

1 Estimating the epidemic reproduction number from temporally aggregated incidence data: a statistical 2 modelling approach and software tool

3
4 Rebecca K Nash¹ MSc, Anne Cori¹ PhD*, Pierre Nouvellet^{1,2} PhD*

5 *Contributed equally

6 ¹MRC Centre for Global Infectious Disease Analysis, Jameel Institute, School of Public Health, Imperial College
7 London

8 ²School of Life Sciences, University of Sussex
9

10 11 Background

12 The time-varying reproduction number (R_t) is an important measure of epidemic transmissibility; it can directly
13 inform policy decisions and the optimisation of control measures. EpiEstim is a widely used software tool that
14 uses case incidence and the serial interval (SI, time between symptoms in a case and their infector) to estimate
15 R_t in real-time. The incidence and the SI distribution must be provided at the same temporal resolution, which
16 limits the applicability of EpiEstim and other similar methods, e.g. for pathogens with a mean SI shorter than
17 the frequency of incidence reporting.

18 Methods

19 We use an expectation-maximisation algorithm to reconstruct daily incidence from temporally aggregated
20 data, from which R_t can then be estimated using EpiEstim. We assess the validity of our method using an
21 extensive simulation study and apply it to COVID-19 and influenza data. The method is implemented in the
22 opensource R package EpiEstim.

23 Findings

24 For all datasets, the influence of intra-weekly variability in reported data was mitigated by using aggregated
25 weekly data. R_t estimated on weekly sliding windows using incidence reconstructed from weekly data was
26 strongly correlated with estimates from the original daily data. The simulation study revealed that R_t was well
27 estimated in all scenarios and regardless of the temporal aggregation of the data. In the presence of weekend
28 effects, R_t estimates from reconstructed data were more successful at recovering the true value of R_t than
29 those obtained from reported daily data.

30 Interpretation

31 R_t can be successfully recovered from aggregated data, and estimation accuracy can even be improved by
32 smoothing out administrative noise in the reported data.

33 Funding

34 MRC doctoral training partnership, MRC centre for global infectious disease analysis, the NIHR HPRU in
35 Modelling and Health Economics, and the Academy of Medical Sciences Springboard, funded by the AMS,
36 Wellcome Trust, BEIS, the British Heart Foundation and Diabetes UK.
37

38 Introduction

39
40 As infectious disease outbreaks become more common, it is increasingly important to rapidly characterise the
41 threat of emerging and re-emerging pathogens.¹ Transmissibility, i.e. a pathogen's ability to spread through a
42 population, can be quantified using the time-varying reproduction number, R_t , defined as the average number
43 of infections that are caused by a primary case at time t of an outbreak. R_t signals whether an outbreak is
44 growing ($R_t > 1$) or declining ($R_t < 1$), and whether current interventions are sufficient to control the spread of
45 the disease.

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

46

47 One of the most popular tools for real-time R_t estimation, the R package EpiEstim, relies on observing the
48 incidence data and supplying an estimated serial interval (SI) distribution – the time between symptom onset
49 in a case and their infector. EpiEstim requires that the SI distribution and incidence data are supplied using
50 the same time units. This can be problematic when daily incidence data is not reported, which is common for
51 many diseases, such as influenza, Zika virus disease, and most notifiable diseases in countries such as the UK
52 and the US.²⁻⁵ Additionally, several studies intentionally aggregate data to reduce the impact of daily reporting
53 variability; administrative noise, such as “weekend effects”, are characterised by a drop in reported cases over
54 weekends, due to reduced care seeking and longer delays in reporting, followed by a peak on Mondays.^{6,7} A
55 commonly used workaround is to aggregate the SI distribution to match the frequency of incidence
56 reporting,^{8,9} however this is not possible if the SI is shorter than the aggregation of data. For example,
57 influenza-like illness is typically reported on a weekly basis, but influenza has an estimated mean SI of 2-4
58 days.^{10,11} Similarly, reporting of COVID-19, which has an estimated SI of 3-7 days, has typically moved from
59 daily to weekly.^{12,13} Therefore, enabling estimation of R_t from temporally aggregated data is critical to ensure
60 methods such as EpiEstim are widely applicable.¹⁴

61

62 In this study, we combine an expectation-maximisation (EM) algorithm with the renewal equation approach
63 implemented in EpiEstim to reconstruct daily incidence from aggregated data and estimate R_t . We assess the
64 performance of the method using influenza and COVID-19 data, in addition to an extensive simulation study.

65

66 **Methods**

67

68 *EpiEstim*

69

70 EpiEstim uses the renewal equation (eq.1), a form of branching process model.¹⁵ In this formulation, the
71 incidence of new symptomatic cases at time t (I_t) is approximated by a Poisson process, where I_{t-s} is the past
72 incidence, and g_s is the probability mass function of the serial interval.

