

Unrepresentative Big Surveys Significantly Overestimate US Vaccine Uptake

Valerie C. Bradley^{1*}, Shiro Kuriwaki^{2*}, Michael Isakov³,
Dino Sejdinovic¹, Xiao-Li Meng⁴, Seth Flaxman^{5,†}

¹Department of Statistics, University of Oxford, Oxford, UK

²Department of Political Science, Stanford University, Stanford, CA, USA

³Harvard College, Harvard University, Cambridge, MA, USA

⁴Department of Statistics, Harvard University, Cambridge, MA, USA

⁵Department of Computer Science, University of Oxford, Oxford, UK

*These authors contributed equally to this work.

†Correspondence to: seth.flaxman@cs.ox.ac.uk.

October 2021

Surveys are a crucial tool for understanding public opinion and behavior, and their accuracy depends on maintaining statistical representativeness of their target populations by minimizing biases from all sources. Increasing data size shrinks confidence intervals but magnifies the impact of survey bias -- an instance of the Big Data Paradox¹. Here we demonstrate this paradox in estimates of first-dose COVID-19 vaccine uptake in US adults: Delphi-Facebook^{2,3} (about 250,000 responses per week) and Census Household Pulse⁴ (about 75,000 per week). By May 2021, Delphi-Facebook overestimated uptake by 17 percentage points and Census Household Pulse by 14, compared to a benchmark from the Centers for Disease Control and Prevention (CDC). Moreover, their large data sizes led to negligible statistical uncertainty on the incorrect estimates. In contrast, an Axios-Ipsos online panel⁵ with about 1,000 responses following survey research best practices⁶ provided reliable estimates and uncertainty. We decompose observed error using a recent analytic framework¹ to explain the inaccuracy in the three surveys. We then analyze the implications for vaccine hesitancy and willingness. We show how a survey of 250,000 respondents can produce an estimate of the population mean that is no more accurate than an estimate from a simple random sample of size 10. Our central message is that data quality matters far more than data quantity, and compensating the former with the latter is a mathematically provable losing proposition.

Governments, businesses, and researchers rely on survey data to inform the provision of government services⁷, steer business strategy, and guide response to the COVID-19 pandemic^{8,9}.

With the ever-increasing volume and accessibility of online surveys and organically-collected data, the line between traditional survey research and Big Data is becoming increasingly blurred¹⁰. Large datasets enable analysis of fine-grained subgroups, which are in high-demand

36 for designing targeted policy interventions¹¹. However, counter to common intuition¹², larger
37 sample sizes alone do not ensure lower error. Instead, small biases are *compounded* as sample
38 size increases¹.

39 We see initial evidence of this in the discrepancies in estimates of first-dose COVID-19
40 vaccine uptake, willingness, and hesitancy from three online surveys in the US. Two of them —
41 Delphi-Facebook’s COVID-19 symptom tracker^{2,3} ($n \approx 250,000$ per week and with over 4.5
42 million responses from January to May 2021) and the Census Bureau’s Household Pulse survey⁴
43 ($n \approx 75,000$ per survey wave and with over 600,000 responses from January to May 2021) —
44 have large enough sample sizes to render standard uncertainty intervals negligible, yet report
45 significantly different estimates of vaccination behavior with nearly identically-worded questions
46 (Table 1). For example, Delphi-Facebook’s state-level estimates for willingness to receive a
47 vaccine from the end of March 2021 are 8.5 percentage points lower on average than those from
48 the Census Household Pulse (Extended Data Fig. 1A), with differences as large as 16 percentage
49 points.

50 The US Centers for Disease Control and Prevention (CDC) compiles and reports vaccine
51 uptake statistics from state and local offices¹³. These figures serve as a rare external benchmark,
52 permitting us to compare survey estimates of vaccine uptake to those from the CDC. The CDC
53 has noted the discrepancies between their own reported vaccine uptake and that of the Census
54 Household Pulse^{14,15}, and we find even larger discrepancies between the CDC and Delphi-
55 Facebook data (Fig 1a). In contrast, the Axios-Ipsos’ Coronavirus Tracker⁵ ($n \approx 1,000$ responses
56 per wave, and over 10,000 responses from January to May 2021) tracks the CDC benchmark
57 well. None of these surveys use the CDC benchmark to adjust or assess their estimates of
58 vaccine uptake, thus by examining patterns in these discrepancies, we can infer each survey’s

59 accuracy and *statistical representativeness*, a nuanced concept that is critical for the reliability of
60 survey findings¹⁶⁻¹⁹.

61 **The Big Data Paradox in vaccine uptake**

62 We focus on the Delphi-Facebook and Census Household Pulse surveys because their
63 large sample sizes (each greater than 10,000 respondents²⁰) present the opportunity to examine
64 the Big Data Paradox¹ in surveys. The Census Household Pulse is an experimental product
65 designed to rapidly measure pandemic-related behavior. Delphi-Facebook has stated that the
66 intent of their survey is to make comparisons over space, time, and subgroups and that point
67 estimates should be interpreted with caution³. However, despite these intentions, Delphi-
68 Facebook has reported point estimates of vaccine uptake in its own publications^{11,21}.

69 Delphi-Facebook and Census Household Pulse surveys persistently overestimate vaccine
70 uptake relative to the CDC's benchmark (Fig. 1a). Despite being the smallest survey by an order
71 of magnitude, Axios-Ipsos' estimates track well with the CDC rates (see Fig 1a), and their 95%
72 confidence intervals contain the benchmark estimate from the CDC in 10 out of 11 surveys (an
73 empirical coverage probability of 91%).

74 One might hope that estimates of changes in first-dose vaccine uptake are correct, even if
75 each snapshot is biased. Unfortunately, errors have increased over time, from just a few
76 percentage points in January 2021 to 4.2 percentage points (Axios-Ipsos), 14 percentage points
77 (Census Household Pulse), and 17 percentage points (Delphi-Facebook) by mid-May 2021
78 (Fig. 1b). For context, for a state near the herd immunity threshold (70-80% based on recent

79 estimates²²), a discrepancy of 10 percentage points in vaccination rates could be the difference
80 between containment and uncontrolled exponential growth in new SARS-CoV-2 infections.

81 Conventional statistical formulas for uncertainty further mislead when applied to biased
82 big surveys because as sample size increases, bias (rather than variance) dominates estimator
83 error. Fig. 1a shows 95% confidence intervals for vaccine uptake based on each survey's reported
84 sampling standard errors and weighting design effects²³. Axios-Ipsos has the widest confidence
85 intervals, but also the smallest design effects (1.1-1.2), suggesting that its accuracy is driven
86 more by minimizing bias in data collection rather than post-survey adjustment. Census
87 Household Pulse's 95% confidence intervals are widened by large design effects (4.4-4.8) but
88 they are still too narrow to include the true rate of vaccine uptake in almost all survey waves.
89 The confidence intervals for Delphi-Facebook are vanishingly small, driven by large sample size
90 and moderate design effects (1.4-1.5), and give us essentially zero chance of being even close to
91 the truth.

92 One benefit of such large surveys might be to compare estimates of spatial and
93 demographic subgroups²⁴⁻²⁶. However, in March of 2021, Delphi-Facebook and Census
94 Household Pulse over-estimated CDC state-level vaccine uptake by 16 and 9 percentage points,
95 respectively (Extended Data Fig. 1G-H), and by equal or larger amounts by May 2021 (Extended
96 Data Fig. 2G-H). Relative estimates were no better than absolute estimates in March of 2021:
97 there is barely any agreement in a survey's estimated state-level rankings with the CDC (a
98 Kendall rank correlation of 0.31 for Delphi-Facebook in Extended Data Fig. 1I, 0.26 for Census
99 Household Pulse in Extended Data Fig. 1J) but they improved in May of 2021 (correlations 0.78
100 and 0.74 in Extended Data Fig. 2I-J). Among 18-64 year-olds, both Delphi-Facebook and Census
101 Household Pulse overestimate uptake, with errors increasing over time (Extended Data Fig. 6).

102 These examples illustrate a mathematical fact. That is, when biased samples are large,
103 they are doubly misleading: they produce confidence intervals with incorrect centers and
104 substantially underestimated widths. This is the Big Data Paradox¹: *the larger the data size, the*
105 *surer we fool ourselves* when we fail to account for bias in data collection.

