

1 **ARTICLE TYPE**

2 Theory and Methods

3 **TITLE**

4 Missing data matters in participatory syndromic surveillance systems: comparative evaluation
5 of missing data methods when estimating disease burden

6 **AUTHORS**

7 Kristin Baltrusaitis¹, Craig Dalton^{2,3}, Sandra Carlson², Laura F. White⁴

8 **INSTITUTIONS**

9 ¹ Center for Biostatistics in AIDS Research, Harvard T.H. Chan School of Public Health, Boston,
10 United States of America

11 ²Hunter New England Population Health, Wallsend, Australia

12 ³Hunter Medical Research Institute, Newcastle, Australia

13 ⁴Department of Biostatistics, Boston University School of Public Health, Boston University,
14 Boston, United States of America

15 **CORRESPONDING AUTHOR**

16 Kristin Baltrusaitis

17 651 Huntington Avenue, Boston, MA 02115

18 kbaltrus@sdac.harvard.edu

19 **WORD COUNT**

20 3200

21 **KEY WORDS**

22 Participatory Syndromic Surveillance, Digital Epidemiology, Missing Data, Missing Not at
23 Random, Multiple Imputation, Influenza Surveillance

24 **ABSTRACT**

25 **Introduction**

26 Traditional surveillance methods have been enhanced by the emergence of online participatory
27 syndromic surveillance systems that collect health-related digital data. These systems have
28 many applications including tracking weekly prevalence of Influenza-Like Illness (ILI), predicting
29 probable infection of Coronavirus 2019 (COVID-19), and determining risk factors of ILI and
30 COVID-19. However, not every volunteer consistently completes surveys. In this study, we
31 assess how different missing data methods affect estimates of ILI burden using data from
32 FluTracking, a participatory surveillance system in Australia.

33 **Methods**

34 We estimate the incidence rate, the incidence proportion, and weekly prevalence using five
35 missing data methods: available case, complete case, assume missing is non-ILI, multiple
36 imputation (MI), and delta (δ) MI, which is a flexible and transparent method to impute missing
37 data under Missing Not at Random (MNAR) assumptions. We evaluate these methods using
38 simulated and FluTracking data.

39 **Results**

40 Our simulations show that the optimal missing data method depends on the measure of ILI
41 burden and the underlying missingness model. Of note, the δ -MI method provides estimates of
42 ILI burden that are similar to the true parameter under MNAR models. When we apply these

43 methods to FluTracking, we find that the δ -MI method accurately predicted complete, end of
44 season weekly prevalence estimates from real-time data.

45 **Conclusion**

46 Missing data is an important problem in participatory surveillance systems. Here, we show that
47 accounting for missingness using statistical approaches leads to different inferences from the
48 data.

49 **INTRODUCTION**

50 Over the past two decades, traditional surveillance methods have been enhanced by the
51 emergence of online participatory syndromic surveillance systems that collect health-related
52 digital data in near-real time.[1] Through these systems, participants volunteer to regularly
53 report syndromic, health information via online or mobile communication technologies. These
54 systems complement traditional healthcare-based surveillance systems by reducing the time
55 delay associated with visiting a healthcare provider and capturing individuals who do not seek
56 medical care.[2] The first of these systems, de Grote Griepmeting, or the Great Influenza
57 Survey, started in 2003 in the Netherlands and Belgium,[3] and multiple systems throughout
58 Europe,[4] Australia (AU),[5] the United States (US),[6] Mexico, and Brazil have followed. More
59 recently, participatory syndromic surveillance systems have been deployed to track symptoms
60 of Coronavirus 2019 (COVID-19) in the community.[7]

61 Participatory syndromic surveillance systems are being used to forecast weekly prevalence of
62 Influenza-Like Illness (ILI),[9–13] produce age-specific attack rates of ILI,[14–16] predict
63 probable infection of COVID-19,[17] determine risk factors of ILI and COVID-19,[11,17] estimate
64 influenza vaccine effectiveness,[11,18,19] and assess healthcare-seeking behavior.[20–22]

65 These systems actively engage the public in reporting and providing timely information about
66 disease trends within the community, while providing a mechanism for community members to
67 become "citizen-scientists".[2,8] However, because not every participant reports regularly the
68 population at risk can be temporally inconsistent and include systematic biases.[23] For
69 example, in post-influenza season surveys, over 30% of users from the U.S. system, Flu Near
70 You (FNY), reported they were at least slightly more likely to report if they have had symptoms.
71 Furthermore, during the influenza season, FNY users were significantly less likely to consistently
72 submit reports if they reported symptoms during their first report following registration.[24] Of
73 note, the reporting habits of users vary by system. For example, approximately 22% of FNY
74 users submitted at least half of the symptom reports during the 2013-2014 influenza season,[6]
75 whereas the AU system, FluTracking, reported that during the 2017 influenza season, of the
76 participants who completed a survey during the first four survey weeks, 69% completed all
77 available surveys, and 82% completed over 90% of available surveys.[5]
78 The most common approach for addressing the inconsistencies in the user reporting is to select
79 a cohort of "active users", where the definition of "active user" varies by system and study, and
80 assume that all missing reports were asymptomatic.[11,14–16] However, no study has assessed
81 how this deterministic assumption affects estimates or systematically compared approaches
82 through simulations.
83 Here, we draw upon statistical methods for missing data and assess how different missing data
84 methods, such as complete case and imputation-based methods, affect estimates of ILI burden
85 using both simulated data and data from FluTracking.

86 **METHODS**

87 **FluTracking**

88 FluTracking is an online health surveillance system of influenza in AU and, as of 2018, New
89 Zealand and Hong Kong in 2021. Launched in 2006, the FluTracking system has grown to include
90 over 150,000 participants during the peak weeks of COVID-19 surveillance across Australia and
91 New Zealand in 2020.[25] At registration, FluTracking users provide basic demographic
92 information, including month and year of birth, sex, postcode of residence, Indigenous status,
93 highest level of education, and whether or not they work directly with patients in a healthcare
94 setting. After registration, users complete weekly surveys for themselves and other household
95 members about ILI symptoms including fever, cough, and/or sore throat. Users who report any
96 symptoms are asked follow-up questions about absenteeism from work or normal duties, visits
97 to health care providers, and results of laboratory tests. All users are also asked about influenza
98 vaccination. Symptom surveys are sent every Monday, however, unlike other participatory
99 surveillance systems, FluTracking participants have the option to complete missed surveys up to
100 five weeks previous, called "retrospective reports".

