response in grid-style question
Skipping	Skipping questions when given opportunity	Whether subject picks the “Skip” option in a multiple choice opinion question about policy opinions
Attention Check	Noticing embedded “attention check”	Whether subject answers a factual short-item question with “yes,” rather than the true answer, as directed in a reading prompt

## 18 C Response Quality by Survey Topic

19 In the body of the paper, we present plots including the standardized coefficients for  
20 volunteer response quality across all tests, including demographic covariates such as age, in-  
21 come, education, race, frequency of religious service attendance, religious tradition, political  
22 ideology, and party identification and dummy variable for the survey topic. We report re-  
23 sults using the combined survey data across studies. Here, we provide additional background  
24 information on test design and present coefficient plots separated by survey (Figure A). We  
25 also present the same coefficient plot for the bivariate regression of each quality test on volun-  
26 teer status, without controlling for demographic covariates (Figure B). Due to the embedded  
27 attention check in the secularism version of the survey, we can only present results for time  
28 investment into reading the article and answering subsequent questions for the foreign eco-  
29 nomic policy survey. Likewise, results for the attention check are limited to the secularism  
30 survey.

Figure A: Standardized “Volunteer” Coefficients for All Response Quality Tests, by survey with covariates

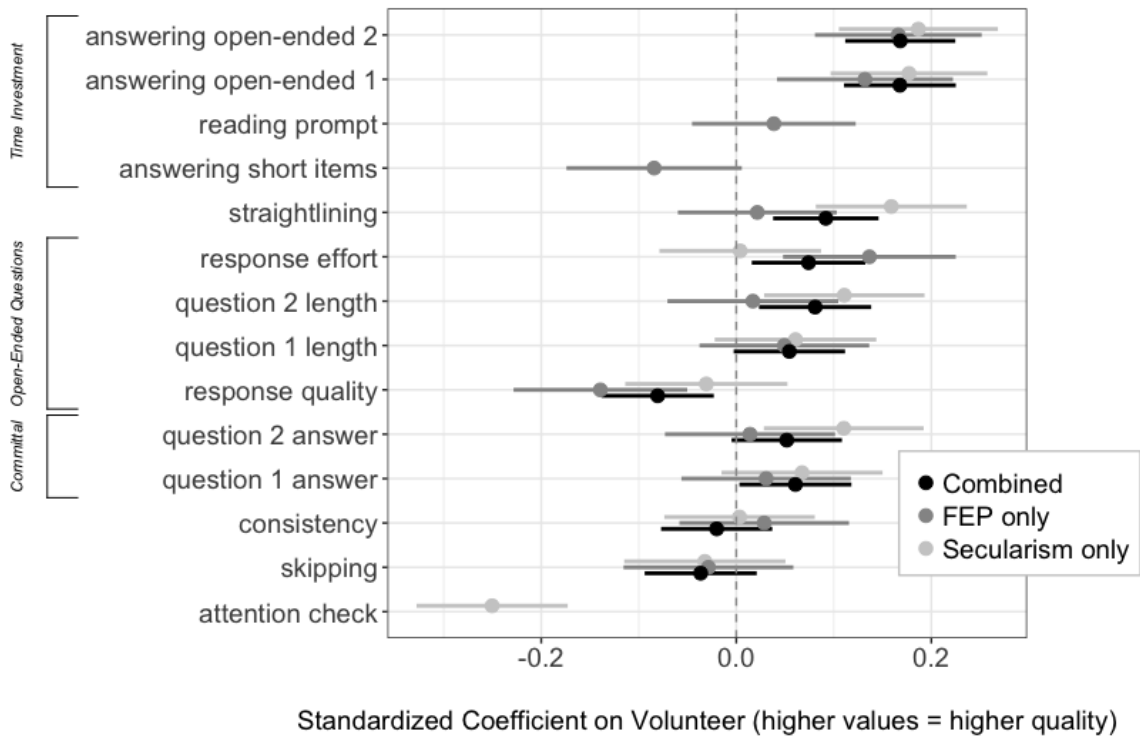
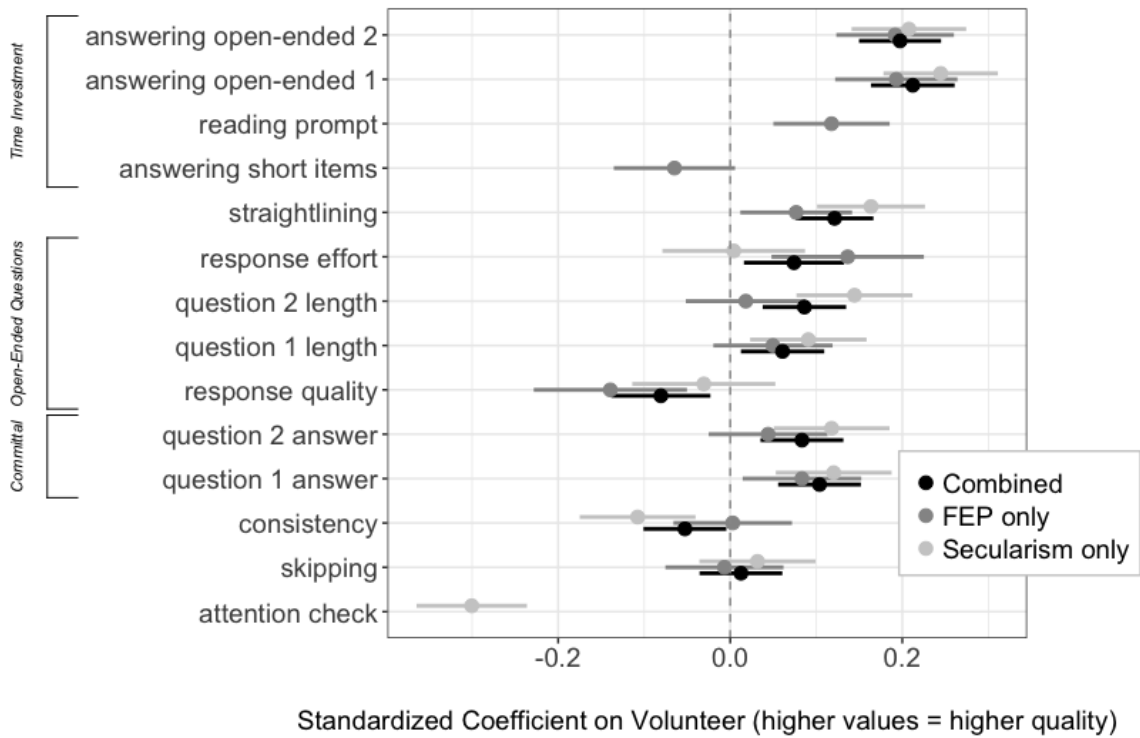