106 **A framework for quantifying data quality**

107 While it is well-understood that traditional confidence intervals capture only survey
108 sampling errors²⁷ (and not total error), the traditional survey framework lacks analytic tools for
109 quantifying nonsampling errors separately from sampling errors. A recently formulated statistical
110 framework¹ permits us to exactly decompose total error of a survey estimate into three
111 components:

$$112 \quad \text{Total Error} = \mathbf{Data\ Quality\ Defect} \times \mathbf{Data\ Scarcity} \times \mathbf{Inherent\ Problem\ Difficulty}. \quad (1)$$

113 This framework has been applied to COVID-19 case counts²⁸ and election forecasting²⁹.
114 Its full application requires ground-truth benchmarks or their estimates from independent
115 sources¹.

116 Specifically, **Total Error** is the difference between the observed sample mean \bar{Y}_n as an
117 estimator of the ground truth, the population mean \bar{Y}_N . The **Data Quality Defect** is measured
118 using $\hat{\rho}_{Y,R}$, called *data defect correlation (ddc)*¹, which quantifies total bias (from any source),
119 measured by the correlation between the event that an individual's response is recorded and its
120 value, Y . The impact of data quantity is captured by **Data Scarcity**, which is a function of the
121 sample size n and the population size N , measured as $\sqrt{\frac{N-n}{n}}$, and hence what matters for error is

122 the relative sample size, i.e., how close n is to N , rather the absolute sample size n . The third
123 factor is **Inherent Problem Difficulty**, which measures the population heterogeneity (via
124 standard deviation σ_Y of Y), because the more heterogeneous a population is, the harder it is to
125 estimate its average well. Mathematically, Equation 1 is given by $\bar{Y}_n - \bar{Y}_N = \hat{\rho}_{Y,R} \times \sqrt{\frac{N-n}{n}} \times \sigma_Y$.
126 Incidentally, this expression was inspired by the Hartley-Ross inequality for biases in ratio
127 estimators published in *Nature* in 1954.³⁰ More details on the decomposition are provided in the
128 Methods (“Calculation and interpretation of ddc ”), where we also present a generalization for
129 weighted estimators.

130 **Decomposing error in COVID surveys**

131 While the data defect correlation ddc is not directly observed, COVID-19 surveys present
132 a rare case in which it can be deduced because all other terms in Equation 1 are known (see
133 “Calculation and interpretation of ddc ” in the Methods for an in-depth explanation). We apply
134 this framework to the aggregate error shown in Fig. 1b, and the resulting components of error
135 from the right-hand side of Equation 1 are shown in Fig. 1c-e.

136 We use the CDC’s report of the cumulative count of first doses administered to US adults
137 as the benchmark^{8,13}, \bar{Y}_N . This benchmark may suffer from administrative delays and slippage in
138 how the CDC centralizes information from states^{31–34}. As a sensitivity analysis to check the
139 robustness of our findings to further misreporting, we present our results with sensitivity
140 intervals under the assumption that CDC’s reported numbers suffer from $\pm 5\%$ and $\pm 10\%$ error.
141 These scenarios were chosen based on analysis of the magnitude by which the CDC’s initial
142 estimate for vaccine uptake by a particular day increases as the CDC receives delayed reports of

143 vaccinations that occurred on that day (Extended Data Fig. 3 and Supplementary Information
144 A.2). That said, these scenarios may not capture latent systemic issues affecting CDC vaccination
145 reporting.

146 The **Total Error** of each survey’s estimate of vaccine uptake (Fig. 1b) increases over
147 time for all studies, most markedly for Delphi-Facebook. The **Data Quality Defect**, measured by
148 the *ddc*, also increases over time for Census Household Pulse and for Delphi-Facebook (Fig. 1c).
149 The *ddc* for Axios-Ipsos is much smaller and steady over time, consistent with what one would
150 expect from a representative sample. The **Data Scarcity** ($\sqrt{\frac{N-n}{n}}$) for each survey is roughly
151 constant across time (Fig. 1d). **Inherent Problem Difficulty** is a population quantity common to
152 all three surveys which peaks when the benchmark vaccination rate approaches 50% in April
153 2021 (Fig. 1e). Therefore, the decomposition suggests that the increasing error in estimates of
154 vaccine uptake in Delphi-Facebook and Census Household Pulse is primarily driven by
155 increasing *ddc*, which captures the overall impact of the bias in coverage, selection, and
156 response.

157 Equation 1 also yields a formula for the bias-adjusted effective sample size n_{eff} , which is
158 the size of a simple random sample that we would expect to exhibit the same level of Mean
159 Square Error (MSE) as what was actually observed in a given study with a given *ddc*. Unlike the
160 classical effective sample size²³, this quantity captures the impact of bias as well as that of an
161 increase in variance from weighting and sampling. Details for this calculation are in Methods
162 (“Error decomposition with survey weights”).

163 For estimating the US vaccination rate, Delphi-Facebook has a bias-adjusted effective
164 sample size of less than 10 in April 2021, a 99.99% reduction from the raw average weekly

165 sample size of 250,000 (Fig. 2). The Census Household Pulse also suffers from over 99%
166 reductions in effective sample size by May 2021. A simple random sample would have
167 controlled estimation errors by controlling *ddc*. However, once this control is lost, small
168 increases in *ddc* beyond what is expected in simple random samples can result in drastic
169 reductions of effective sample sizes for large populations¹.

170 Comparing study designs

171 Understanding *why* bias occurs in some surveys but not others requires an understanding
172 of the sampling strategy, modes, questionnaire, and weighting scheme of each survey. Table 1
173 compares the design of each survey (more details in the “Additional survey methodology” in
174 Methods and Extended Data Table 1).

175 All three surveys are conducted online and target the US adult population, but vary in
176 respondent recruitment methods³⁵. The Delphi-Facebook survey recruits respondents from active
177 Facebook users (the Facebook Active User Base, or FAUB) using daily unequal-probability
178 stratified random sampling². The Census Bureau uses a systematic random sample to select
179 households from the subset of the Census’ Master Address File (MAF) for which they have
180 obtained either cell phone or email contact information (approximately 81% of all households on
181 the MAF)⁴.

182 In comparison, Axios-Ipsos relies on inverse response propensity sampling from Ipsos’
183 online KnowledgePanel. Ipsos recruits panelists using an address-based probabilistic sample
184 from USPS’s Delivery Sequence File (DSF)⁵. The DSF is similar to the Census’ MAF. Unlike
185 the Census Household Pulse, potential respondents are not limited to the subset for whom email

186 and phone contact information is available. Furthermore, Ipsos provides internet access and
187 tablets to recruited panelists who lack home internet access. In 2021, this “offline” group
188 typically comprises 1% of the final survey (Extended Data Table 1).

189 All three surveys weight on age and gender, i.e., assign larger weights to respondents of
190 underrepresented age by gender subgroups and smaller weights to those of overrepresented
191 subgroups^{2,4,5} (Table 1) Axios-Ipsos and Census Household Pulse also weight on education and
192 race/ethnicity. Axios-Ipsos additionally weights to the composition of political partisanship
193 measured with the ABC News/Washington Post poll in 6 of the 11 waves we study. Education, a
194 known correlate of propensity to respond to surveys³⁶ and social media use³⁷, are notably absent
195 from Delphi-Facebook’s weighting scheme, as is race/ethnicity. As noted before, none of the
196 surveys use the CDC benchmark to adjust or assess estimates of vaccine uptake.

197 Explanations for error

198 Table 2 illustrates some consequences of these design choices. Axios-Ipsos samples
199 mimic the actual breakdown of education attainment among US adults even before weighting,
200 while those of Census Household Pulse and Delphi-Facebook do not. After weighting, Axios-
201 Ipsos and Census Household Pulse match the population benchmark, by design. Delphi-
202 Facebook does not explicitly weight on education, and hence the education bias persists in their
203 weighted estimates: those without a college degree are underrepresented by nearly 20 percentage
204 points. The story is similar for race/ethnicity. Delphi-Facebook’s weighting scheme does not
205 adjust for race/ethnicity, and hence their weighted sample still over-represents White adults by 8
206 percentage points, and under-represents Black and Asian proportions by around 50 percent of
207 their size in the population (Table 2).

208 The overrepresentation of White adults and people with college degrees explains part of
209 the error of Delphi-Facebook. The racial groups that Delphi-Facebook under-represents tend to
210 be more willing and less vaccinated in the samples (Table 2). In other words, re-weighting the
211 Delphi-Facebook survey to upweight racial minorities will bring willingness estimates closer to
212 Household Pulse and the vaccination rate closer to CDC. The three surveys also report that
213 people *without* a 4-year college degree are less likely to have been vaccinated compared to those
214 *with* a degree (Table 2 and Supplemental Information Table 1). If we assume that vaccination
215 behaviors do not differ systematically between non-respondents and respondents *within* each
216 demographic category, under-representation of less-vaccinated groups would contribute to the
217 bias found here. However, this alone cannot explain the discrepancies in all the outcomes.
218 Census Household Pulse weights on both race and education⁴ and still over-estimates vaccine
219 uptake by over ten points in late May of 2021 (Fig. 1b).