101 For this study, ILI is defined as report of both fever and cough, with or without sore throat.

102 Descriptive statistics for participant characteristics, including age, sex, household status, and
103 influenza vaccination status, are displayed as median (25th percentile, 75th percentile) for
104 continuous variables and n (%) for categorical variables for of all participants who submitted at
105 least one symptom report during the 2016, 2017, and 2018 influenza seasons. Although
106 FluTracking collects data from the beginning of May through mid-October, we use reports
107 submitted during the southern hemisphere's influenza surveillance season, defined as

108 Morbidity and Mortality weeks 25 through 41, or approximately late-June through mid-
109 October.

110 Measures of Influenza-Like Illness Burden

111 As recommended by the World Health Organization (WHO), we use the incidence rate (IR) as
112 one measure of ILI burden.[26] The IR is equal to the number of incident ILI reports, defined as
113 a report of ILI in which ILI was not reported the previous week, divided by the total person-time
114 reported by participants:

$$115 \quad IR = \frac{\sum_{i=1}^N \sum_{t=1}^T Y_{it}}{\sum_{i=1}^N \sum_{t=1}^T (1 - R_{it})} \times 10,000 \quad (1)$$

116 where

$$Y_{it} = \begin{cases} 1, & \text{if ILI reported at time } t \\ 0, & \text{otherwise} \end{cases}$$

$$117 \quad R_{it} = \begin{cases} 0, & \text{if } Y_{it} \text{ observed} \\ 1, & \text{otherwise} \end{cases},$$

118 and t indexes the week numbers and ranges from 25 to 41.

119 The rate is expressed per 10,000 person weeks. Because person-time at risk is unavailable for
120 most routine influenza surveillance data, we also present the incidence proportion (IP) for
121 comparability across systems. The IP is equal to the number of participants who reported ILI at
122 least once during the influenza season divided by the total number of participants:

$$123 \quad IP = \frac{\sum_{i=1}^N Q_i}{N} \quad (2)$$

124 where

$$125 \quad Q_i = \begin{cases} 1, & \text{if } \sum_{t=1}^T Y_{it} > 0 \\ 0, & \text{otherwise} \end{cases}.$$

126 The 95% Confidence Intervals (CI) for these estimates are given by

127
$$95\% CI = \left(\frac{IR \text{ or } IP}{e^{1.96/\sqrt{d}}}, IR \text{ or } IP \times e^{1.96/\sqrt{d}} \right), \quad (3)$$

128 where d is the number of cases.[27,28] We calculate these measures of ILI burden for the
129 overall population and by age group (<5, 5-17, 18-49, 50+ years). Finally, we present the weekly
130 prevalence (WP) of ILI at each week, which is calculated by dividing the number of ILI reports by
131 the total number of reports observed,

132
$$WP_t = \frac{\sum_{i=1}^N Y_{it}}{\sum_{i=1}^N (1-R_{it})}. \quad (4)$$

133 Because WP is estimated in near-real time, all retrospective reports were assumed missing
134 when calculating estimates. We assess and compare these measures of ILI burden across three
135 influenza seasons (2016, 2017, and 2018).

136 **Missing Data Methods**

137 We assess five methods that account for missing data:

- 138 1. Available case
- 139 2. Complete case
- 140 3. Assume all missing reports are non-ILI
- 141 4. Multiple imputation (MI)
- 142 5. MI with delta (δ) adjustment

143 The first method uses all available cases (i.e., all reports submitted or select all Y_{it} for $R_{it} = 0$),
144 whereas the second method includes only complete cases (i.e., reports from individuals who
145 submitted all reports or select all individuals, i , where $R_{it} = 0$ for all t). These methods assume
146 that data is Missing Completely at Random (MCAR). In other words, the causes of the missing
147 data are unrelated to the data. The third method assumes that all missing reports are non-ILI

148 reports (i.e., $P(Y_{it} = 1 | R_{it} = 1) = 0$), similar to past studies. The next two methods use MI
149 methods to produce 10 point estimates, which are aggregated using Rubin's rules with a log
150 transformation to account for non-normality.[29] For the first MI method (method 4), we fit a
151 model assuming Missing at Random (MAR),

$$152 \quad \text{logit}[P(Y_t = 1 | \mathbf{X}_t)] = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_{3t} + \beta_4 Y_1 + \dots + \beta_{2+t} Y_{t-1} \quad (5)$$

153 where

$$X_1 = \text{age}, X_2 = \text{sex}, X_{3t} = \text{vaccination status at week } t.$$

154 Under MAR, the probability that a value is missing depends only on the observed data and not
155 the value itself. Because this assumption may not be valid for this data, we also perform MI
156 using an δ adjustment, which is a flexible and transparent method to impute missing data
157 under Missing Not at Random (MNAR) assumptions.[30] For MNAR, the probability that a value
158 is missing depends on the unobserved data. The δ -MI method (method 5) uses (5), however,
159 prior to imputing the missing data, a fixed quantity, δ , is added to the linear predictor of the
160 regression model,

$$161 \quad \text{logit}[P(Y_t = 1 | \mathbf{X}_t)] = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_{3t} + \beta_4 Y_1 + \dots + \beta_{2+t} Y_{t-1} + \delta R_t. \quad (6)$$

162 Here, δ represents the difference in log-odds of ILI between participants who did not report
163 compared with participants who did report.

164 Estimation of delta

165 Because the log-odds of ILI for participants who did not report are unknown, we use the
166 retrospective reports to estimate this value and create annual $\tilde{\delta}$ estimates. In other words,

$$167 \quad \tilde{\delta} = \log \left[\frac{\text{odds}(P(Y=1 | \tilde{R}=1))}{\text{odds}(P(Y=1 | \tilde{R}=0))} \right] \quad (7)$$

168 where,

$$\tilde{R} = \begin{cases} 0, & \text{if } Y \text{ is report for same week} \\ 1, & \text{if } Y \text{ is retrospective report} \end{cases}$$

169 Negative values of $\tilde{\delta}$ indicate that participants were less likely to report ILI for retrospective
170 reports compared with reports submitted during the same week. Data are analyzed using R,
171 version 3.3.2, and both MI methods are fit using the Multivariate Imputation by Chained
172 Equations (MICE) package.[31,32]

173 **Simulations**

174 We evaluate the missing data methods assuming three missingness scenarios: MCAR, MAR, and
175 MNAR, using simulated data. The data is simulated using a three-step process. First, 1000
176 FluTracking populations (n=30,000 each) are simulated using the characteristics (age group, sex,
177 and vaccination status) of the 2016 influenza season participant population. Simulated
178 participants are assigned an age group, sex, vaccination status, and 17 weeks of symptom
179 reports, Y_{it} . These weekly symptom reports are simulated using a multinomial distribution,
180 where n , which is Poisson distributed with an age-group specific mean, represents the total
181 number of ILI reports for the participant and $\mathbf{q} = \{q_1, \dots, q_{17}\}$ is the vector of weekly
182 percentage of the 2016 sentinel general practitioner consultations that were ILI as reported by
183 AU's Department of Health.[33] We simulate 17 missingness indicators, R_{it} , to reflect
184 distribution of FluTracking participant reports (Supplemental Figure 1), using three missingness
185 scenarios:

186 1. MCAR

187
$$R_i \sim \text{Binomial}(n = 17, p_i) \quad (7)$$

$$p_i = \begin{cases} 0.05 & \text{with probability } 0.8 \\ \text{Uniform}(0.1, 0.95) & \text{with probability } 0.2 \end{cases}$$