73

$$I_t \sim Pois \left(R_t \sum_{s=1}^t I_{t-s} g_s \right) \quad (1)$$

74

75 With EpiEstim, R_t can be assumed to remain constant within user defined time windows, which smooth out
76 estimates.

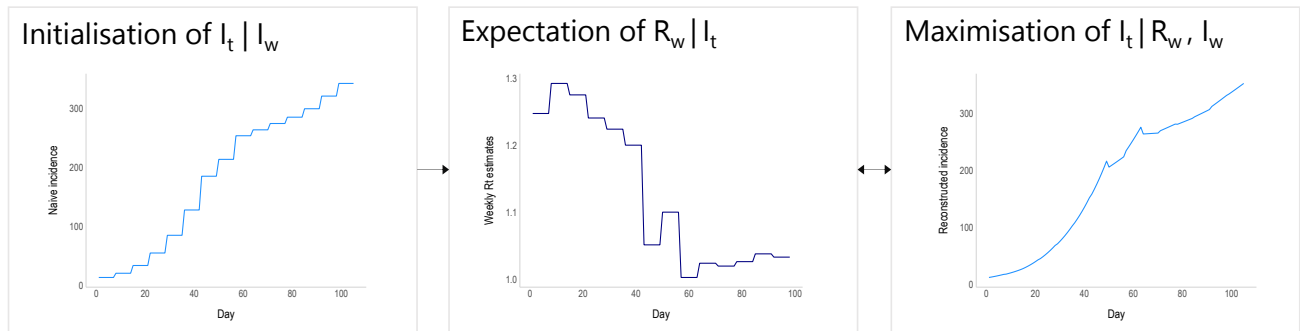
77

78 *Extending EpiEstim for coarsely aggregated data*

79

80 We extended EpiEstim to estimate R_t from aggregated incidence data (I_w), where each aggregation window
81 (w) is >1 day, whilst still conditioning on an assumed serial interval distribution (g_s). We use an EM algorithm
82 to iteratively reconstruct daily incidence (I_t) from I_w , and in turn estimate R_t . We present the method with
83 weekly data in mind, but the method and software can be applied to any temporal aggregation (Figure 1 &
84 appendix pp22-24). The algorithm involves three steps: initialisation, expectation, and maximisation.

85



86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

Figure 1. Schematic of the EM algorithm approach used to reconstruct daily incidence (I_t) from weekly aggregated incidence data (I_w). The algorithm is initialised with a naive disaggregation of the weekly incidence (assuming constant daily incidence throughout the aggregation window). The resulting daily incidence is then used to estimate the reproduction number for each aggregation window, in this case for each week, R_w . R_w is converted into a growth rate (see eq. 2), which is in turn used to reconstruct daily incidence data, whilst ensuring that if I_t were to be reaggreated it would still sum to the original weekly totals. The process cycles between the expectation and maximisation steps until convergence.

Initialisation

The algorithm is initialised with naively disaggregated incidence data. For weekly data, the total incidence for each week is split evenly over 7 days (allowing for non-integers).

Expectation

The current reconstructed I_t is used to estimate the expected reproduction number for each aggregation window, R_w , obtained as the posterior mean from EpiEstim.¹⁶

Maximisation

Conditional on R_w , we reconstruct the most likely I_t . First, R_w is translated into a daily growth rate for that week (r_w), using Wallinga and Lipsitch's method:¹⁷

$$R_w = \frac{1}{\sum_{s=0}^{+\infty} e^{-r_w s} g(s)} \quad (2)$$

I_t for that week is then computed assuming exponential growth, with a multiplying constant k_w ensuring that when reaggreated, the reconstructed I_t matches the original I_w :

$$I_t = k_w \exp(r_w a_t) \quad (3)$$

$$k_w = \frac{I_w}{\sum_{n=1}^7 \exp(r_w n)} \quad (4)$$

where t is time (in days) and a_t is an index representing the day of the aggregation window, e.g. taking values 1 to 7.

The process is repeated iteratively until convergence, at which point I_t can be used to estimate the full posterior distribution of R_t using EpiEstim. For this final step, R_t can be estimated on any time window.

117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

Case studies

We chose datasets where incidence data was available daily, and then artificially aggregated them to weekly counts. R_t was estimated from daily incidence that was reconstructed from weekly aggregated data using our new approach, and compared to R_t estimates obtained from the reported daily incidence using the original EpiEstim R package. All R_t estimates were made using both daily and weekly sliding time windows, and we refer to those estimates as daily R_t estimates and weekly R_t estimates respectively.

We considered three characteristics: 1) mean R_t estimates, 2) uncertainty in the R_t estimates, and 3) the classification of R_t as increasing, uncertain or declining (appendix pp8-9). To compare the performance of this approach to the original method, we assessed the correlations between each of the three characteristics when using the reported and reconstructed incidence.

The priors for R_w and R_t were set to a mean and standard deviation of 5.