Figure B: Standardized “Volunteer” Coefficients for All Response Quality Tests, by survey without covariates





## 31 D Design and Results Details for All Tests

32 **Time Investment:** One measure of response quality is time investment in answering  
33 survey questions. In general, we consider longer time investment to be a reflection of more  
34 careful and higher quality responses. Our primary test involves a reading prompt followed  
35 by several short questions. We chose this format because reading prompts take more time  
36 than standard survey questions and subjects have no restrictions on how much time they  
37 need to spend reading the prompt. This test thus enables us to identify potential variation in  
38 time spent reading the prompt and answering questions, the latter of which are not available  
39 until a respondent confirms they have finished reading.

40 For the secularism questionnaire the prompt is an excerpt of a *New York Times* op-ed  
41 about the rise of secularism and atheism in American society. For the foreign economic  
42 policy study, respondents were invited to read an excerpt from *Foreign Policy Magazine*  
43 on Chinese global economic activities. The prompts are structurally similar in terms of  
44 their number and length of paragraphs, their emphasis on providing objective facts on the  
45 topic, and their overall length. Each prompt is about 400 words. While we include a time  
46 investment test in both surveys, we also embed an attention check in the secularism survey,  
47 invalidating comparisons across the two survey instruments for time investment. Thus, we  
48 present results only for the foreign economic policy survey.

49 For our primary test of subjects' investment of time into their survey responses, we mea-  
50 sure the number of seconds respondents spend reading the article prompt before clicking  
51 "next" to answer questions about it. We also measure the number of seconds spent respond-  
52 ing to five short-answer or multiple choice questions about the article. Two additional tests  
53 of time investment are the number of seconds respondents spend answering two open-ended  
54 questions later in the survey. For all of our time investment tests, the dependent variable  
55 is the number of seconds spent before clicking to the next page. Positive coefficients would

56 mean volunteers spend more time on the task, and we interpret more seconds spent as an  
57 indication of higher response quality. For this analysis, we trim the outer 5th percentile of  
58 time to eliminate extreme outliers.