188 2. MAR

189 $R_{it} \sim \text{Bernoulli}(p_{it})$ (8)

$$p_{it} = \frac{e^{Z_i + \gamma R_{it-1}}}{1 + e^{Z_i + \gamma R_{it-1}}}$$

$$Z_i = \begin{cases} \log\left(\frac{0.05}{1 - 0.05}\right) & \text{with probability } 1 - W_i \\ \text{Uniform}\left(\log\left(\frac{0.1}{1 - 0.1}\right), \log\left(\frac{0.95}{1 - 0.95}\right)\right) & \text{with probability } W_i \end{cases}$$

$$W_{it} = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_7 X_7}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_7 X_7}}$$

$\mathbf{X} = \{\text{age group } (X_1 - X_5), \text{sex, vaccination status}\}$

$\gamma = -0.2$ (i.e., R_{it} is less likely to be observed if $R_{it-1} = 1$)

190 3. MNAR

191 $R_{it} \sim \text{Bernoulli}(p_{it})$ (9)

$$p_{it} = \frac{e^{Z_i + \gamma R_{it-1} + \delta Y_{it}}}{1 + e^{Z_i + \gamma R_{it-1} + \delta Y_{it}}}$$

$\delta = 0.3$ (i.e., R_{it} is more likely to be observed if $Y_{it} = 1$)

192 The values of β_0 through β_7 are estimated for FluTracking data using multivariate logistic
 193 regression, similar to [23]. We define δ equal to 0.3. However, we also present a sensitivity
 194 analysis that assesses how varying the MNAR assumption affects IR estimates (δ equal to 0.3
 195 0.8, 1.3, 3.0, and 5.0). Finally, each of the missing data methods described in the previous
 196 section is applied to produce overall and age-specific IR and IP estimates, and the overall WP
 197 estimates for each simulated dataset. The estimates from each missing data method are
 198 compared to the original simulation parameters qualitatively using violin plots and

199 quantitatively Normalized Root Mean Square Errors (NRMSE) normalized by the original
 200 parameter,

$$201 \quad NRMSE = \frac{\sqrt{\frac{\sum_i^t (\hat{B} - B)^2}{T}}}{B} \quad (10)$$

202 for $T = 1000$. In this equation, B represents the measure of ILI burden (i.e., IR, IP, or WP).

203 RESULTS

204 Simulations

205 Incidence Rate

206 As shown in Figure 1, under each missingness scenario, assuming that all missing reports are
 207 non-ILI underestimates the IR and results in the largest NRMSE (Table 1). Under MCAR and MAR
 208 scenarios, available case and MI methods have smaller NRMSEs compared with the complete
 209 case and δ -MI methods. However, under MNAR, IR estimates using the δ -MI method have
 210 smaller NRMSEs compared with all other methods.

211 **Table 1.** Normalized Root Mean Square Errors, expressed as percentage, by age group for
 212 Incidence Rates (IR) and Incidence Proportions (IP), and overall for Weekly Prevalence (WP)
 213 under Missing Completely at Random (MCAR), Missing At Random (MAR), and Missing Not at
 214 Random (MNAR) scenarios.

Measure of ILI Burden	Missingness Scenario	Age Group	Missing Data Method				
			Available Case	Assume Missing Are Non-ILI	Complete Case	Multiple Imputation	δ Multiple Imputation
Incidence Rate	MCAR	Overall	0.44	14.02	1.52	0.73	2.51
		<5	1.94	14.05	6.82	2.80	3.31
		5-17	1.17	14.07	4.00	1.69	2.87
		18-49	0.65	14.01	2.45	1.02	2.57
		50+	0.67	14.06	2.43	1.05	2.68
	MAR	Overall	0.56	14.59	1.73	0.83	2.57
		<5	2.23	19.47	7.41	3.45	4.36

		5-17	1.26	16.72	4.13	1.98	3.24
		18-49	0.68	15.71	2.45	1.18	2.84
		50+	0.64	12.08	2.34	1.05	2.28
	MNAR	Overall	2.08	12.50	1.61	2.95	0.81
		<5	3.65	16.95	7.57	5.02	3.24
		5-17	2.91	14.45	4.25	3.71	1.77
		18-49	2.64	13.52	2.78	3.19	1.11
		50+	2.18	10.22	2.68	2.65	0.91
Incidence Proportion	MCAR	Overall	-	12.98	1.44	0.64	2.55
		<5	-	12.46	6.05	2.53	3.14
		5-17	-	12.84	3.63	1.58	2.83
		18-49	-	12.95	2.25	0.96	2.59
		50+	-	13.16	2.32	0.96	2.74
	MAR	Overall	-	13.54	1.61	0.69	2.62
		<5	-	17.64	6.69	2.98	4.19
		5-17	-	15.46	3.79	1.76	3.20
		18-49	-	14.65	2.30	1.04	2.91
		50+	-	11.31	2.20	0.96	2.33
	MNAR	Overall	-	11.58	1.51	2.42	0.87
		<5	-	15.29	6.78	4.08	3.00
		5-17	-	13.32	3.89	3.02	1.75
		18-49	-	12.59	2.54	2.66	1.17
		50+	-	9.56	2.48	2.26	0.94
Weekly Proportion	MCAR	Overall	2.17	14.14	-	3.01	4.04
	MAR	Overall	2.17	14.37	-	3.10	3.97
	MNAR	Overall	2.88	12.21	-	4.13	2.98

215

216 Incidence Proportion

217 Under all missingness scenarios, the complete case and imputation methods outperform the

218 available case method (Figure 2, Table 1). For MCAR and MAR scenarios, the MI method has the

219 smallest NRMSEs, the δ -MI method underestimates the IP, and the complete case method has

220 the largest variation in estimates. Similar to IR, the δ -MI method is the best approach when

221 data are MNAR. Estimates for the available case method are not shown because this method is

222 the same as assuming that all missing reports are non-ILI when estimating the IP.

223 Weekly Prevalence

224 When estimating the WP, the available case method outperforms the other methods under
225 each missingness scenario (Figure 3, Table 1). The method that assumes all missing reports are
226 non-ILI underestimates the WP under each missingness scenario, and the δ -MI method
227 outperforms MI when data are MNAR. Because the complete case population is unknown until
228 the end of the season, we did not assess this missing data method.

229 Sensitivity Analysis

230 Sensitivity Analysis results for calculation of overall and age group specific IRs under five MNAR
231 scenarios are shown in Supplemental Figure 2. As the value of δ increases (i.e., ILI reports
232 become less likely to be missing), the method that assumes missing reports are non-ILI
233 becomes a better missing data method.

234 FluTracking

235 During the 2016, 2017, and 2018 influenza seasons, 29,671, 32,778, and 43,389 unique
236 participants submitted at least one symptom report between week 25 and week 41,
237 respectively. Across all influenza seasons, approximately 60% identified as female, and the
238 median age of participants ranged from 47 to 49 years. The largest age group was 50+ years,
239 followed by 18 to 49, 5 to 17, and finally <5. More than half were primary users who submitted
240 reports on their own behalf, and 59%, 61%, and 67% of participants reported that they received
241 the influenza vaccination during the 2016, 2017, and 2018 influenza seasons, respectively
242 (Table 2).