Influenza

We obtained a five-week subset of a dataset (11th December 2009 – 14th January 2010) on US active component military personnel (employed by the military as their full-time occupation) that made an outpatient visit to a permanent military treatment facility describing a respiratory-related illness. This daily incidence by date of presentation at a clinic was originally obtained by Riley et al. from the Armed Forces Health Surveillance Center and were digitally extracted for use here.¹⁸ We used a mean SI of 3.6 days and SD of 1.6 days.¹⁰

COVID-19

Incidence of UK COVID-19 cases and deaths were taken from the UK government website.¹⁹ For COVID-19 cases, we obtained ninety-seven weeks of data (21st February 2020 to 30th December 2021) for incidence by date of specimen, which is the date that a sample was taken from an individual which later tested positive. For COVID-19 deaths, we used ninety-six weeks of data (2nd March 2020 to 2nd January 2022) for incidence by date of death within twenty-eight days of a positive test. We assumed a mean SI of 6.1 days and SD of 4.2 days.¹²

Simulation study

We considered scenarios where R_t either remained constant or varied over time, with a stepwise or gradual change. For each scenario, one hundred seventy-day epidemic trajectories were simulated using a Poisson branching process as implemented in the R package projections.²⁰ Daily datasets were aggregated weekly and used to estimate R_t using the proposed method; these values were compared to R_t estimates obtained from simulated daily data using the original EpiEstim R package. We explored the impact of weekend effects on R_t estimates, the ability to supply alternative temporal aggregations of data e.g., three-day, ten-day, or two-weekly aggregations, and finally, the number of iterations required to reach convergence when reconstructing daily incidence data. The full simulation study description and details can be found in the appendix.

Role of funding source

162 The funders of the study had no role in the study design, data collection, data analysis, data interpretation,
163 or writing of the report.

164

165 Results

166

167 Hereafter, we refer to reported and reconstructed incidence data, these are the reported daily incidence
168 and the daily incidence that has been reconstructed from weekly aggregated data, respectively.

169

170 *Influenza*

171

172 The reconstructed incidence of influenza was much smoother than the reported incidence, which showed
173 clear weekend effects and lower reported cases on two public holidays, both occurring on Fridays (Figure 2A
174 & appendix p8). Considering weekly sliding R_t first, there was a high correlation in both the mean R_t estimates
175 derived from each dataset ($R^2 = 0.91$, Figure 2C & appendix p2) and their associated uncertainty ($R^2 = 0.93$ &
176 Figure 2C). The overall agreement in the classification of R_t reached 81.8% (see methods & appendix p9).