220 Delphi-Facebook and Census Household Pulse may be unrepresentative with respect to
221 political partisanship, which has been found to be correlated with vaccine behavior³⁸ and with
222 survey response³⁹, and thus may contribute to observed bias. However, neither Delphi-Facebook
223 nor Census Household Pulse collects partisanship of respondents; Census agencies are prohibited
224 from asking about political preference. Moreover, no unequivocal population benchmark for
225 partisanship exists.

226 Rurality may also contribute to the errors because it correlates with vaccine status⁸ and
227 home internet access⁴⁰. Neither the Census Household Pulse nor Delphi-Facebook weights on
228 sub-state geography, which may mean that adults in more rural areas are less likely to be
229 vaccinated and also underrepresented in the surveys, leading to overestimation of vaccine uptake.

230 Axios-Ipsos weights to metropolitan status and also recruits a fraction of its panelists
231 from an “offline” population of individuals without Internet access. We find that *dropping* these
232 offline respondents ($n = 21$, or 1 percent of the sample) in their March 22, 2021 wave *increases*
233 Axios-Ipsos’ overall estimate of the vaccination rate by 0.5 percentage points, thereby increasing
234 the total error (Extended Data Table 2). However, this offline population is simply too small to
235 explain the entirety of the difference in accuracy between Axios-Ipsos and either the Census
236 Household Pulse (6 percentage points) or Delphi-Facebook (14 percentage points), in this time
237 period.

238 Careful recruitment of panelists is at least as important as weighting. Weighting on
239 observed covariates alone cannot explain or correct the discrepancies we observe. For example,
240 reweighting Axios-Ipsos March 22, 2021 wave using only Delphi-Facebook’s weighting
241 variables (age group and gender) increased the error in their vaccination estimates by 1
242 percentage point, but this estimate with Axios-Ipsos data is still more accurate than that from
243 Delphi-Facebook during the same period (Extended Data Table 2). The Axios-Ipsos estimate
244 with Delphi-Facebook weighting overestimated vaccination by 2 percentage points, whereas
245 Delphi-Facebook overestimated it by 11 percentage points.

246 The key implication is that there is no silver bullet: every small part of panel recruitment,
247 sampling, and weighting matters for controlling the data quality measured as the correlation
248 between an outcome and response, what we call the *ddc*. In multi-stage sampling, which includes
249 for instance the selection of participants followed by non-response, bias in even a *single* step can
250 substantially impact the final result (see Methods “Population size in multi-stage sampling”,
251 Extended Data Table 3). A *total quality control* approach – inspired by the Total Survey Error
252 framework⁴¹ – is a better strategy than trying to prioritize some components over others in order

253 to improve data quality. This emphasis is merely a reaffirmation of the best practice for survey
254 research as advocated by the American Association for Public Opinion Research⁶: “The quality
255 of a survey is best judged not by its size, scope, or prominence, but by how much attention is
256 given to [preventing, measuring and] dealing with the many important problems that can arise.”⁴²

257 **Addressing common misperceptions**

258 The three surveys discussed in this article demonstrate a seemingly paradoxical
259 phenomenon – the two larger surveys that we studied are far more statistically confident, yet also
260 far more biased, than the smaller, more traditional Axios-Ipsos poll. These findings are
261 paradoxical only when we fall into the trap of the long-held, but incorrect, intuition that
262 estimation errors necessarily decrease in larger datasets¹².

263 A limitation of our vaccine uptake analysis is that we only examine *ddc* with respect to an
264 outcome for which a benchmark is available: first-dose vaccine uptake. One might hope that
265 surveys biased on vaccine uptake are not biased on other outcomes, for which there may not be
266 benchmarks to expose their biases. However, the absence of evidence of bias for the remaining
267 outcomes is not evidence of its absence. In fact, mathematically, when a survey is found to be
268 biased with respect to one variable, it implies that the entire survey fails to be *statistically*
269 *representative*. The theory of survey sampling relies on statistical representativeness for all
270 variables achieved via probabilistic sampling⁴³. Indeed, Neyman’s original introduction of
271 probabilistic sampling showed the limits of purposive sampling, which attempted to achieve
272 overall representativeness via enforcing it only on a set of variables^{18,44}.

273 In other words, when a survey loses its overall statistical representativeness (e.g., through
274 bias in coverage or nonresponse), which is difficult to repair (e.g., via weighting or modeling on
275 observable characteristics) and almost impossible to verify⁴⁵, researchers who wish to use the
276 survey for scientific studies must supply other reasons to justify the reliability of their survey
277 estimates, such as evidence about the independence between the variable of interest and the
278 factors that are responsible for the unrepresentativeness. Furthermore, scientific journals that
279 wish to publish studies based on unrepresentative surveys¹⁷, especially those with large sizes
280 such as Delphi-Facebook (biased with respect to vaccination status (Fig. 1), race and education
281 (Table 2)) need to ask for reasonable effort from the authors to address the unrepresentativeness.
282 A simple acknowledgment of the potential bias is insufficient for alerting about potentially
283 seriously flawed findings, as we reveal in this article.

284 Some may argue that bias is a necessary trade-off for having data that is sufficiently large
285 for conducting highly granular analysis, such as county-level estimation of vaccine hesitancy²⁶.
286 While high-resolution inference is important, we warn that this is a double-edged argument. A
287 highly biased estimate with a misleadingly small confidence interval can do more damage than
288 having no estimate at all. We further note that bias is not limited to population point estimates,
289 but also affects estimates of changes over time (contrary to published guidance³) – both Delphi-
290 Facebook and Census Household Pulse significantly overestimate the *slope* of vaccine uptake
291 relative to that of the CDC benchmark (Fig. 1b).

292 The accuracy of our analysis also relies on the accuracy of the CDC's estimates of
293 COVID vaccine uptake. However, if the selection bias in the CDC's benchmark is significant
294 enough to alter our results, then that itself would be yet another example of the Big Data
295 Paradox.

296

Discussion

297 This is not the first time that the Big Data Paradox has reared its head: Google Trends
298 predicted more than twice the number of influenza-like illnesses than the CDC in February
299 2013⁴⁶. This analysis demonstrates that the Big Data Paradox applies not only to organically-
300 collected Big Data, like Google Trends, but also to surveys. Delphi-Facebook is “the largest
301 public health survey ever conducted in the United States”⁴⁷. The Census Household Pulse is
302 conducted in collaboration between the US Census Bureau and eleven statistical government
303 partners, all with enormous resources and survey expertise. Both studies take steps to mitigate
304 selection bias, yet overestimate vaccine uptake by double digits. As we demonstrated, the impact
305 of bias is magnified as relative sample size increases.

306 In contrast, Axios-Ipsos records only about 1,000 responses per wave, but makes
307 additional efforts to prevent selection bias. Small surveys can be just as wrong as large surveys
308 in expectation – of the three other small-to-medium online surveys additionally analyzed, two
309 also miss the CDC vaccination benchmark (Extended Data Fig. 5). The overall lesson is that
310 investing in data quality (particularly during collection, but also in analysis) minimizes error
311 more efficiently than does increasing data quantity. Of course, a sample size of 1,000 may be too
312 small (i.e., leading to unhelpfully large confidence intervals) for the kind of 50-state estimates
313 given by big surveys. However, small area methods that borrow information across subgroups⁴⁸
314 can perform better with better quality, albeit small, data, and it is an open question whether that
315 approach would outperform the large, biased surveys.

316 There are approaches to correct for these biases in both probability and nonprobability
317 samples alike. For COVID-19 surveys in particular, since June 2021, the AP-NORC multi-mode

318 panel has weighted their COVID-19 related surveys to the CDC benchmark, so that the weighted
319 *ddc* for vaccine uptake is zero by design⁴⁹. More generally, there is an extensive literature on
320 approaches for making inferences from data collected from nonprobability samples⁵⁰⁻⁵². Other
321 promising approaches include integrating surveys of varying quality^{53,54}, and leveraging the
322 estimated *ddc* in one outcome to correct bias in others under several scenarios (Supplemental
323 Information D).