243 **Table 2.** Descriptive statistics of the FluTracking cohort and participant reporting habits during
244 the 2016, 2017, and 2018 influenza seasons. Continuous variables are displayed as median (25th
245 percentile, 75th percentile) and categorical variables are displayed as n (%).

Variable		2016	2017	2018
Participants	n	29,671	32,778	43,389
Sex	Male	11,153 (37.59)	12,665 (38.64)	17,086 (39.38)
	Female	17,267 (58.19)	19,277 (58.81)	25,561 (58.91)
	Other	<5 (0)	5 (0.02)	19 (0.04)
	Unknown	>1,245 (4.22)	831 (2.54)	723 (1.67)
Age (years)	Median (Q1, Q3)	47 (31, 58)	48 (31, 59)	49 (32, 61)
Age Group (years)	<5	963 (3.25)	1,062 (3.24)	1,506 (3.47)
	5 to 17	3,387 (11.42)	3,813 (11.63)	4,961 (11.43)
	18 to 49	11,812 (39.81)	12,306 (37.54)	15,667 (36.11)
	50+	13,509 (45.53)	15,597 (47.58)	21,255 (48.99)
Household	Primary user	17,525 (59.06)	19,170 (58.48)	24,945 (57.49)
	Household member	12,146 (40.94)	13,608 (41.52)	18,444 (42.51)
Vaccinated	Yes	17,526 (59.07)	19,966 (60.91)	29,081 (67.02)
	No	12,145 (40.93)	12,812 (39.09)	14,308 (32.98)
Reports	Total	452,627 (100)	498,465 (100)	665,935 (100)
	Non-ILI	442,817 (97.83)	486,380 (97.58)	655,187 (98.39)
	ILI	9810 (2.17)	12,085 (2.42)	10,748 (1.61)
Reports Submitted Within One Week	Total	398,496 (88.04)	432,733 (86.81)	585,295 (87.89)
	Non-ILI	389,706 (97.79)	422,139 (97.55)	575,611 (98.35)
	ILI	8790 (2.21)	10 594 (2.45)	9684 (1.65)
Retrospective Reports	Total	54,131 (11.96)	66,096 (13.26)	80,640 (12.11)
	Non-ILI	53,111 (98.12)	64,605 (97.74)	79,576 (98.68)
	ILI	1020 (1.88)	1491 (2.26)	1064 (1.32)
Reports per Week Reports per Participant	Median (Q1, Q3)	26,570 (26,470, 26,860)	29,370 (29,200, 29,520)	39,150 (38,550, 39,800)
	Median (Q1, Q3)	17 (16, 17)	17 (16, 17)	17 (16, 17)
delta ($\tilde{\delta}$)		-0.158	-0.082	-0.226
OR ($e^{\tilde{\delta}}$)		0.85	0.92	0.80

247 The descriptive statistics of the reporting habits of participants during the 2016, 2017, and 2018
248 influenza seasons are also shown in Table 2. The total number of symptom reports submitted
249 during the influenza season increased from approximately 450,000 in 2016, to almost 500,000
250 in 2017, and finally to over 650,000 in 2018. The median number of weekly reports also
251 increased from approximately 26,000 in 2016 to 39,000 in 2018. During each influenza season,
252 the median number of reports per participant was 17 (16, 17), indicating that more than half of
253 participants submitted a symptom report each week (Supplemental Figure 1). While most
254 reports were submitted within one week of the symptom report date, a larger proportion of
255 these reports were ILI compared to the retrospective reports, resulting in a negative $\tilde{\delta}$.

256 Although the exact value of $\tilde{\delta}$ varies by season, the corresponding ORs of reporting ILI for a
257 retrospective report compared with reporting ILI for a report submitted the same week are all
258 less than 1. The OR was 0.85 in 2016, 0.92 in 2017, and 0.80 in 2018.

259 As shown in Supplemental Figure 1, most FluTracking participants register before week 25,
260 however, the percentage of participants lost to follow-up increases as the season progresses, as
261 shown by the bright green bars. The percentage of retrospective reports (dark green) is fairly
262 consistent through the season, but the proportion of missing reports (light green) appears to
263 increase until mid-season, at which point it slowly decreases as more participants are lost to
264 follow-up (bright green). These patterns are consistent across all influenza seasons.

265 Incidence Rate

266 Supplemental Table 1 and Supplemental Figure 3 display overall and age-group specific IRs and
267 95% CIs, expressed as number of ILI reports per 10,000 person weeks, by influenza season.
268 Although the 2017 influenza season had higher IRs compared with the 2016 and 2018 seasons,

269 the general patterns in estimates are consistent across all seasons and age groups. The method
270 that assumes that all missing reports are non-ILI has the lowest IR estimates, whereas available
271 case and MI methods have the highest IR estimates. As expected, IR estimates from the δ -MI
272 method are slightly less than estimates from the MI method without the δ adjustment,
273 reflecting that missing reports are less likely to be ILI. In most age groups, IR estimates from the
274 complete case method are similar or slightly greater than estimates from the method that
275 assumes all missing reports are non-ILI.

276 Incidence Proportion

277 The overall and age-group specific IPs and 95% CIs for each influenza season are shown in
278 Supplemental Table 2 and Supplemental Figure 4. Similar to IR estimates, IPs estimates from
279 the method that assumes all missing reports are non-ILI and complete case method are less
280 than IP estimates from the MI and δ -MI methods. However, the differences in IP estimates
281 appear to be less pronounced compared with the differences in IR estimates.

282 Weekly Prevalence

283 Near-real time WP estimates are shown in Supplemental Figure 5. We also present the
284 complete data with the retrospective reports for comparison. The method that assumes all
285 missing reports are non-ILI results in WP estimates lower than the other methods. WP
286 estimates from the MI method are slightly larger than WP estimates from the available case
287 method and the δ -MI method. The estimates from these two methods are similar to the
288 complete data.

289 **DISCUSSION**

290 As we have shown, participatory surveillance systems, due to their design, suffer from

291 substantial missing data that is likely informative. We have described five ways to handle
292 missing data, drawing upon traditional statistical practice, and shown how these perform under
293 different missingness scenarios. Our simulations show that the optimal missing data method
294 depends on the measure of ILI burden and the underlying missingness model. Of note, the δ -MI
295 method provides estimates of ILI burden that are similar to the true parameter under MNAR
296 scenarios. When we apply these methods to a participatory surveillance system in Australia, we
297 find that the δ -MI method accurately predicted end of season WP estimates from real-time
298 data.

299 In 2020, Liu et al. estimated WP of ILI in AU using a Bayesian approach that adjusts FluTracking
300 estimates by the probability of an individual reporting in each week given their past
301 reporting.[34] Similar to our approach, they show that estimates of ILI burden are affected by
302 models that correct for user behavior and that users are more likely to report when ill.
303 However, their weekly behavior-adjusted estimates of prevalence were less than the WP
304 estimates that assume all missing reports were non-ILI during the 2017 influenza season. This
305 difference may indicate that their model assumes that the weekly percentage of ILI in the
306 FluTracking population is greater than the weekly percentage of ILI in the general AU
307 population. Their approach also differs because they do not leverage the information within
308 retrospective reports.