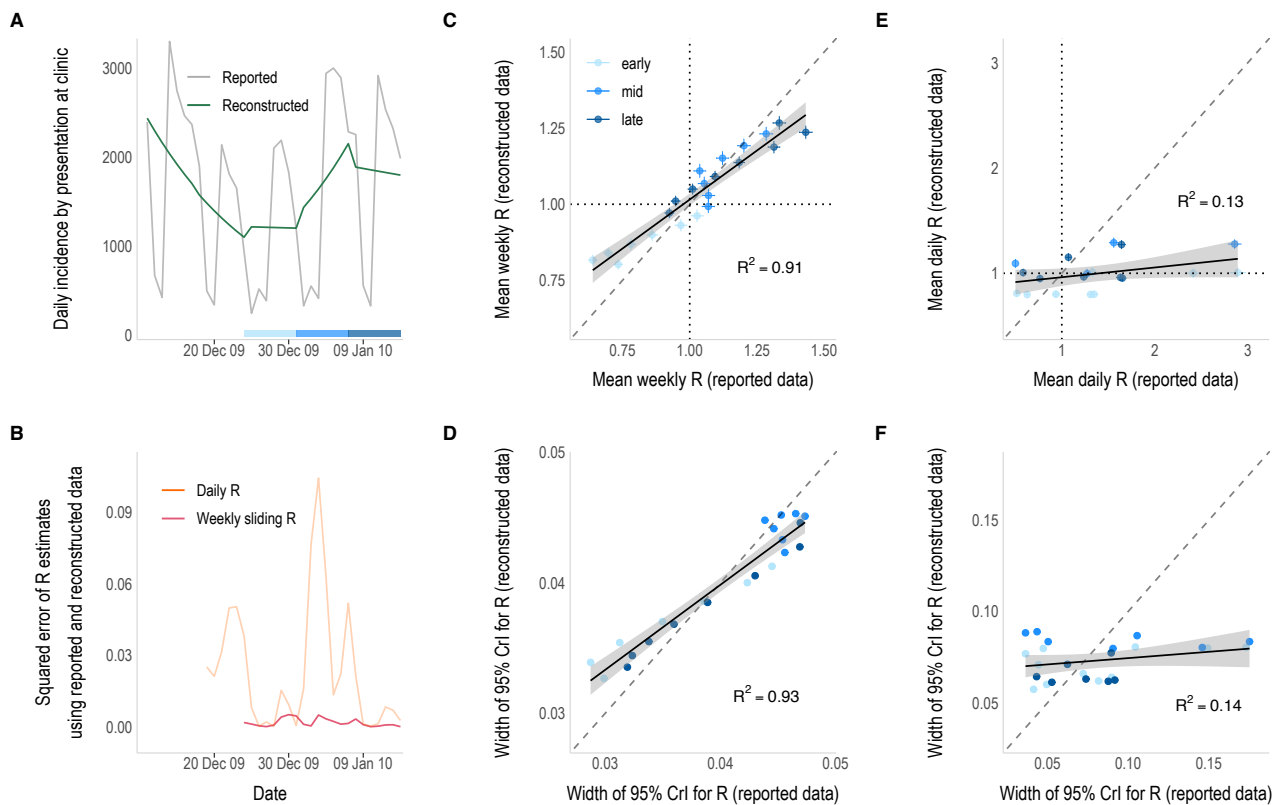
177

178 In contrast, mean daily R_t estimates differed markedly depending on whether the reported or reconstructed
179 data were used, with an R^2 of 0.13 and much higher mean R_t and uncertainty in estimates obtained from
180 reported data (Figure 2E-F). Higher mean R_t estimates coincided with large peaks in the reported daily
181 incidence (typically on Mondays), as daily R_t estimates were not smoothed and therefore more affected by
182 intra-weekly variability (appendix p2). The overall agreement in the classification of daily R_t estimates was
183 much lower, with only 44.4% agreement (appendix p9).

184

185 In this case study, the greatest differences in R_t estimates tended to correspond to time periods when the
186 reported and reconstructed incidence data were most dissimilar (Figure 2B & appendix p3). There was no
187 apparent pattern in the estimates with regard to the outbreak phase, i.e. early, mid or late-phase, but this is
188 likely due to this dataset being a snapshot of incidence taken from within an established epidemic (Figure 2).

189



190
191

192 Figure 2. R_t estimates from daily incidence that was either reported or reconstructed from weekly aggregated
193 influenza data. A) The reported (grey) and reconstructed (green) daily incidence of influenza by date of
194 presentation at a military clinic. B) Squared error of the daily (orange) and weekly sliding (pink) R_t estimates
195 that were made from reconstructed daily data compared to those obtained from the reported daily data. R_t
196 estimation starts on the first day of the second aggregation window (day 8 – 18th December 2009) and is
197 plotted on the last day of the time window used for estimation (i.e., starting on day 9 (19th December) for daily
198 estimates and day 14 (24th December) for weekly estimates). Note: the x-axis is shared with the incidence plot
199 above. C & E) Correlation between the weekly sliding (C) and daily (E) mean R_t estimates using reconstructed
200 data (y-axis) and reported daily data (x-axis). Vertical and horizontal lines depict the 95% credible intervals and
201 dotted lines show the threshold of $R_t = 1$. D & F) Correlation between the uncertainty in the weekly sliding (D)
202 and daily (F) R_t estimates, defined as the width of the 95% credible intervals, using the reconstructed (y-axis)
203 and reported (x-axis) daily data. The colour of the points in panels C-F correspond to the epidemic phase, i.e.
204 the early (19th – 30th December for daily estimates, or 24th – 30th December for weekly sliding estimates),
205 middle (31st December – 6th January) or late (7th – 14th January) phase of the data, shown by the strip in panel
206 A. Solid lines show the linear model fit with 95% confidence intervals (grey shading). Dashed lines represent
207 the $x = y$ line.