324 While more needs to be done to fully examine the nuances of large surveys, organically
325 collected administrative datasets, and social media data, we hope this first comparative study of
326 *ddc* highlights the alarming implications of the *Big Data Paradox* – how large sample sizes
327 magnify the impact of seemingly small defects in data collection, leading to overconfidence in
328 incorrect inferences.

329

331 **Main Text Tables**

	Axios-Ipsos	Census Household Pulse	Delphi-Facebook
Recruitment mode	Address-based mail sample to Ipsos KnowledgePanel	SMS and email	Facebook Newsfeed
Interview mode	Online	Online	Online
Average size	1,000/wave	75,000/wave	250,000/week
Sampling frame	Ipsos KnowledgePanel; internet/tablets provided to ~5% of panelists who lack home internet	Census Bureau’s Master Address File (individuals for whom email / phone contact information is available)	Facebook active users
Vaccine uptake question	“Do you personally know anyone who has already received the COVID-19 vaccine?”	“Have you received a COVID-19 vaccine?”	“Have you had a COVID-19 vaccination?”
Vaccine uptake definition	“Yes, I have received the vaccine”	“Yes”	“Yes”
Other vaccine uptake response options	“Yes, a member of my immediate family,” “Yes, someone else,” “No”	“No”	“No,” “I don’t know”
Weighting variables	Gender by age, race, education, Census region, metropolitan status, household income, partisanship.	Education by age by sex by state, race/ethnicity by age by sex by state, household size	Stage 1: age, gender “other attributes which we have found in the past to correlate with survey outcomes” to FAUB; Stage 2: state by age by gender

Table 1

Figure Legends

Fig 1 | Errors in estimates of vaccine uptake. **a.** Estimates of vaccine uptake for US adults in 2021 compared to CDC benchmark data, plotted by end date of each survey wave. Points indicate each study's weighted estimate of first-dose vaccine uptake, and intervals are 95% CIs using reported standard errors and design effects. Delphi-Facebook has $n = 4,525,633$ across 19 waves, Census Household Pulse has $n = 606,615$ across 8 waves, and Axios-Ipsos has $n = 11,421$ across 11 waves. Delphi-Facebook's CIs are too small to be visible. **b.** Total error $\bar{Y}_n - \bar{Y}_N$, **c.** data defect correlation $\hat{\rho}_{Y,R}$, **d.** data scarcity $\sqrt{\frac{(N-n)}{n}}$, **e.** inherent problem difficulty σ_Y . Shaded bands represent scenarios of +/-5% (darker) and +/-10% (lighter) error in the CDC benchmark relative to reported values (points). **b - e** comprise the decomposition in Equation 1.

Fig 2 | Bias-adjusted effective sample size. The bias-adjusted effective sample size of an estimate (different from the classic Kish effective sample size) is the size of a simple random sample which would have the same Mean Square Error as the observed estimate. Effective sample sizes are shown on the \log_{10} scale. The original sample size was $n = 4,525,633$ across 19 waves for Delphi-Facebook, $n = 606,615$ across 8 waves for Census Household Pulse, $n = 11,421$ across 11 waves for Axios-Ipsos. Shaded bands represent scenarios of +/-5% error in the CDC benchmark relative to point estimates based on actual reported values.

Table 1 | Comparison of survey designs. Comparison of key design choices across Axios-Ipsos, Census Household Pulse, and Delphi-Facebook studies. All surveys target the US adult population. See Extended Data Table 1 for additional comparisons.

360 **Table 2 | Composition of survey respondents by educational attainment and**
361 **race/ethnicity.** Axios-Ipsos: wave ending March 22, 2021, $n = 995$. Census Household Pulse:
362 wave ending March 29, 2021, $n = 76,068$. Delphi-Facebook: wave ending March 27, 2021, $n =$
363 181,949. Benchmark uses the 2019 US Census American Community Survey (ACS), composed
364 of roughly 3 million responses. Right-most column shows estimates of vaccine uptake (Vax),
365 willingness (Will) and hesitancy (Hes) from the Census Household Pulse of the same wave.

366

367

368

Calculation and interpretation of *ddc*

The mathematical expression for Equation 1 is given here for completeness:

$$\bar{Y}_n - \bar{Y}_N = \hat{\rho}_{Y,R} \times \sqrt{\frac{N-n}{n}} \times \sigma_Y \quad (2)$$

The first factor $\hat{\rho}_{Y,R}$ is called the *data defect correlation* (*ddc*)¹. It is a measure of data quality represented by the correlation between the recording indicator R ($R = 1$ if an answer is recorded and $R = 0$ otherwise) and its value, Y . Given a benchmark, the *ddc* $\hat{\rho}_{Y,R}$ can be calculated by substituting known quantities into Equation 2. In the case of a single survey wave of a COVID-19 survey, n is the sample size of the survey wave, N is the population size of US adults from US Census estimates⁵⁵, \bar{Y}_n is the survey estimate of vaccine uptake, and \bar{Y}_N is the estimate of vaccine uptake for the corresponding period taken from the CDC's report of the cumulative count of first doses administered to US adults^{8,13}. We calculate $\sigma_Y = \sqrt{\bar{Y}_N(1 - \bar{Y}_N)}$ because Y is binary (but Equation 2 is not restricted to binary Y).

We calculate $\hat{\rho}_{Y,R}$ by using *total* error $\bar{Y}_n - \bar{Y}_N$, which captures not only selection bias but also any measurement bias (e.g., from question wording). However, with this calculation method, $\hat{\rho}_{Y,R}$ lacks the direct interpretation as a correlation between Y and R , and instead becomes a more general index of data quality directly related to classical design effects (see Methods section “Bias-adjusted effective sample size”).

387 It is important to point out that the increase in ddc does not necessarily imply that the
 388 response mechanisms for Delphi-Facebook and Census Household Pulse have changed over
 389 time. The correlation between a changing *outcome* and a steady response mechanism could
 390 change over time, hence changing the value of ddc . For example, as more individuals become
 391 vaccinated, and vaccination status is driven by individual behavior rather than eligibility, the
 392 correlation between vaccination status and propensity to respond could increase even if
 393 propensity to respond for a given individual is constant. This would lead to large values of ddc
 394 over time, reflecting the *increased impact* of the same response mechanism.

395 Error decomposition with survey weights

396 The data quality framework given by Equations 1 and 2 are a special case of a more
 397 general framework for assessing the actual error of a weighted estimator $\bar{Y}_w = \frac{\sum_i w_i R_i Y_i}{\sum_i w_i R_i}$, where w_i
 398 is the survey weight assigned to individual i . It is shown in Meng¹ that

$$399 \quad \bar{Y}_w - \bar{Y}_N = \hat{\rho}_{Y,R_w} \times \sqrt{\frac{N-n_w}{n_w}} \times \sigma_Y, \quad (3)$$

400 where $\hat{\rho}_{Y,R_w} = \text{Corr}(Y, R_w)$ is the finite population correlation between Y_i and $R_{w,i} =$
 401 $w_i R_i$ (over $i = 1, \dots, N$). The “hat” on ρ reminds us that this correlation depends on the specific
 402 realization of $\{R_i, i = 1, \dots, N\}$. The term n_w is the classical “effective sample size” due to
 403 weighting²³, i.e., $n_w = \frac{n}{(1+CV_w^2)}$, where CV_w is the coefficient of variation of the weights for all
 404 individuals in the observed sample, that is, the standard deviation of weights normalized by their
 405 mean. It is common for surveys to rescale their weights to have mean 1, in which case CV_w^2 is
 406 simply the sample variance of W .

407 When all weights are the same, Equation 3 reduces to Equation 2. In other words, the *ddc*
408 term $\hat{\rho}_{Y,R_w}$ now also takes into account the impact of the weights as a means to combat the
409 selection bias represented by the recording indicator R . Intuitively, if $\hat{\rho}_{Y,R} = \text{Corr}(Y, R)$ is high
410 (in magnitude), then some Y_i 's have a higher chance of entering our data set than others, thus
411 leading to a sample average that is a biased estimator for the population average. Incorporating
412 appropriate weights can reduce $\hat{\rho}_{Y,R}$ to $\hat{\rho}_{Y,R_w}$, with the aim to reduce the impact of the selection
413 bias. However, this reduction alone may not be sufficient to improve the accuracy of \bar{Y}_w because
414 the use of weight necessarily reduces the sampling fraction $f = \frac{n}{N}$ to $f_w = \frac{n_w}{N}$ as well since $n_w <$
415 n . Equation 3 precisely describes this trade off, providing a formula to assess when the reduction
416 of *ddc* is significant to outweigh the reduction of the effective sample size.