309 While AU estimates of IRs and IPs for ILI are not currently available, laboratory confirmed
310 influenza is a nationally notifiable disease in AU.[35] In 2016 the age-specific rate estimates of
311 laboratory confirmed influenza ranged from 24.67 to 123.72 per 10,000 population.[35] As
312 expected, these estimates are lower than the ILI estimates from FluTracking because the

313 reported number of cases represent only a proportion of the total cases in the community, that
314 is, only those cases for which health care was sought, a test conducted and a diagnosis made,
315 followed by a notification to health authorities and does not include the many other viral and
316 bacterial causes of ILL. But the patterns in age-group specific estimates are similar, as are the
317 weekly trends.

318 This study has several limitations. FluTracking has a highly engaged user population with a
319 smaller proportion of missing reports compared with other systems. Because imputation
320 methods may not be feasible with large amounts of missing data, other systems may need to
321 select a cohort of highly engaged users before implementing MI methods. FluTracking is also
322 unique because it provides users with the opportunity to complete missing surveys. This system
323 accommodation not only adds approximately 10% more weekly reports, but also provides a
324 way to estimate δ for imputation. Additionally, in our MNAR simulation model, because the
325 value of δ is known, the resulting δ -MI model can be properly parameterized. In reality, this
326 value is unknown, and the estimated value is based only on retrospective reports, which may
327 be affected by recall bias. However, our sensitivity analysis shows that under modest changes
328 in δ , for example, increasing δ (OR) from 0.3 (1.35) to 0.8 (2.22) or 1.3 (3.67), the δ -MI method
329 still outperforms the method that assumes missing reports are non-ILL. Finally, only one value of
330 δ was used for imputation for all age groups and sex. In the future, this parameter can be easily
331 updated and adapted during an influenza season to provide real time, demographic-specific
332 estimates. A Bayesian approach to estimation can also be used.

333 National estimates of influenza burden in the population are essential to understand the overall
334 global burden of influenza disease.[37] Alternative data sources, such as FluTracking, have the

335 potential to complement traditional sentinel systems by capturing a population not routinely
336 included among the other healthcare-based systems and identifying healthcare seeking
337 behaviors and testing rates among ill persons.[2] While these surveillance systems now play an
338 important role in surveillance, data analysis needs to be handled carefully to draw appropriate
339 and useful conclusions. The δ -MI method provides a straightforward, interpretable approach to
340 appropriately handling missing data in participatory surveillance systems.

341

342 **ABBREVIATIONS**

343 AU: Australia

344 CI: Confidence Interval

345 COVID-19: Coronavirus 2019

346 FNY: Flu Near You

347 ILL: Influenza-Like Illness

348 IP: Incidence Proportion

349 IR: Incidence Rate

350 MAR: Missing at Random

351 MCAR: Missing Completely at Random

352 MI: Multiple Imputation

353 MICE: Multivariate Imputation by Chained Equations

354 MMWR: Morbidity and Mortality Weekly Report

355 MNAR: Missing Not at Random

356 NRMSE: Normalized Root Mean Square Error

357 US: United States of America

358 WHO: World Health Organization

359 WP: Weekly Prevalence

360 **ACKNOWLEDGEMENTS**

361 The authors acknowledge all the participants who contributed their time and information to the
362 FluTracking system.

363 **COMPETING INTERESTS**

364 The authors have no competing interests to declare.

365 **FUNDING**

366 The FluTracking surveillance system is funded by the Australian Government

367 Department of Health.

368 **REFERENCES**

- 369 1. Choi J, Cho Y, Shim E, Woo H. Web-based infectious disease surveillance systems and
370 public health perspectives: a systematic review. *BMC Public Health* **2016**; 16:1238.
371 Available at: [http://bmcpublichealth.biomedcentral.com/articles/10.1186/s12889-016-](http://bmcpublichealth.biomedcentral.com/articles/10.1186/s12889-016-3893-0)
372 [3893-0](http://bmcpublichealth.biomedcentral.com/articles/10.1186/s12889-016-3893-0).
- 373 2. Smolinski MS, Crawley AW, Olsen JM, Jayaraman T, Crawley AW. Participatory Disease
374 Surveillance: Engaging Communities Directly in Reporting, Monitoring, and Responding
375 to Health Threats Corresponding Author: **2017**; 3.
- 376 3. Marquet RL, Bartelds AIM, van Noort SP, et al. Internet-based monitoring of influenza-
377 like illness (ILI) in the general population of the Netherlands during the 2003-2004
378 influenza season. *BMC Public Health* **2006**; 6:242. Available at:
379 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1609118&tool=pmcentrez&](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1609118&tool=pmcentrez&rendertype=abstract)
380 [rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1609118&tool=pmcentrez&rendertype=abstract).
- 381 4. Paolotti D, Carnahan A, Colizza V, et al. Web-based participatory surveillance of
382 infectious diseases: The InfluenzaNet participatory surveillance experience. *Clin.*
383 *Microbiol. Infect.* **2014**; 20:17–21. Available at: [http://dx.doi.org/10.1111/1469-](http://dx.doi.org/10.1111/1469-0691.12477)
384 [0691.12477](http://dx.doi.org/10.1111/1469-0691.12477).
- 385 5. Moberley S, Carlson SJ, Durrheim DN, Dalton CB, DN. FluTracking: Weekly online

- 386 community based surveillance of influenza-like illness in Australia, 2017 Annual Report.
387 Communicable Diseases Intelligence **2019**
- 388 6. Smolinski MS, Crawley AW, Baltrusaitis K, et al. Flu Near You: Crowdsourced Symptom
389 Reporting Spanning 2 Influenza Seasons. *Am. J. Public Health* **2015**; 105:2124–2130.
390 Available at: <http://ajph.aphapublications.org/doi/10.2105/AJPH.2015.302696>.
- 391 7. Segal E, Zhang F, Lin X, et al. Building an international consortium for tracking
392 coronavirus health status. *Nat. Med.* **2020**; <https://doi.org/10.1038/s41591-020-0929-x>
- 393 8. Kullenberg C, Kasperowski D. What is citizen science? - A scientometric meta-analysis.
394 *PLoS One* **2016**; 11:1–16. Available at: <http://dx.doi.org/10.1371/journal.pone.0147152>.
- 395 9. Santillana M, Nguyen AT, Dredze M, Paul MJ, Nsoesie EO, Brownstein JS. Combining
396 Search, Social Media, and Traditional Data Sources to Improve Influenza Surveillance.
397 *PLOS Comput. Biol.* **2015**; 11:e1004513. Available at:
398 <http://dx.plos.org/10.1371/journal.pcbi.1004513>.
- 399 10. Perrotta D, Bella A, Rizzo C, Paolotti D. Participatory Online Surveillance as a
400 Supplementary Tool to Sentinel Doctors for Influenza-Like Illness Surveillance in Italy.
401 *PLoS One* **2017**; 12:e0169801. Available at:
402 <http://dx.plos.org/10.1371/journal.pone.0169801>.
- 403 11. van Noort SP, Codeço CT, Koppeschaar CE, van Ranst M, Paolotti D, Gomes MGM. Ten-
404 year performance of Influenzanet: ILL time series, risks, vaccine effects, and care-seeking
405 behaviour. *Epidemics* **2015**; 13:28–36. Available at:
406 <http://linkinghub.elsevier.com/retrieve/pii/S1755436515000638>.
- 407 12. Brownstein JS, Chu S, Marathe A, et al. Combining Participatory Influenza Surveillance