208

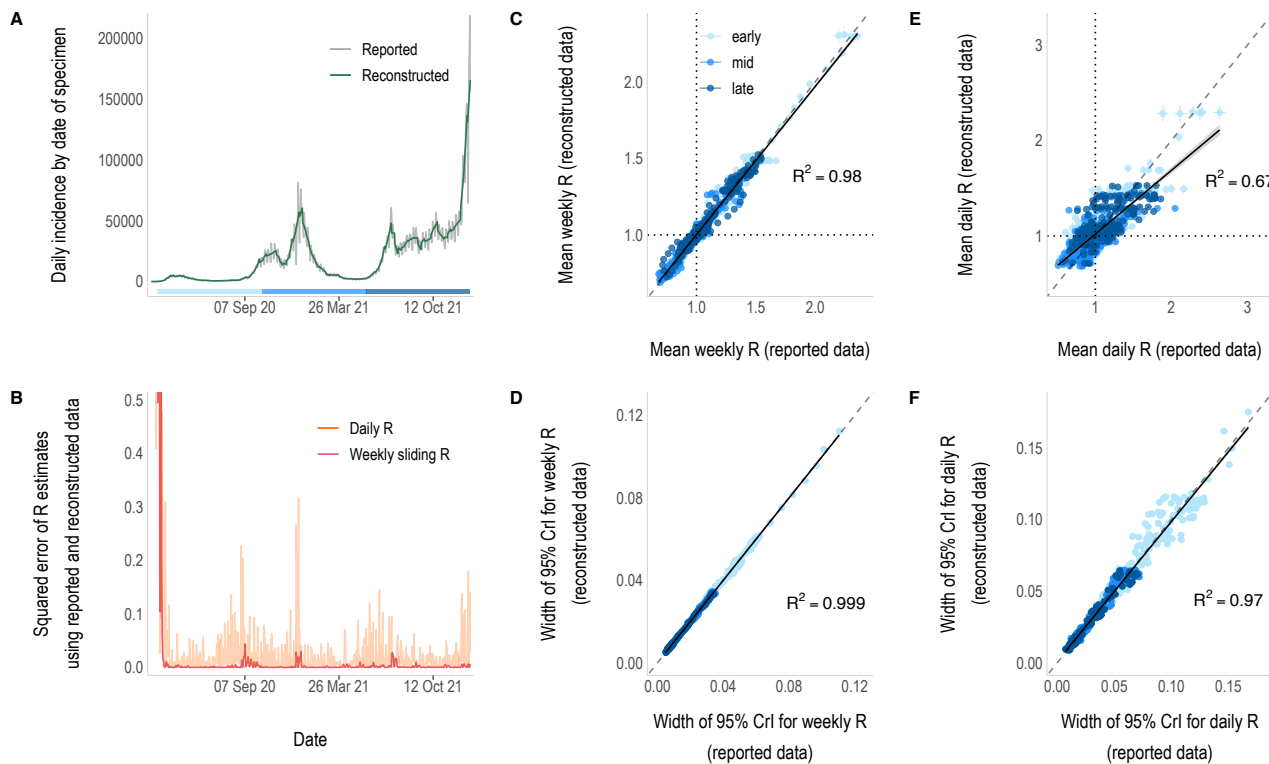
209

210 *COVID-19 cases*

211

212 The reconstructed incidence of COVID-19 smoothed out intra-weekly variability, caused by factors such as
213 weekend effects (Figure 3A & appendix pp7-8). Weekly sliding R_t estimates obtained from reconstructed and
214 reported incidence were similar, both in their means ($R^2 = 0.98$) and their level of uncertainty ($R^2 = 0.99$, Figure
215 3C-D & appendix p4). Mean daily R_t estimates were less well correlated ($R^2 = 0.67$), although the difference is

216 less marked than in the influenza case study (Figure 3E), and the uncertainty in the estimates was similar
 217 across both approaches ($R^2 = 0.97$, Figure 3F). Most of the discrepant R_t estimates and higher levels of
 218 uncertainty coincide with the early phase of the outbreak when incidence was lower (Figure 3E-F). Outside of
 219 periods of low incidence, the largest differences in R_t estimates tended to correspond to time periods with
 220 greater disparities between the reported and reconstructed incidence data (Figure 3B & appendix p5). The
 221 overall agreement in the classification of R_t estimates was higher than for influenza, with 74.4% and 94.9%
 222 agreement for daily and weekly sliding R_t estimates respectively (appendix p9).
 223
 224