417 Measuring the correlation between Y and R is not a new idea in survey statistics (though
418 note that *ddc* is the population correlation between Y and R , not the sample correlation), nor is
419 the observation that as sample size increases, error is dominated by bias instead of variance^{56,57}.
420 The new insight is that *ddc* is a general metric to index the *lack of* representativeness of the data
421 we observe, regardless of whether or not the sample is obtained via a probabilistic scheme, or
422 weighted to mimic a probabilistic sample. As discussed in the the section on addressing common
423 missperception, any single *ddc* deviating from what is expected under representative sampling
424 (e.g., probabilistic sampling) is sufficient to establish the sample is not representative (but the
425 converse is not true). Furthermore, the *ddc* framework refutes the common belief that increasing
426 sample size necessarily improves statistical estimation^{1,58}.

427 Bias-adjusted effective sample size

428 By matching the mean-squared error of \bar{Y}_w with the variance of the sample average from
429 simple random sampling, Meng¹ derives the following formula for calculating a *bias-adjusted*
430 *effective sample size*, or n_{eff} :

$$431 \quad n_{\text{eff}} = \frac{n_w}{N - n_w} \times \frac{1}{E[\hat{\rho}_{Y,R_w}^2]}$$

432 Given an estimator \bar{Y}_w with expected total Mean Squared Error (MSE) T due to data
433 defect, sampling variability, and weighting, this quantity n_{eff} represents the size of a simple
434 random sample such that its mean \bar{Y}_N , as an estimator for the same population mean \bar{Y}_N , would
435 have the identical MSE T . The term $E[\hat{\rho}_{Y,R_w}^2]$ represents the amount of selection bias (square)
436 expected on average from a particular recording mechanism R and a chosen weighting scheme.

437 For each survey wave, we use $\hat{\rho}_{Y,R_w}^2$ to approximate $E[\hat{\rho}_{Y,R_w}^2]$. This estimation is unbiased
438 by design, since we use an estimator to estimate its expectation. Therefore, the only source of
439 error is the sampling variation, which is typically negligible for large surveys, such as for
440 Delphi-Facebook and the Census Household Pulse surveys. This estimation error may have more
441 impact for smaller traditional surveys, such as Axios-Ipsos' survey, an issue we will investigate
442 in subsequent work.

443 We compute $\hat{\rho}_{Y,R_w}$ by using the benchmark \bar{Y}_N , namely, via solving Equation 3 for $\hat{\rho}_{Y,R_w}$,

$$444 \quad \hat{\rho}_{Y,R_w} = \frac{Z_w}{\sqrt{N}}, \quad \text{where} \quad Z_w = \frac{\bar{Y}_w - \bar{Y}_N}{\sqrt{\frac{1 - f_w}{n_w} \sigma_Y^2}}$$

445 We introduce this notation Z_w because it is the quantity that determines the well-known
446 survey efficiency measure, the so-called *design effect*, which is the variance of Z_w for a
447 probabilistic sampling design²³ (when we assume the weights are fixed). For the more general
448 setting where \bar{Y}_w may be biased, we replace the variance by MSE, and hence the bias-adjusted
449 design effect $D_e = E[Z_w^2]$, which is the MSE relative to the benchmark measured in the unit of
450 the variance of an average from a simple random sample of size n_w . Hence $D_I \equiv E[\hat{\rho}_{Y,R_w}^2]$,
451 which was termed as *data defect index*¹, is simply the bias-adjusted design effect *per unit*,
452 because $D_I = \frac{D_e}{N}$.

453 Furthermore, because Z_w is the standardized actual error, it captures any kind of error
454 inherited in \bar{Y}_w . This observation is important because when Y is subject to measurement errors,
455 $\frac{Z_w}{\sqrt{N}}$ no longer has the simple interpretation as a correlation. But because we estimate D_I by $\frac{Z_w^2}{N}$
456 directly, our effective sample size calculation is still valid even when Equation 3 does not hold.

457 **Asymptotic behavior of *ddc***

458 As shown in Meng¹, for any probabilistic sample without selection biases, the *ddc* is on
459 the order of $\frac{1}{\sqrt{N}}$. Hence the magnitude of $\hat{\rho}_{Y,R}$ (or $\hat{\rho}_{Y,R_w}$) is small enough to cancel out the impact
460 of $\sqrt{N-n}$ (or $\sqrt{N-n_w}$) in the data scarcity term on the actual error, as seen in Equation 2 (or
461 Equation 3). However, when a sample is unrepresentative, e.g. when those with $Y = 1$ are more
462 likely to enter the dataset than those with $Y = 0$, then $\hat{\rho}_{Y,R}$ can far exceed $\frac{1}{\sqrt{N}}$ in magnitude. In
463 this case, error will increase with \sqrt{N} for a fixed *ddc* and growing population size N (Equation
464 2). This result may be counter-intuitive in the traditional survey statistics framework, which often

465 considers how error changes as sample size n grows. The *ddc* framework considers a more
466 general setup, taking into account individual response behavior, including its impact on sample
467 size itself.

468 As an example of how response behavior can shape both total error and the number of
469 respondents n , suppose individual response behavior is captured by a logistic regression model

$$470 \quad \text{logit}[\Pr(R = 1|Y)] = \alpha + \beta Y. \quad (4)$$

471 This is a model for a response propensity score. Its value is determined by α , which
472 drives the overall sampling fraction $f = \frac{n}{N}$, and by β , which controls how strongly Y influences
473 whether a participant will respond or not.

474 In this logit response model, when $\beta \neq 0$, $\hat{\rho}_{Y,R}$ is determined by individual behavior, not
475 by population size N . In Supplemental Information B.1, we prove that *ddc* cannot vanish as N
476 grows, nor can the observed sample size n ever approach 0 or N for a given set of (finite and
477 plausible) values of $\{\alpha, \beta\}$, because there will always be a non-trivial percentage of non-
478 respondents. For example, an f of 0.01 can be obtained under this model for either $\alpha =$
479 $-0.46, \beta = 0$ (no influence of individual behavior on response propensity), or for $\alpha = -3.9, \beta =$
480 -4.84 . However, despite the same f , the implied *ddc* and consequently the MSE will differ. For
481 example, the MSE for the former (no correlation with Y) is 0.0004, while the MSE for the latter
482 (a -4.84 coefficient on Y) is 0.242, over 600 times larger.

483 See Supplemental Information B.2 for the connection between *ddc* and a well-studied non-
484 response model from econometrics, the Heckman selection model⁵⁹.

485 Population size in multi-stage sampling

486 We have shown that the asymptotic behavior of error depends on whether the data
487 collection process is driven by individual response behavior or by survey design. The reality is
488 often a mix of both. Consequently, the relevant “population size” N depends on when and where
489 the representativeness of the sample is destroyed, i.e., when the individual response behaviors
490 come into play. Real-world surveys that are as complex as the three surveys we analyze here
491 have multiple stages of sample selection.

492 Extended Data Table 3 takes as an example the sampling stages of the Census Household
493 Pulse, which has the most extensive set of documentation among the three surveys we analyze.
494 As we have summarized (Table 1 and Extended Data Table 1 the Census Household Pulse (1)
495 first defines the sampling frame as the reachable subset of the Master Address File, (2) takes a
496 random sample of that population to prompt (send a survey questionnaire), and (3) waits for
497 individuals to respond to that survey. Each of these stages reduces the desired data size, and the
498 corresponding *population size* is the intended sample size from the prior stage (in notation, $N_s =$
499 n_{s-1} , for $s = 2,3$). For example, for stage 3, the population size N_3 is the size of the intended
500 sample size n_2 from the second stage, i.e., the sampling stage, because only the sampled
501 individuals have a chance to respond.

502 Although all stages contribute to the grand *ddc*, the stage that dominates is the *first stage*
503 *at which the representativeness of our sample is destroyed*— whose size will be labeled as the
504 *dominating population size (dps)*—when the relevant population size decreases dramatically at
505 each step. However, we must bear in mind that *dps* refers to the worse case scenario, when biases
506 accumulate, instead of (accidentally) canceling each other out.