- 408 with Modeling and Forecasting: Three Alternative Approaches. *JMIR Public Heal. Surveill.*
409 **2017**; 3:e83. Available at: <http://publichealth.jmir.org/2017/4/e83/>.
- 410 13. Moss R, Zarebski AE, Carlson SJ, McCaw JM. Accounting for Healthcare-Seeking
411 Behaviours and Testing Practices in Real-Time Influenza Forecasts. *Trop. Med. Infect. Dis.*
412 **2019**; 4,12, doi:10.3390/tropicalmed4010012.
- 413 14. Patterson-Lomba O, Van Noort S, Cowling BJ, et al. Utilizing syndromic surveillance data
414 for estimating levels of influenza circulation. *Am. J. Epidemiol.* **2014**; 179:1394–1401.
- 415 15. Chunara R, Goldstein E, Patterson-lomba O, Brownstein JS. Estimating influenza attack
416 rates in the United States using a participatory cohort. **2015**; :1–5.
- 417 16. Stockwell MS, Reed C, Vargas CY, et al. MoSAIC: Mobile surveillance for acute respiratory
418 infections and influenza-like illness in the community. *Am. J. Epidemiol.* **2014**; 180:1196–
419 1201.
- 420 17. Menni C, Valdes AM, Freidin MB, et al. Real-time tracking of self-reported symptoms to
421 predict potential COVID-19. *Nat. Med.* **2020**; <https://doi.org/10.1038/s41591-020-0916-2>
- 422 18. Carlson SJ, Durrheim DN, Dalton CB. FluTracking provides a measure of field influenza
423 vaccine effectiveness, Australia, 2007–2009. *Vaccine* **2010**; 28:6809–6810. Available at:
424 <http://dx.doi.org/10.1016/j.vaccine.2010.08.051>.
- 425 19. Debin M, Colizza V, Blanchon T, Hanslik T, Turbelin C, Falchi A. Effectiveness of 2012 –
426 2013 influenza vaccine against influenza-like illness in general population Estimation in a
427 French web-based cohort. **2014**; 10:536–543.
- 428 20. Tilston NL, Paolotti D, Ealden T. Internet-based surveillance of Influenza-like-illness in the
429 UK during the 2009 H1N1 influenza pandemic. *BMC Public Health* **2010**; 10:1–9.

- 430 21. Peppa M, John Edmunds W, Funk S. Disease severity determines health-seeking
431 behaviour amongst individuals with influenza-like illness in an internet-based cohort.
432 BMC Infect. Dis. **2017**; 17:1–13.
- 433 22. Baltrusaitis K, Reed C, Sewalk K, Brownstein JS, Crawley AW, Biggerstaff M. Health-care
434 seeking behavior for respiratory illness among Flu Near You participants in the United
435 States during the 2015-16 through 2018-19 influenza season. JID. **2020**;
- 436 23. Chunara R, Wisk LE, Weitzman ER. Denominator Issues for Personally Generated Data in
437 Population Health Monitoring. Am. J. Prev. Med. **2017**; 52:549–553. Available at:
438 <http://dx.doi.org/10.1016/j.amepre.2016.10.038>.
- 439 24. Baltrusaitis K, Santillana M, Crawley AW, Chunara R, Smolinski M, Brownstein JS.
440 Determinants of Participants' Follow-Up and Characterization of Representativeness in
441 Flu Near You, A Participatory Disease Surveillance System. JMIR public Heal. Surveill.
442 **2017**; 3:e18. Available at:
443 <http://publichealth.jmir.org/2017/2/e18/>[http://www.ncbi.nlm.nih.gov/pubmed/283](http://www.ncbi.nlm.nih.gov/pubmed/28389417)
444 89417.
- 445 25. FluTracking. Available at: <https://info.FluTracking.net/>. Accessed 19 April 2021.
- 446 26. (WHO) WHO. A Manual for Estimating Disease Burden Associated with Seasonal
447 Influenza. [http://apps.who.int/iris/bitstream/10665/178801/](http://apps.who.int/iris/bitstream/10665/178801/1/9789241549301_eng.pdf?ua=1&ua=1)
448 [1/9789241549301_eng.pdf?ua=1&ua=1](http://apps.who.int/iris/bitstream/10665/178801/1/9789241549301_eng.pdf?ua=1&ua=1). Accessed November 28, 2017. Who **2015**; :124.
449 Available at:
450 http://www.who.int/influenza/resources/publications/manual_burden_of_disease/en/.
- 451 27. Kirkwood BR, Sterne J. Essentials of Medical Statistics. 2nd ed. Blackwell Science Ltd,

- 452 2003.
- 453 28. Giesecke J. Modern Infectious Disease Epidemiology. 2nd ed. London, United Kingdom:
454 Taylor & Francis Ltd, 2002.
- 455 29. Rubin DB. Multiple Imputation for Nonresponse in Surveys. New York: John Wiley and
456 Sons, 2004.
- 457 30. Leacy FP, Floyd S, Yates TA, White IR. Analyses of sensitivity to the missing-at-random
458 assumption using multiple imputation with delta adjustment: Application to a
459 tuberculosis/HIV prevalence survey with incomplete HIV-status data. *Am. J. Epidemiol.*
460 **2017**; 185:304–315.
- 461 31. R Core Team (R Foundation for Statistical Computing). R: A Language and Environment
462 for Statistical Computing. 2016; Available at: <https://www.r-project.org/>.
- 463 32. Buuren S van, Groothuis-Oudshoorn K. **mice**: Multivariate Imputation by Chained
464 Equations in R. *J. Stat. Softw.* **2011**; 45. Available at: <http://www.jstatsoft.org/v45/i03/>.
- 465 33. AUSTRALIAN INFLUENZA Laboratory Confirmed Influenza Activity. **2016**; :1–10. Available
466 at: <http://www.health.gov.au>.
- 467 34. Lui D, Mitchell L, Cope RC, Carlson SJ, and Ross JV. Elucidating User Behaviors in a Digital
468 Health Surveillance System to Correct Prevalence Estimates. *Epidemics.* **2020**; 33
- 469 35. Sullivan SG, Kate Pennington JR, Franklin LJ, et al. A Summary of Influenza Surveillance
470 Systems in Australia, 2015. *Commun. Dis. Intell.* **2016**; :1–51. Available at:
471 [http://www.health.gov.au/internet/main/publishing.nsf/Content/cda-surveil-ozflu-](http://www.health.gov.au/internet/main/publishing.nsf/Content/cda-surveil-ozflu-flucurr.htm/$File/Influenza-Surveillance-Systems-Paper.pdf)
472 [flucurr.htm/\\$File/Influenza-Surveillance-Systems-Paper.pdf](http://www.health.gov.au/internet/main/publishing.nsf/Content/cda-surveil-ozflu-flucurr.htm/$File/Influenza-Surveillance-Systems-Paper.pdf).
- 473 36. National Notifiable Diseases Surveillance System. 2016. Available at:

474 <http://www.health.gov.au/nndssdata>. Accessed 15 July 2018.