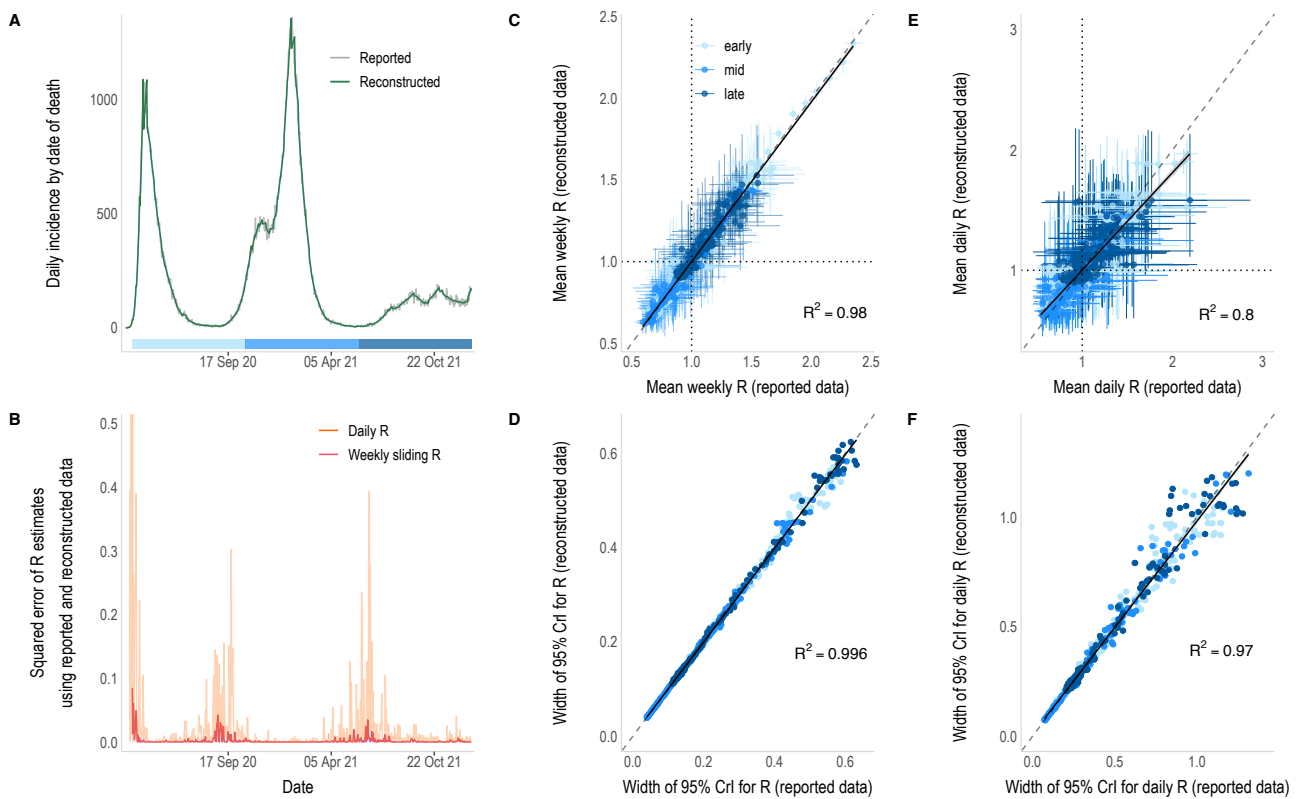
225
 226
 227 Figure 3. R_t estimates from daily incidence that was either reported or reconstructed from weekly aggregated
 228 COVID-19 case data. A) The reported (grey) and reconstructed (green) daily incidence of COVID-19 by date of
 229 specimen. B) Squared error of the daily (orange) and weekly sliding (pink) R_t estimates made from
 230 reconstructed data compared to those obtained from the reported daily data. R_t estimation starts on the first
 231 day of the second aggregation window (day 8 – 28th February 2020) and is plotted on the last day of the time
 232 window used for estimation (i.e., starting on day 9 (29th February) for daily estimates and day 14 (5th March)
 233 for weekly estimates). Note: the x-axis is shared with the incidence plot above and the y-axis has been limited
 234 to 0.5 for clarity. C & E) Correlation between the weekly sliding (C) and daily (E) mean R_t estimates using
 235 reconstructed (y-axis) and reported (x-axis) daily data, excluding the first 30 days due to low incidence. Vertical
 236 and horizontal lines depict the 95% credible intervals and dotted lines show the threshold of $R_t = 1$. D & F)
 237 Correlation between the uncertainty in the weekly sliding (D) and daily (F) R_t estimates, defined as the width
 238 of the 95% credible intervals, using the reconstructed (y-axis) and reported (x-axis) daily data. The colour of
 239 the points in panels C-F correspond to the epidemic phase, i.e. the early (21st March – 12th October 2020),
 240 middle (13th October 2020 – 22nd May 2021) or late (23rd May – 30th December 2021) phase of the data, shown
 241 by the strip in panel A. Solid lines show the linear model fit with 95% confidence intervals (grey shading).
 242 Dashed lines represent the $x = y$ line.

243

244 *COVID-19 deaths*

245

246 The reported incidence of COVID-19 deaths was much less influenced by day-to-day variation. The
 247 reconstructed daily incidence was more similar to the observed daily data than in the previous case studies
 248 (Figure 4A). Both weekly and daily R_t estimates obtained from weekly data were highly consistent with those
 249 obtained from daily observations ($R^2 = 0.98$ and $R^2 = 0.80$ respectively, Figure 4C & 4E). The overall agreement
 250 in R_t classifications for daily estimates was the highest of all case studies at 85.8%, and 93.3% for weekly R_t
 251 estimates (appendix p9). Discrepancies between the two mostly coincide with periods of particularly low
 252 incidence of deaths (Figure 4B & appendix p7). The overall lower incidence of COVID-19 deaths compared to
 253 COVID-19 cases means there is greater uncertainty in R_t estimates in this case study (Figure 4D, 4F & appendix
 254 p6). However, there was minimal difference in the uncertainty of estimates obtained from daily and weekly
 255 data (Figure 4D & 4F).
 256



257

258

259 Figure 4. R_t estimates from daily incidence that was either reported or reconstructed from weekly aggregated
 260 COVID-19 death data. A) The reported (grey) and reconstructed (green) daily incidence of COVID-19 by date
 261 of death within 28 days of a positive test. B) Squared error of the daily (orange) and weekly sliding (pink) R_t
 262 estimates that were made from reconstructed data compared to those obtained from the reported daily data.
 263 R_t estimation starts on the first day of the second aggregation window (day 8 – 9th March 2020) and is plotted
 264 on the last day of the time window used for estimation (i.e., starting on day 9 (10th March) for daily estimates
 265 and day 14 (15th March) for weekly estimates). Note: the x-axis is shared with the incidence plot above and
 266 the y-axis has been limited to 0.5 for clarity. C & E) Correlation between the weekly sliding (C) and daily (E)
 267 mean R_t estimates using reconstructed (y-axis) and reported daily data (x-axis), excluding the first 30 days due
 268 to low incidence. Vertical and horizontal lines depict the 95% credible intervals and dotted lines show the