507 For example, if the 20 percent of the MAF excluded from the Census Household Pulse
508 sampling frame (because they had no cell phone or email contact information) is not
509 representative of the US Adult population, then the dps is N_1 , or 255 million adults contained in
510 144 million households. Then the increase in bias for given ddc is driven by the rate of $\sqrt{N_1}$
511 where $N_1 = 2.55 \times 10^8$ and is large indeed (with $\sqrt{2.5 \times 10^8} \approx 15,000$). In contrast, if the the
512 sampling frame is representative of the target population and the outreach list is representative of
513 the frame (and hence representative of the US adult population) but there is non-response bias,
514 then dps is $N_3 = 10^6$ and the impact ddc is amplified by the square root of that number ($\sqrt{10^6} =$
515 1,000). In contrast, Axios-Ipsos reports a response rate of about 50%, and obtains a sample of
516 $n = 1000$, so the dps could be as small as $N_3 = 2000$ (with $\sqrt{2000} \approx 45$).

517 This decomposition is why our comparison of the surveys is consistent with the *Law of*
518 *Large Populations* (estimation error increases with \sqrt{N}), *even though all three surveys ultimately*
519 *target the same US Adult Population*. Given our existing knowledge about online-offline
520 populations⁴⁰ and our analysis of Axios-Ipsos' small "offline" population, Census Household
521 Pulse may suffer from unrepresentativeness at Stage 1 of Extended Data Table 3 where $N = 255$
522 million, and Delphi-Facebook may suffer from unrepresentativeness at the initial stage of
523 starting from the Facebook User Base. In contrast, the main source of unrepresentativeness for
524 Axios-Ipsos maybe at a later stage where the relevant population size is orders of magnitude
525 smaller.

526 **CDC estimates of vaccination rates**

527 The CDC benchmark data used in our analysis was downloaded from the CDC’s COVID
528 data tracker¹³. We employ the cumulative count of people who have received at least one dose of
529 COVID-19 vaccine reported in the “Vaccination Trends” tab. This data set contains vaccine
530 uptake counts for all US residents (not only adults). However, the surveys of interest estimate
531 vaccine uptake among adults. The CDC receives age-group-specific data on vaccine uptake from
532 all states except for Texas on a daily basis, which is also reported cumulatively over time.

533 Therefore, we must impute the number of adults who have received at least one dose on
534 each day. We assume Texas is exchangeable with the rest of the states in terms of the age
535 distribution for vaccine uptake. Under this assumption, for each day, we use the age group
536 vaccine uptake data from all states except for Texas to calculate the proportion of cumulative
537 vaccine recipients who are 18 or older, then we multiply that number by the total number of
538 people who have had at least one dose to estimate the number of US *adults* who have received at
539 least one dose.

540 The CDC performs a similar imputation for the 18+ numbers reported in their COVID
541 data tracker. However the CDC’s imputed 18+ number is available only as a snapshot and not a
542 historical time series, hence the need for our imputation. See Supplemental Information for
543 details of the imputation implementation.

544 **Additional survey methodology**

545 The Census Household Pulse and Delphi-Facebook surveys are the first of their kind for
546 each organization, while Ipsos has maintained their online panel for 12 years.

547 *Question wording*

548 All three surveys ask whether respondents have received a COVID-19 vaccine. See
549 Extended Data Table 1. Delphi-Facebook and Census Household Pulse ask similar questions
550 (“Have you had / received a COVID-19 vaccination / vaccine?”). Axios-Ipsos asks “Do you
551 personally know anyone who has already received the COVID-19 vaccine?,” and respondents
552 are given response options including “Yes, I have received the vaccine.” The Axios-Ipsos
553 question wording might pressure respondents to conform to their communities’ modal behavior
554 and thus misreport their true vaccination status, or may induce acquiescence bias from the
555 multiple “yes” options presented. This pressure may exist both in high- and low-vaccination
556 communities, so its net impact on Axios-Ipsos’ results is unclear. Nonetheless, Axios-Ipsos’
557 question wording does differ from that of the other two surveys, and may contribute the observed
558 differences in estimates of vaccine uptake across surveys.

559 *Population of Interest*

560 All three surveys target US adult population, but with different sampling and weighting
561 schemes. Household Pulse sets the denominator of their percentages as the household civilian,
562 non-institutionalized population in the United States of 18 years of age or older, excluding
563 Puerto Rico or the island areas. Axios-Ipsos designs samples to representative of the US general
564 adult population 18 or older. For Facebook, the US target population reported in weekly
565 contingency tables is the US adult population, excluding Puerto Rico and other US territories.
566 For the CDC Benchmark, we define the denominator as the US 18+ population, excluding Puerto
567 Rico and other US territories. To estimate the size of the total US population, we use the US
568 Census Bureau Annual Estimates of the Resident Population for the United States and Puerto

569 Rico, 2019⁵⁵. This is also what the CDC uses as the denominator in calculating rates and
570 percentages of the US population⁶⁰.

571 Axios-Ipsos and Delphi-Facebook generate target distributions of the US adult population
572 using the Current Population Survey (CPS), March Supplement, from 2019 and 2018,
573 respectively. Census Household Pulse uses a combination of 2018 1-year American Community
574 Survey (ACS) estimates and the Census Bureau's Population Estimates Program (PEP) from July
575 2020. Both the CPS and ACS are well-established large surveys by the Census and the choice
576 between them is largely inconsequential.

577 *Axios-Ipsos Data*

578 The Axios-Ipsos Coronavirus tracker is an ongoing, bi-weekly tracker intended to
579 measure attitudes towards COVID-19 of adults in the US. The tracker has been running since
580 March 13, 2020 and has released results from 45 waves as of May 28, 2021. Each wave
581 generally runs over a period of 4 days. The Axios-Ipsos data used in this analysis was scraped
582 from the topline PDF reports released on the Ipsos website⁵. The PDF reports also contain Ipsos'
583 design effects, which we have confirmed are calculated as 1 plus the variance of the (scaled)
584 weights.

585 *Census Household Pulse Data*

586 The Census Household Pulse is an experimental product of the US Census Bureau in
587 collaboration with eleven other federal statistical agencies. We use the point estimates presented
588 in Data Tables, as well as the standard errors calculated by the Census Bureau using replicate
589 weights. The design effects are not reported, however we can calculate it as $1 + CV_w^2$, where CV_w
590 is the coefficient of variation of the individual-level weights included in the microdata²³.

591 *Delphi-Facebook COVID symptom survey*

592 The Delphi-Facebook COVID symptom survey is an ongoing survey collaboration
593 between Facebook, the Delphi Group at Carnegie Mellon University (CMU), and the University
594 of Maryland². The survey is intended to track COVID-like symptoms over time in the US and in
595 over 200 countries. We use only the US data in this analysis. The study recruits respondents
596 using a daily stratified random samples recruiting a cross-section of Facebook Active Users.
597 New respondents are obtained each day, and aggregates are reported publicly on weekly and
598 monthly frequencies. The Delphi-Facebook data used here was downloaded directly from
599 CMU's repository for weekly contingency tables with point estimates and standard errors.

600 **Ethical compliance**

601 According to HRA decision tools (<http://www.hra-decisiontools.org.uk/research/>), our study is
602 considered Research, and according to the NHS REC review tool (<http://www.hra->
603 [decisiontools.org.uk/ethics/](http://www.hra-decisiontools.org.uk/ethics/)), we do not need NHS Research Ethics Committee (REC) review, as we
604 used only (1) publicly available, (2) anonymized, and (3) aggregated data outside of clinical
605 settings.

606 **Data availability**

607 Raw data is deposited in the Harvard Dataverse, at <https://doi.org/10.7910/DVN/GKBUUK>. Data
608 was collected from publicly available repositories of survey data by downloading it directly or
609 using APIs.

610 **Code availability**

611 Code to replicate the findings is available in the repository <https://github.com/vcbradley/ddc->
612 [vaccine-US](#). The main decomposition of the *ddc* is available on the package “*ddi*” from the
613 Comprehensive R Archive Network (CRAN).

614 **Acknowledgments**

615 We thank Frauke Kreuter, Alex Reinhart, and the Delphi Group at Carnegie Mellon
616 University, Facebook’s Demography and Survey Science group; Frances Barlas, Chris Jackson,
617 Catherine Morris, Mallory Newall, and the Public Affairs team at Ipsos; and Jason Fields and
618 Jennifer Hunter Childs at the US Census Bureau for productive conversations about their
619 surveys. We further thank the Delphi Group at CMU for their help in computing weekly design
620 effects for their survey, the Ipsos team for providing data on their “offline” respondents, and the
621 CDC for responding to our questions. Susan Paddock, other participants at the JPSM 2021
622 lecture (delivered by Meng), and Steve Finch provided helpful comments, which we greatly
623 appreciate. We thank the anonymous reviewers for their constructive comments, which
624 substantially improved our work. We thank Ariel Edwards-Levy for a tweet which originally
625 inspired our interest in this topic, and Rick Born for suggesting more intuitive terms used in
626 Equation 1. V.B. is funded by the University of Oxford’s Clarendon Fund and the EPSRC and
627 MRC through the OxWaSP CDT programme (EP/L016710/1). X-L. M acknowledges partial
628 financial support by NSF. S.F. acknowledges the support of the EPSRC (EP/V002910/1).