475 37. Yang W, Lipsitch M, Shaman J. Inference of seasonal and pandemic influenza
476 transmission dynamics. Proc. Natl. Acad. Sci. **2015**; 112:2723–2728. Available at:
477 <http://www.pnas.org/content/112/9/2723.abstract>

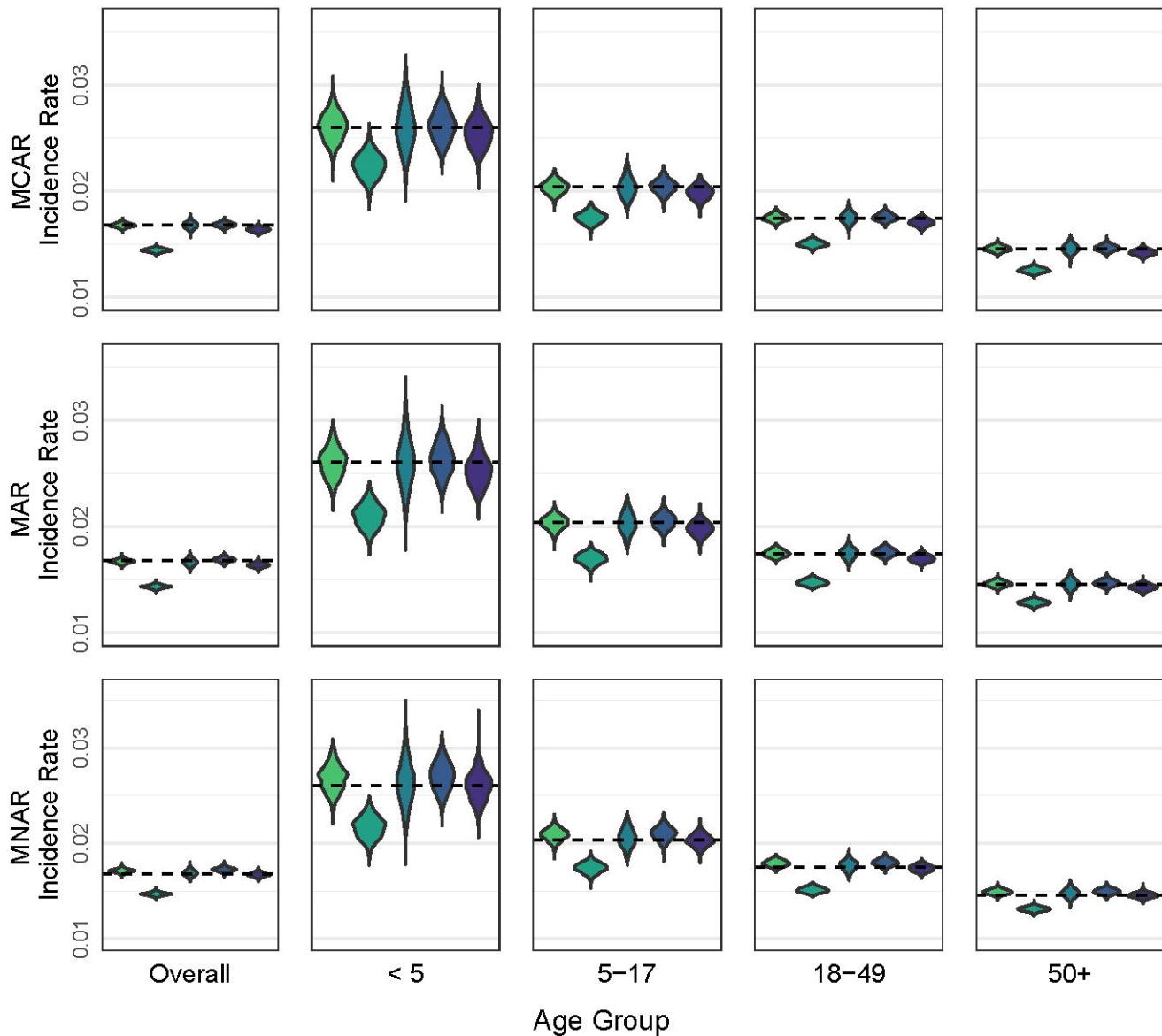
478 **FIGURE LEGENDS**

479 **Figure 1.** Distributions of overall and age-group specific Incidence Rates simulated under three
480 missingness scenarios: Missing Completely at Random (MCAR), Missing at Random (MAR), and
481 Missing Not at Random (MNAR). Dotted line represents the original simulated parameter.

482 **Figure 2.** Distributions of overall and age-group specific Incidence Proportions simulated under
483 three missingness scenarios: Missing Completely at Random (MCAR), Missing at Random
484 (MAR), and Missing Not at Random (MNAR). Dotted line represents the original simulated
485 parameter.