59 Both with and without controls, volunteers spend more time than paid subjects on the  
60 reading prompts, though the difference becomes statistically insignificant with controls. For  
61 a random half of subjects, we also included a reminder to respondents to take their time  
62 reading, as they would not be permitted to click “back” to view the article. We do not analyze  
63 that effect here. Volunteers spent slightly less time responding to the block of questions about  
64 the article they had read. Volunteers spent more time responding to open-ended questions  
65 on both surveys than paid subjects, controlling for subject characteristics.

66 **Straightlining:** We examine whether paid and volunteer subjects have significantly dif-  
67 ferent propensities to engage in straightlining. Straightlining is a well-known phenomenon in  
68 survey research in which respondents rush through a survey and provide the same response  
69 for many questions without actually reading and considering the question content. More  
70 complex forms of straightlining include patterned or random responses, which are consid-  
71 erably more difficult for investigators to detect. With well-designed survey questions that  
72 naturally induce variation, less straightlining signals a higher quality response.

73 We design a test for straightlining that presents respondents with several matrix-style  
74 question blocks. This type of questioning is arguably especially vulnerable to straightlining.  
75 This is because several grid-style questions are presented on a single page, and answer choices  
76 to each question are located in close proximity to each other. Figure C offers an example of  
77 a question block from the foreign economic policy survey. We include seven of these blocks  
78 in each survey.

79 Question blocks within each survey vary in terms of their directionality, that is, the  
80 extent to which choosing the same answer for each question would reflect consistent, logical  
81 attitudes. Including bidirectionality in some of these question blocks allows us to detect

Figure C: Sample Question Matrix used in Straightlining Tests

To what extent do you agree or disagree with each of the following statements about America?

	Strongly disagree	Disagree	Neutral	Agree	Strongly agree
In the U.S., our people are not perfect, but our culture is superior to others.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I would rather be a citizen of America than of any other country in the world.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
The world would be a better place if people from other countries were more like Americans.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

82 straightlining in instances where a respondent with consistent preferences should not choose  
 83 the same answer category. For example, in one of the question blocks in the foreign economic  
 84 policy survey, we ask respondents whether they thought both “border walls” and “more  
 85 open borders” are beneficial for American interests. In this case, vertical straightlining  
 86 behavior, in which a respondent chooses the same preference for every question, would be  
 87 strong evidence that a subject is rushing through a survey. While respondents see multiple  
 88 blocks, some with bidirectional response items and others with unidirectional response items,  
 89 our measure of straightlining is constructed including only those blocks with bidirectional  
 90 responses, i.e., those on which a respondent paying attention would *not* be expected to  
 91 straightline naturally.

92 Arguably the most straight-forward test of straightlining is to simply create a variable  
 93 that measures whether a subject engaged in vertical straightlining—meaning he or she an-  
 94 swered the same response category for every question—for each question block. This type of  
 95 answer behavior, if detected, is perhaps the most egregious form of straightlining and thus

106 represents a useful first step. Thus, we use a simple binary operationalization of whether  
107 a respondent engaged in vertical straightlining by offering a uniform response category for  
108 every item within a straightlining block. Since each of our surveys had three total bidirec-  
109 tional grid-style question matrices, our dependent variable in this test is the proportion of  
110 these three questions in which the subject did offer a single uniform response category.

111 As depicted in Figure B, the volunteer coefficient for both the economic policy and  
112 secularism surveys is statistically distinguishable from the null in the bivariate regressions,  
113 in the direction of a higher quality response. In a multivariate setup (Figure A), volunteers  
114 were less likely to engage in straightlining in the secularism survey than paid respondents,  
115 but there was no significant difference between volunteer and paid respondents in the survey  
116 related to economic policy.

117 **Open-Ended Investment:** We next test the possibility that subjects motivated by  
118 different incentives vary in the quality of their open-ended survey answers. We include  
119 multiple open-ended response questions in each survey, which provide subjects with the  
120 opportunity to expand on their other answers, leave feedback for the research team, or  
121 otherwise write additional content. Researchers relying on open-ended response data may  
perceive higher quality responses as those which are longer and include more interesting  
content. Researchers might also value open-ended responses that provide suggestions for  
improving the study.

122 In the body of the manuscript and Figure A, we present results on open-ended response  
123 quality based on a human-coded, subjective, composite quality score with several dimensions.  
124 Two undergraduate students coded each open-ended response. They provided a subjective  
125 1-5 ranking of the amount of effort they perceived the subjects invested in the item. They  
126 also provide binary 0-1 score for whether each response was long, topical, and complete.  
127 These last three dimensions of quality are aggregated into a subjective quality score. One of  
128 the authors also coded all responses on which the two coders had a difference of more than

122 2 points in the 5-point effort scale or a disagreement on the binary score for any of the three  
123 quality dimensions. For both effort and the composite quality score, higher values represent  
124 higher quality responses.