629 *Author contributions*

630 V.B. and S.F. conceived and formulated the research questions. V.B. and
631 S.K. contributed equally to data analysis, writing, and visualization. X-L.M. conceived and
632 formulated the methodology. All authors contributed to methodology, writing, visualization,
633 editing, and data analysis. S.F. supervised the work.

634 *Competing Interests*

635 Authors have no competing interests, financial or otherwise.

636

References

- 638 1. Meng, X.-L. Statistical paradises and paradoxes in big data (I): Law of large populations,
639 big data paradox, and the 2016 US presidential election. *The Annals of Applied Statistics*
640 **12**, 685–726 (2018).
- 641 2. Barkay, N. *et al.* Weights and methodology brief for the COVID-19 Symptom Survey by
642 University of Maryland and Carnegie Mellon University, in partnership with Facebook.
643 1–7 (2020).
- 644 3. Kreuter, F. & *et al.* [Partnering with Facebook on a university-based rapid turn-around global survey](#). *Survey*
645 *Research Methods* **14**, 159–163 (2020).
- 646 4. Fields, J. *et al.* Design and operation of the 2020 Household Pulse survey. (2020).
- 647 5. Jackson, C., Newall, M. & Yi, J. Axios Ipsos Coronavirus Index. (2021).
- 648 6. Public Opinion Research (AAPOR), A. A. for. [Best practices for survey research](#). (2021).
- 649 7. Hastak, M., Mazis, M. B. & Morris, L. A. The role of consumer surveys in public policy
650 decision making. *Journal of Public Policy & Marketing* **20**, 170–185 (2001).
- 651 8. B.P. Murthy, *et al.* Disparities in COVID-19 vaccination coverage between urban and
652 rural counties: United States, December 14, 2020 – April 10, 2021. *Morbidity and*
653 *Mortality Weekly Report* (2021) doi:[10.15585/mmwr.mm7020e3](https://doi.org/10.15585/mmwr.mm7020e3).
- 654 9. Arrieta, A., Gakidou, E., Larson, H., Mullany, E. & Troeger, C. Through understanding
655 and empathy, we can convince women to get the COVID-19 vaccine. *Think Global*
656 *Health* (2021).
- 657 10. Japac, L. *et al.* [AAPOR report on Big Data](#). *American Association of Public Opinion Researchers*
658 (2015).
- 659 11. Reinhart, A., Kim, E., Garcia, A. & LaRocca, S. Using the COVID-19 Symptom Survey
660 to track vaccination uptake and sentiment in the United States. (2021).
- 661 12. Mayer-Schönberger, V. & Cukier, K. *Big data: A revolution that will transform how we*
662 *live, work, and think*. (Houghton Mifflin Harcourt, 2013).
- 663 13. CDC. Trends in number of COVID-19 vaccinations. (2021).
- 664 14. Nguyen, K. H. *et al.* Comparison of COVID-19 vaccination coverage estimates from the
665 Household Pulse Survey, Omnibus Panel Surveys, and COVID-19 vaccine administration
666 data, United States, March 2021.
- 667 15. Santibanez, T. A. *et al.* Sociodemographic Factors Associated with Receipt of COVID-19
668 Vaccination and Intent to Definitely Get Vaccinated, Adults aged 18 Years or Above —
669 Household Pulse Survey, United States, April 28–May 10, 2021. (2021).

- 670 16. Kruskal, W. & Mosteller, F. Representative sampling, I: Non-scientific literature.
671 *International Statistical Review/Revue Internationale de Statistique* 13–24 (1979).
- 672 17. Kruskal, W. & Mosteller, F. Representative sampling, II: Scientific literature, excluding
673 statistics. *International Statistical Review/Revue Internationale de Statistique* 111–127
674 (1979).
- 675 18. Kruskal, W. & Mosteller, F. Representative sampling, III: The current statistical
676 literature.
- 677 19. Kruskal, W. & Mosteller, F. Representative sampling, IV: The history of the concept in
678 statistics, 1895-1939. *International Statistical Review/Revue Internationale de Statistique*
679 169–195 (1980).
- 680 20. AAPOR. Margin of Sampling Error/Credibility Interval.
- 681 21. The Delphi Group at Carnegie Mellon University in partnership with Facebook. Topline
682 Report on COVID-19 Vaccination in the United States. (2021).
- 683 22. Haas, E. J. *et al.* Impact and effectiveness of mRNA BNT162b2 vaccine against SARS-
684 CoV-2 infections and COVID-19 cases, hospitalisations, and deaths following a
685 nationwide vaccination campaign in Israel: An observational study using national
686 surveillance data. *The Lancet* (2021).
- 687 23. Kish, L. *Survey Sampling*. (Wiley, 1965).
- 688 24. Institute for Health Metrics and Evaluation (IHME). COVID-19 vaccine hesitancy.
689 (2021).
- 690 25. King, W. C., Rubinstein, M., Reinhart, A. & Mejia, R. J. Time trends and factors related
691 to COVID-19 vaccine hesitancy from january-may 2021 among US adults: Findings from
692 a large-scale national survey. *medRxiv* (2021) doi:[10.1101/2021.07.20.21260795](https://doi.org/10.1101/2021.07.20.21260795).
- 693 26. Centers for Disease Control. Estimates of vaccine hesitancy for COVID-19. (2021).
- 694 27. Groves, R. M. *et al.* *Survey methodology*. vol. 561 (John Wiley & Sons, 2011).
- 695 28. Dempsey, W. [The hypothesis of testing: Paradoxes arising out of reported Coronavirus case-counts](https://arxiv.org/abs/2007.01211). *arXiv* 1–
696 21 (2020).
- 697 29. Isakov, M. & Kuriwaki, S. Towards principled unskewing: Viewing 2020 election polls
698 through a corrective lens from 2016. *Harvard Data Science Review* **2**, (2020).
- 699 30. Hartley, H. O. & Ross, A. Unbiased ratio estimators. *Nature* **174**, 270–271 (1954).
- 700 31. Tiu, A., Susswein, Z., Merritt, A. & Bansal, S. Characterizing the spatiotemporal
701 heterogeneity of the COVID-19 vaccination landscape. (2021).
- 702 32. Groen, J. Sources of Error in Survey and Administrative Data: The Importance of
703 Reporting Procedures. *Journal of Official Statistics* **28**, 173–198 (2012).

- 704 33. Tu, X. M., Meng, X.-L. & Pagano, M. The AIDS epidemic: estimating survival after
705 AIDS diagnosis from surveillance data. *Journal of the American Statistical Association*
706 **88**, 26–36 (1993).
- 707 34. Barnes, O. & Burn-Murdoch, J. COVID response hampered by population data glitches.
708 *Financial Times* (2021).
- 709 35. Kennedy, C. *et al.* Evaluating online nonprobability surveys. *Pew Research Center*
710 (2016).
- 711 36. Kennedy, C. *et al.* [An evaluation of the 2016 election polls in the United States](#). *Public Opinion*
712 *Quarterly* **82**, 1–33 (2018).
- 713 37. Auxier, B. & Anderson, M. Social media use in 2021. *Pew Research Center* (2021).
- 714 38. Gadarian, S. K., Goodman, S. W. & Pepinsky, T. B. Partisanship, health behavior, and
715 policy attitudes in the early stages of the COVID-19 pandemic. *PLOS One* **16**, (2021).
- 716 39. Mercer, A., Lau, A. & Kennedy, C. For Weighting Online Opt-In Samples, What Matters
717 Most? *Pew Research Center* (2018).
- 718 40. Ryan, C. Computer and internet use in the United States: 2016. *American Community*
719 *Survey Reports ACS-39*, (U.S. Census Bureau, Washington, DC, 2017).
- 720 41. Biemer, P. P. & Lyberg, L. E. *Introduction to survey quality*. (John Wiley & Sons, 2003).
- 721 42. Scheuren, F. What is a survey? in (American Statistical Association Alexandria, 2004).
- 722 43. Sukhatme, P. V. *Sampling theory of surveys with applications*. (1954).
- 723 44. Neyman, J. On the two different aspects of the representative method: The method of
724 stratified sampling and the method of purposive selection. *Journal of Royal Statistical*
725 *Society* **97**, 558–625 (1934).
- 726 45. Groves, R. M. Nonresponse rates and nonresponse bias in household surveys. *Public*
727 *opinion quarterly* **70**, 646–675 (2006).
- 728 46. Lazer, D., Kennedy, R., King, G. & Vespignani, A. [The parable of Google Flu: Traps in big data](#)
729 [analysis](#). *Science* **343**, 1203–1205 (2014).
- 730 47. Salomon, J. A. *et al.* The US COVID-19 Trends and Impact Survey, 2020-2021:
731 Continuous real-time measurement of COVID-19 symptoms, risks, protective behaviors,
732 testing and vaccination. (2021).
- 733 48. Park, D. K., Gelman, A. & Bafumi, J. Bayesian multilevel estimation with
734 poststratification: State-level estimates from national polls. *Political Analysis* **12**, 375–
735 385 (2004).
- 736 49. Associated Press-NORC Center for Public Affairs Research. The june 2021 AP-NORC
737 center poll. (2021).