269 threshold of $R_t = 1$. D & F) Correlation between the uncertainty in the weekly sliding (D) and daily (F) R_t
270 estimates, defined as the width of the 95% credible intervals, using the reconstructed (y-axis) and reported
271 daily (x-axis) data. The colour of the points in panels C-F correspond to the epidemic phase, i.e. the early (31st
272 March - 20th October 2020), middle (21st October 2020 – 28th May 2021) or late (29th May 2021 – 2nd January
273 2022) phase of the data, shown by the strip in panel A. Solid lines show the linear model fit with 95%
274 confidence intervals (grey shading). Dashed lines represent the $x = y$ line.

275
276 In all case-studies, incidence reconstructions converged within 10 iterations of the EM algorithm. The overall
277 process of R_t estimation from weekly aggregated data took three seconds or less to run on MacOS (2 GHz
278 Quad-Core Intel Core i5) 16GB RAM (appendix p10); the influenza scenario, with over 57,000 cases, took two
279 seconds to run, whilst the COVID-19 cases and deaths scenarios, with an overall incidence over 149,000 and
280 13 million cases respectively, took three seconds to run.

281
282 *Simulation study*

283
284 The method performed well across all scenarios, successfully estimating R_t from the aggregated simulated
285 data (appendix pp10-20). Convergence of the EM algorithm was quick, with negligible differences in the
286 reconstructed incidence beyond 5 iterations (appendix p22).

287
288 When introducing weekend effects into simulated data, R_t estimates from reconstructed incidence were more
289 successful at recovering the true value of R_t than when using reported incidence (appendix p21). The method
290 can also be successfully applied to other temporal aggregations of data, e.g. three-, ten- or fourteen-day
291 windows (appendix pp22-24).

292
293 Discussion

294
295 Estimates of the time-varying reproduction number (R_t) have frequently been used to inform and guide
296 policymaking during outbreaks, and a commonly used approach to estimate R_t is EpiEstim, which relies on
297 daily incidence data. However, maintaining daily incidence databases requires substantial time and
298 investment in resources, which is not always feasible, particularly for less acute or routinely reported diseases.
299 Therefore, in practice, many diseases are not reported on a daily basis, including influenza and other notifiable
300 diseases in the UK and US.²⁻⁵ As the COVID-19 pandemic persists, daily reporting is also becoming less
301 common.²¹ Coarsely aggregated data can be challenging to deal with in the context of R_t estimation methods,
302 restricting their applications in certain contexts. In this study, we develop a statistical framework and tool that
303 allows R_t estimation from aggregated incidence without introducing bias. Using influenza and COVID-19 data,
304 alongside a simulation study, we demonstrate how a simple expectation-maximisation algorithm approach
305 can rapidly reconstruct daily incidence data and accurately estimate R_t .

306
307 In all case studies, direct comparisons between weekly sliding R_t estimates show that very similar estimates
308 can be made from the reported daily incidence and the reconstructed daily incidence from weekly aggregated
309 data. However, daily R_t estimates are more influenced by noise, such as intra-weekly variability, leading to
310 greater disparities in estimates between datasets. There are clear weekend effects exhibited in the influenza
311 and COVID-19 case data (appendix p8), leading to peaks and troughs in the reported incidence and the
312 resulting daily R_t estimates (Figures 2 & 3, appendix pp2&4). Using reconstructed incidence considerably
313 smoothed the daily R_t estimates, removing the impact of weekend-effects. The overall agreement in the

314 classification of R_t as increasing, uncertain, or declining between estimates made from each dataset rose
315 substantially when some of the variability in the reported data was smoothed by estimating R_t using weekly
316 sliding windows (appendix pp8-9).

317
318 Despite both being affected by weekly periodicity in reporting, concordance of R_t estimates obtained from
319 COVID-19 case data is considerably better than for influenza, perhaps due to the greater quantity of data, with
320 a very strong positive correlation between daily and weekly R_t estimates (Figure 3). This is reflected in the high
321 overall agreement in the classification of R_t estimates obtained from the reported and reconstructed datasets.
322 It is important to note that outlying and much larger R_t estimates obtained from both datasets coincide with
323 the early phase of the epidemic, when incidence was lower and the prior for R_t ($\mu=5$, $\sigma=5$) had more weight
324 on estimates.