125 We also test open-ended response quality by checking whether respondents significantly  
126 differ in the number of characters written in response to open-answer prompts. We argue  
127 that increases in this measure on average indicate a higher quality response or, at least,  
128 greater participant investment in the survey.

129 For all open-ended tests, we differentiate between open-ended questions left blank due to  
130 skipping (in which a subject viewed an open-ended question but did not write anything) and  
131 those left blank because the subject had attrited prior to viewing the open-ended question.  
132 We include skipped questions as zeros in our analysis, but exclude attrited respondents from  
133 the analysis.

134 Figure A provides mixed quantitative support for the idea that respondent motivation  
135 impacts responses to open-ended questions. On both open-ended questions, volunteers wrote  
136 more characters than paid respondents for the secularism survey, though not the foreign  
137 policy survey. For the human-coded scores, on average across the surveys, volunteers scored  
138 higher on effort, but lower on response quality overall.

139 **Non-committal responses:** Our surveys also include tests of commitment. In this  
140 context, commitment refers to subjects' willingness to signal intensity of attitudes in their  
141 reported responses by stating their intent to support (or oppose) a cause with behavior  
142 beyond simply reported survey responses. We employ multiple measures of commitment  
143 in each survey. For example, in the secularization survey we ask whether respondents are  
144 willing to 1) sign a petition and 2) confront an individual about inappropriate conduct.  
145 Subjects can report their intent to partake or abstain from either behavior, or can choose a  
146 less committal answer such as "It depends" or "Unsure."

147 For our tests of commitment, we code a noncommittal answer as one in which a respon-

148 dent does not express intent to participate or abstain from participation. Choosing “unsure”  
149 or “it depends” was taken as a lower-quality response, though we address ambiguity about  
150 interpreting these results in our discussion of findings below.

151 Results on these items were mixed. In the absence of demographic covariates, volun-  
152 teers were somewhat less likely to choose the noncommittal response for all but the second  
153 committal item on the economic policy survey. However, after the inclusion of demographic  
154 covariates, these distinctions are no longer significant in many cases.

155 We note that a response of “unsure” can be interpreted in various ways with respect  
156 to response quality. Although one interpretation of quality is to expect higher-quality re-  
157 sponses to include fewer non-committal responses, alternative interpretations are possible.  
158 For example, more unsure responses might be an expression that respondents are less likely  
159 to engage in cheap talk. Future research could work to examine this distinction in greater  
160 detail.

161 **Inconsistency:** We examine consistency across a subject’s answers within a survey. In  
162 general we perceive higher consistency as a measure of higher response quality: if respondents  
163 report unstable or illogical opinions, preferences, or other responses in a short study, there  
164 is certainly reason to doubt whether such responses reflect actual attitudes (Achen 1975).  
165 Inconsistency alternatively may simply be a proxy for lower levels of attention paid by  
166 subjects to the study content, an equally worrisome possibility.

167 To explore the possibility that paid and unpaid subjects differ in their cognitive invested  
168 in the content of a study, we design a test that asks subjects the same question twice.  
169 The first version of the question is embedded in one of the straightlining blocks discussed  
170 above. The second version is a the same question phrased differently and presented in a  
171 different format, in this case as a standalone multiple choice question later in the survey.  
172 We randomize whether or not a respondent is first reminded that he or she has already been  
173 asked about this issue. We do not analyze that effect here.

174 We score this metric by creating a variable that measures whether there is directional  
175 consistency across the two questions, that is, if a respondent’s reported preferences are in  
176 the same direction, if not to the same degree. To require exact equality in both direction and  
177 degree seemed to be a test so stringent that it was inconsistent with what the literature would  
178 anticipate as reasonably high-quality and consistent (Ansolabehere, Rodden, and Snyder  
179 2008).

180 Figure A depicts the likelihood of volunteer versus paid subjects responding in a consis-  
181 tent direction on the two items. More consistency is seen as an indicator of higher-quality  
182 responses. For the most part, paid and unpaid subjects did not statistically differ on these  
183 items. As one exception, volunteer survey responses were arguably of lower quality in relation  
184 to response consistency on the secularism version of the survey in the bivariate framework.  
185 However, this distinction did not hold once control variables were introduced.