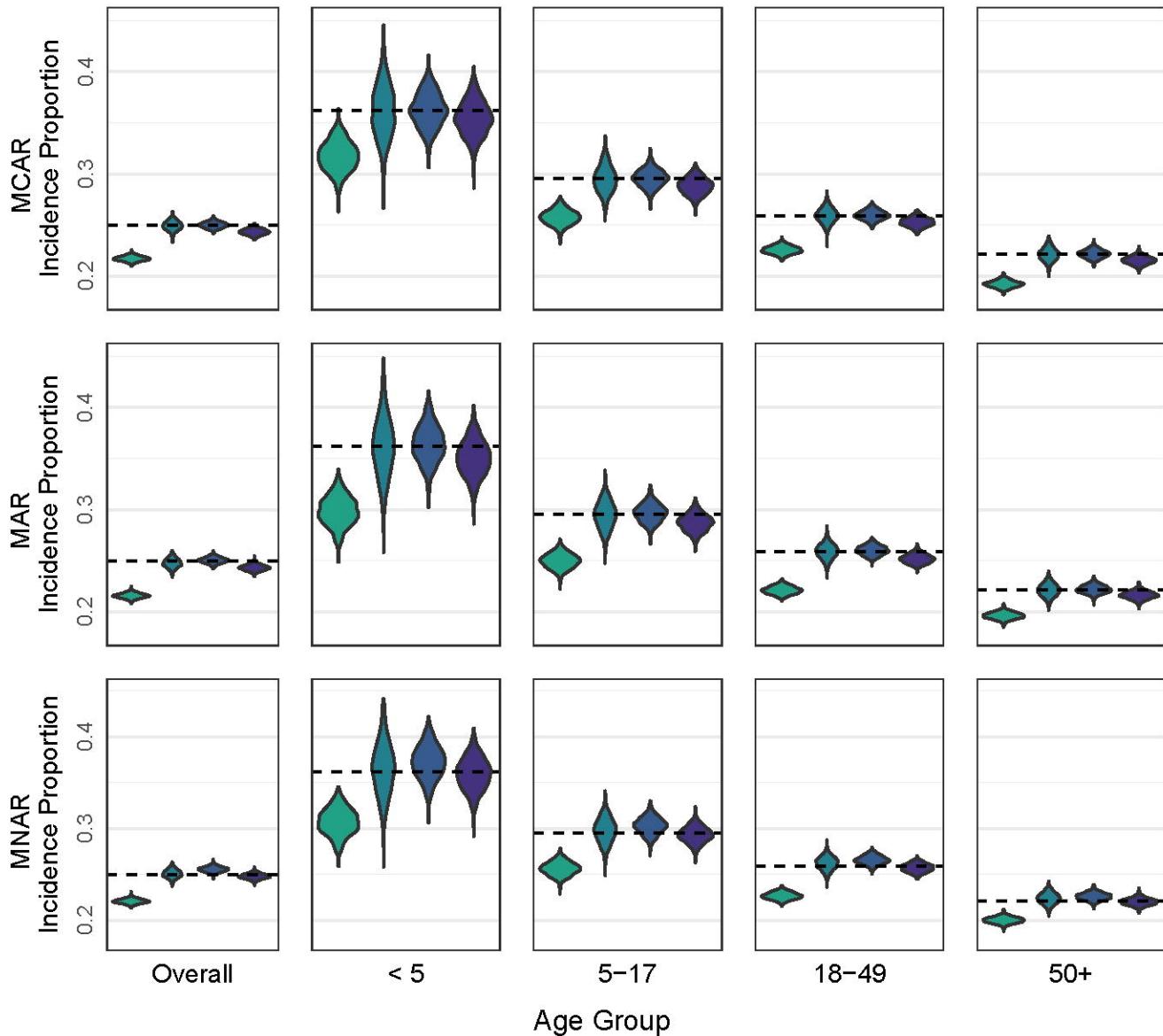
486 **Figure 3.** Time series of Weekly Prevalence simulated under three missingness scenarios:
487 Missing Completely at Random (MCAR), Missing at Random (MAR), and Missing Not at Random
488 (MNAR).

Incidence Rate



Adjustment Method Available Case Missing non-ILI Complete Case MI Delta MI

Incidence Proportion



Adjustment Method



Missing non-ILI



Complete Case

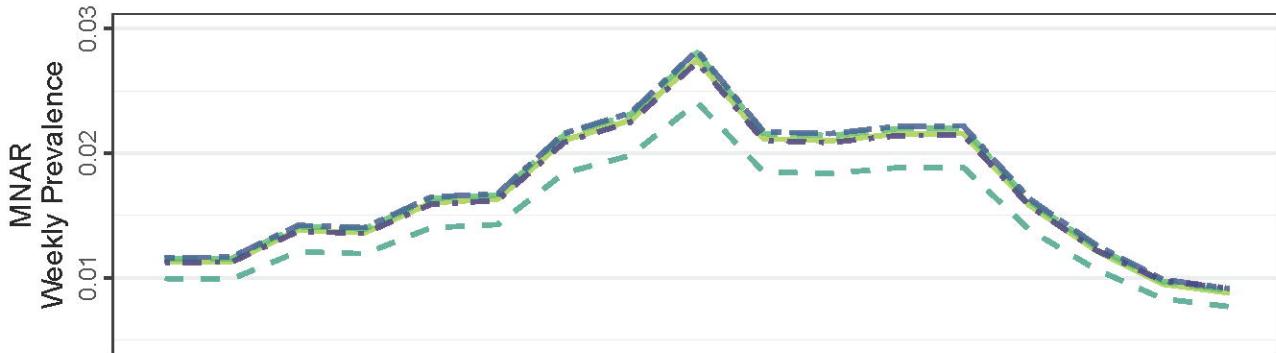
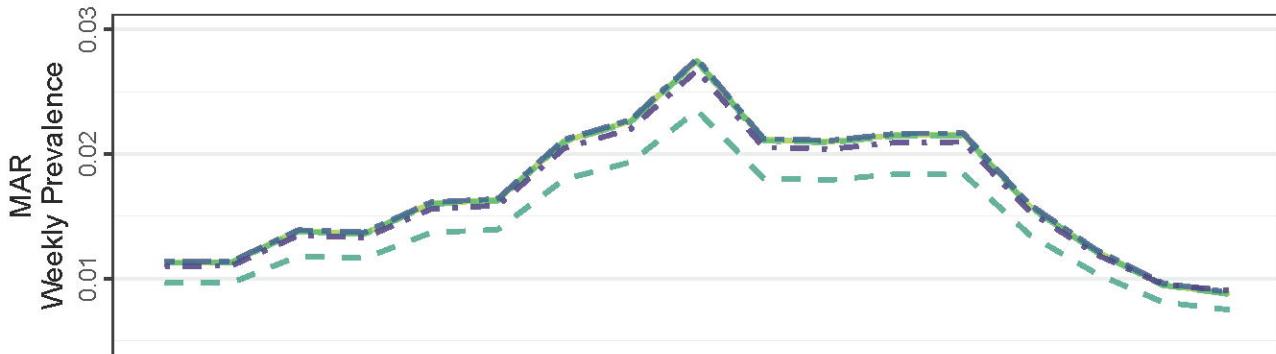
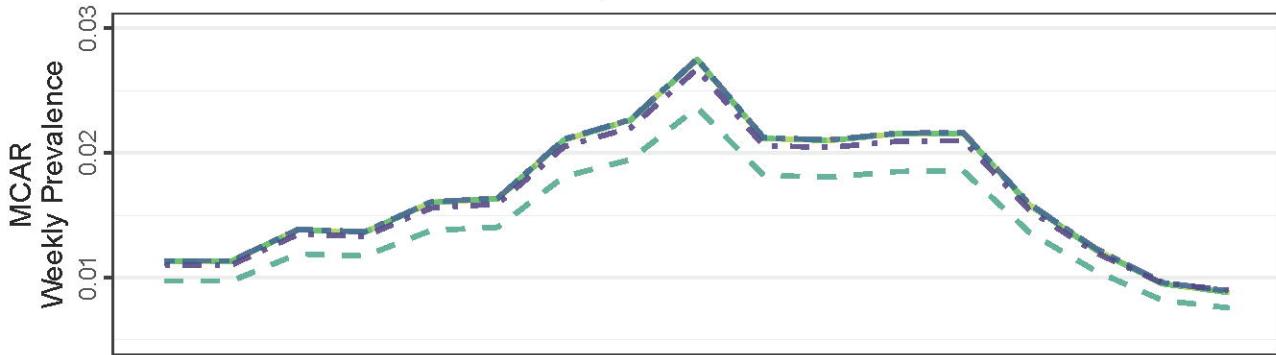


MI



Delta MI

Weekly Prevalence



Week

Adjustment Method Original Available Case Missing non-ILI MI Delta MI