- 738 50. Wang, W., Rothschild, D., Goel, S. & Gelman, A. Forecasting elections with non-
739 representative polls. *International Journal of Forecasting* **31**, 980–991 (2015).
- 740 51. Elliott, M. R. & Valliant, R. Inference for nonprobability samples. *Statistical Science* **32**,
741 249–264 (2017).
- 742 52. Little, R. J., West, B. T., Boonstra, P. S. & Hu, J. Measures of the degree of departure
743 from ignorable sample selection. *Journal of survey statistics and methodology* **8**, 932–
744 964 (2020).
- 745 53. Wiśniowski, A., Sakshaug, J. W., Perez Ruiz, D. A. & Blom, A. G. Integrating
746 probability and nonprobability samples for survey inference. *Journal of Survey Statistics*
747 *and Methodology* **8**, 120–147 (2020).
- 748 54. Yang, S., Kim, J. K. & Song, R. Doubly robust inference when combining probability
749 and non-probability samples with high dimensional data. *Journal of the Royal Statistical*
750 *Society: Series B (Statistical Methodology)* **82**, 445–465 (2020).
- 751 55. Methodology for the United States population estimates: Vintage 2019.
- 752 56. Bethlehem, J. Weighting Nonresponse Adjustments Based on Auxiliary Information.
753 275--288 (New York: Wiley, 2002).
- 754 57. Meng, X.-L. A trio of inference problems that could win you a Nobel prize in statistics (if
755 you help fund it). 537–562 (CRC Press, 2014).
- 756 58. Meng, X.-L. & Xie, X. I got more data, my model is more refined, but my estimator is
757 getting worse! Am I just dumb? *Econometric Reviews* **33**, 218–250 (2014).
- 758 59. Heckman, J. J. Sample selection bias as a specification error. *Econometrica* 153–161
759 (1979).
- 760 60. CDC. Reporting COVID-19 vaccination demographic data. (2021).

761

762

763

764

765

766

Extended Data Legends

Extended Data Fig. 1 | Comparisons of state-level vaccine uptake, hesitancy, and willingness across surveys and the CDC: March 2021. Comparison of state-level point estimates (A-C) and rankings (D-F) for vaccine hesitancy, willingness, and uptake from Delphi-Facebook, and Census Household Pulse. Dotted black lines show agreement and red points show the average of 50 states. Panels G-J compare state-level point estimates and rankings for the same survey waves to CDC benchmark estimates from March 31, 2021. The Delphi-Facebook data is from the week ending March 27, 2021 and the Census Household Pulse is the wave ending March 29, 2021.

Extended Data Fig. 2 | Comparisons of state-level vaccine uptake, hesitancy, and willingness across surveys and the CDC: May 2021. Comparison of state-level point estimates (A-C) and rankings (D-F) for vaccine hesitancy, willingness, and uptake from Delphi-Facebook, and Census Household Pulse. Dotted black lines show agreement and red points show the average of 50 states. Panels G-J compare state-level point estimates and rankings for the same survey waves to CDC benchmark estimates from May 15, 2021. The Delphi-Facebook data is from the wave week ending May 8, 2021 and the Census Household Pulse is the wave ending May 10, 2021.

Extended Data Table 1 | Methodologies of Axios-Ipsos, Census Household Pulse, and Delphi-Facebook studies

Extended Data Fig. 3 | Retroactive adjustment of CDC vaccine uptake figures for April 3-12, 2021, over the 90 days from April 12. Increase is shown as a percentage of the vaccine uptake reported on April 12. Most of the retroactive increases in reported estimates appear to occur in the first 10 days after an estimate is first reported. By about 40 days after the initial estimates for a particular day are reported, the upward adjustment plateaus at around 5-6% of the initial estimate. We use this analysis to guide the choice of 5% and 10% error in the CDC benchmark for our robustness checks.

Extended Data Fig. 4 | Revised estimates of hesitancy and willingness after accounting for survey errors for vaccination uptake. The gray point shows the reported value at the last point of the time series. Each line shows a different scenario for what might be driving the error in uptake estimate, derived using hypothetical *ddc* values for willingness and hesitancy based on the observed *ddc* value for uptake. *Access* scenario: willingness suffers from at least as much, if not more, bias than uptake. *Hesitancy* scenario: hesitancy suffers from at least as much, if not more, bias than uptake. *Uptake* scenario: the error is split roughly equally between hesitancy and willingness. See Supplementary Information D for more details.

Extended Data Fig. 5 | Vaccination Rates compared with CDC benchmark for four online polls. Ribbons indicate traditional 95 percent confidence intervals which are twice the standard error reported by the poll. Data for Progress asks "As of today, have you been vaccinated for Covid-19?"; Morning Consult asks "Have you gotten the vaccine, or not?"; Harris Poll asks "Which of the following best describes your mindset when it comes to getting the COVID-19 vaccine when

810 it becomes available to you?". See Supplementary Information C.3 for more details on each
811 survey and discussion of differences. Gray line is the CDC benchmark.

812
813 **Extended Data Fig. 6 | Survey error by Age Group (18-64 year-olds, and those 65 and over. a.**
814 Estimates of vaccine uptake from Delphi-Facebook (blue) and Census Household Pulse (green)
815 for each 18-64 year-olds (left) and those 65 or older (right). Bounds on the CDC's estimate of
816 vaccine uptake for those groups are shown in gray. The CDC receives vaccination-by-age data
817 only from some jurisdictions. We do know, however, the total number of vaccinations in the
818 U.S. Therefore, we calculate the bounds by allocating all the vaccine doses for which age is
819 unknown to either 18-64 or 65+. **b.** Unweighted *ddc* for each Delphi-Facebook and Census
820 Household Pulse calculated for the 18-64 group using the bounds on the CDC's estimates of
821 uptake. *ddc* for 65+ is not shown due to large uncertainty in the bounded CDC estimates of
822 uptake.

823
824 **Extended Data Table 2 | Contribution of offline recruitment and weighting schemes to**
825 **discrepancies between surveys.** A portion of each Axios-Ipsos wave is recruited from a
826 population with no stable internet connection; Ipsos KnowledgePanel provides tablets to these
827 respondents. In the Axios-Ipsos March 22, 2020 wave, the offline panelists (n = 21) were 24
828 percentage points less likely to be vaccinated than online panelists (n = 974). Weighting the
829 same Axios-Ipsos data (n = 995) to the age and gender target distribution implied by Delphi-
830 Facebook's weights make the vaccination estimates higher by 1 percentage point. However,
831 this number is still lower than Delphi-Facebook's (responses from March 14-20 2020, n =
832 249,954) own estimate of 46%. During this time period, the CDC benchmark vaccination rate
833 was 35.2%. This suggests that the recruitment of offline respondents and different weighting
834 schemes each explains only a small portion of the discrepancy between the two data sources.

835
836 **Extended Data Table 3 | Example of multi-stage population selection.** The *Law of Large*
837 *Populations* described in Methods section "Population size in multi-stage sampling" shows that
838 the population size at the sampling stage where simple random sampling breaks down will
839 dominate the error. This table explains these stages with a concrete example, using the Census
840 Household Pulse. Population and sample sizes for three stages (stage number denoted s = 1, 2,
841 or 3) of sampling of the Census Household Pulse survey data collection process. Approximate
842 sample sizes based on the March 24, 2021 wave. "m" stands for millions and "hh" stands for
843 household. The final row compares the total adult population in the US (255 million adults,
844 made up of 144 million households) to the sample size in one wave of the household pulse. For
845 illustration, we have ignored the impact of unequal sampling probabilities on the sample sizes
846 at each stage.

847