186 **Skipping:** We design a simple test to detect a subject’s propensity to skip questions.  
187 Because question skipping creates missing data, which can create bias if not corrected (King  
188 et al. 2001), subjects who skip fewer questions are typically more desirable than those who  
189 skip more. We design this test by creating a question about individual’s preferences for  
190 a certain policy. For example, in the secularization survey we ask a question related to  
191 churches’ rights to engage in political activities, and for the foreign economic policy survey  
192 we ask about whether the government should encourage more free trade. For these questions  
193 respondents could choose, “I’m not sure. Skip.” as an answer choice beyond the standard  
194 support-oppose scale. To measure skipping, we simply create a variable that receives the  
195 value of “1” if a respondent chose to skip the question.

196 It is worth noting that the choice to skip in survey questions can be interpreted in  
197 multiple directions in relation to response quality. On the one hand, skipping may represent  
198 taking an “easy way out,” wherein subjects avoid engaging with a cognitively challenging  
199 question or are simply rushing to finish and, thus, may represent low quality responses. On

200 the other hand, skipping may genuinely represent uncertainty among subjects and may be  
201 more desirable than other response strategies, such as choosing an answer at random. For  
202 subjects who are aware that they are uninformed on a particular policy issue, skipping may  
203 represent a reasonable, high-quality choice. There are no significant differences between paid  
204 and unpaid respondents in this test of skipping.

205 **Attention Check:** We also embed an attention check, or “screener” in our secularism  
206 survey instrument. As noted above, we task subjects with reading an article and responding  
207 to questions about that article. In order to ensure respondents were actually reading the  
208 article, not merely opening the prompt and then doing something else before proceeding, we  
209 included a sentence in the middle of the article that asked respondents to reply “yes” to an  
210 open-ended question in the question block following the article, rather than answering the  
211 question. We include this test only in the text of the article in the secularism survey.

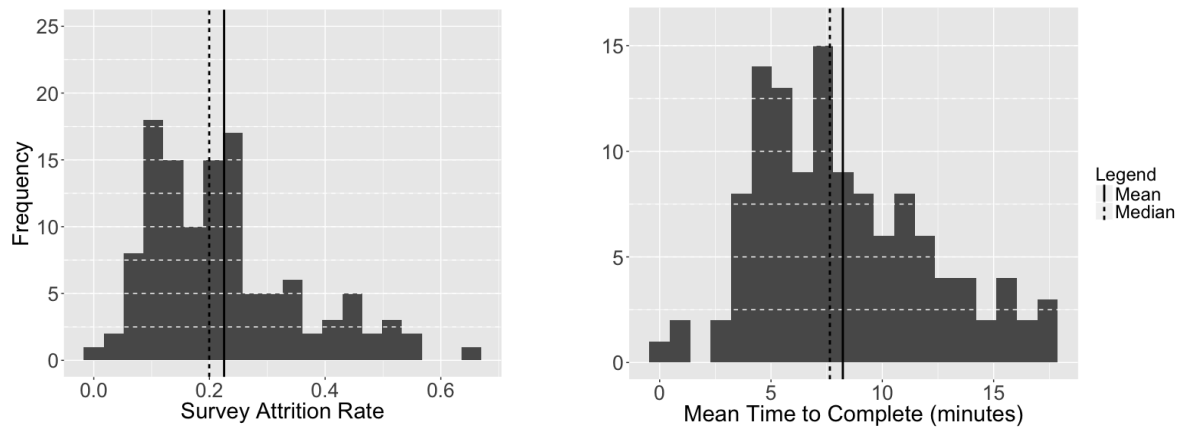
212 In the manuscript, we discuss ways in which volunteer subjects may potentially exhibit  
213 less propensity for attrition than paid subjects. Figure D displays the mean and median  
214 attrition rates and completion times for DLABSS studies.

## 215 E Attrition and Response Time

216 Figure D shows the frequency of attrition rates by survey and the average time to complete  
217 by survey across 120 studies hosted on DLABSS.



Figure D: DLABSS Study Attrition Rates and Completion Times



*Left figure is total attrition by study and right figure is mean time to complete by study for 120 studies hosted by DLABSS.*

## References

- Achen, Christopher H. 1975. “Mass political attitudes and the survey response.” *American Political Science Review* 69(04): 1218–1231.
- Ansolabehere, Stephen, Jonathan Rodden, and James M Snyder. 2008. “The strength of issues: Using multiple measures to gauge preference stability, ideological constraint, and issue voting.” *American Political Science Review* 102(2): 215–232.
- King, Gary, James Honaker, Anne Joseph, and Kenneth Scheve. 2001. “Analyzing incomplete political science data: An alternative algorithm for multiple imputation.” *American Political Science Review* 95(1): 49–69.