325
326 During the early stages of epidemics, despite there being far fewer deaths than cases, death data can
327 sometimes be considered more reliable.^{22,23} For example, case reporting is affected by surveillance system
328 quality and the robustness of testing practices, which can vary considerably over the course of an epidemic,
329 especially early on. COVID-19 incidence by date of death is much less influenced by administrative noise in the
330 data (appendix p8), and the reconstructed incidence is most similar to the reported daily incidence of any case
331 study. Therefore, the greatest differences in R_t estimates from death data coincide with periods of low
332 incidence (appendix p7) when uncertainty increases. Weekly sliding R_t estimates are equally as correlated as
333 those from COVID-19 case data, but daily R_t estimates are the most strongly correlated of any dataset (Figure
334 4). Additionally, there is very high overall agreement in the classification of daily and weekly R_t (appendix p9).
335 This provides further support that differences between daily R_t estimates for influenza and COVID-19 cases is
336 likely due to the reconstructed incidence smoothing out weekly periodicity in reporting.

337
338 To investigate further, weekend effects were artificially introduced to data in the simulation study (appendix
339 p21). We have shown that, when using reported incidence, R_t estimates are all strongly influenced by weekend
340 effects (regardless of the smoothing time-window). Reconstructing daily incidence from weekly data
341 completely removes the effect of noise from resulting R_t values, greatly improving the accuracy of estimates.
342 This demonstrates that it may be beneficial to artificially aggregate daily data, as has been done in previous
343 studies.^{6,7} However, we did assume quite an extreme level of administrative noise, so in instances where the
344 pattern is less prominent, it may have less of an impact on estimates. Disentangling important temporal trends
345 in R_t from noise in the data can be difficult, and if aggregated data is used it will be at the cost of reduced
346 temporal resolution in R_t estimates.

347
348 This can be seen when the method is applied to data aggregated over longer timescales, such as ten- to
349 fourteen-days (appendix pp22-24). This approach requires two layers of smoothing: 1) the incidence is
350 smoothed over each aggregation window during the reconstruction process and 2) R_t estimates are smoothed
351 by the sliding window chosen by the user. If a change in R_t occurs at the end of an aggregation window (i.e. on
352 the last day), such as a sudden decrease in R_t due to a strict lockdown, that change is detected with a lag,
353 corresponding to the length of the sliding window used for R_t estimation (appendix p23). However, if the event
354 occurs mid-aggregation window, then in addition to the usual lag caused by the sliding window, estimates will
355 be affected by the smoothing of the incidence within the aggregation window during reconstruction (appendix
356 p24). The change in R_t will seem more gradual over the period that data are aggregated over and will appear
357 to start earlier than in reality (corresponding to the first day of the aggregation window). It is important for
358 users to keep this in mind, particularly when using longer aggregations of data.

359

360 Another consideration is that the reconstructed incidence can have discontinuities in the borders between
361 aggregation windows (appendix pp11-12). This occurs because in reconstructing daily incidence we impose
362 that, if it were to be re-aggregated, it would match the original data. Methods that simply fit smoothing splines
363 to weekly data, inferring daily case counts from the daily difference in cumulative counts, are not affected by
364 this.^{24,25} To circumvent this problem, we recommend that sliding windows used to estimate R_t are at least
365 equal to or longer than the length of aggregation windows to reduce the impact of discontinuities on estimates
366 (appendix pp22-24).

367

368 Alternative approaches include modelling frameworks implemented in the Epidemia and EpiNow2 R
369 packages.^{6,22,26} Daily infections are modelled as a latent process, back-calculated from observed data on cases
370 or deaths, depending on an appropriate infection to observation distribution. In addition, Epidemia integrates
371 further information, such as the infection ascertainment rate (for cases) or the infection fatality rate (for
372 deaths).²² This facilitates a ‘nowcasting’ approach, allowing users to estimate R_t directly from the unobserved
373 infections, but they typically require more data (e.g. incidence of deaths and cases), more assumptions (e.g.
374 delay distributions and ascertainment rates), and are much more computationally intensive, which can be a
375 barrier to the adoption of such methods by users.¹⁴

376

377 Here, R_t estimates are based on a single daily incidence reconstruction, meaning R_t can be estimated very
378 rapidly from aggregated data, which is particularly desirable during real-time outbreak analysis.¹⁴ A potential
379 downside is that uncertainty in R_t estimates could be underestimated. However, the simulation study showed
380 that the 95% credible interval of estimates encompassed the correct value of R_t the majority of the time, and
381 we found no substantial indication that this approach detrimentally affected our characterisation of the
382 uncertainty.

383

384 Given that this method is directly derived from EpiEstim, it relies on similar assumptions and caveats.^{15,27} As
385 time of infection is more difficult to observe than symptom onset, the SI is typically used as an approximation
386 of the generation time in the renewal equation, which may introduce bias.²⁸ The SI, the level of undetected
387 cases, and the reporting rate are assumed to remain constant, which is often not the case in practice. Factors
388 such as changes in population immunity, and the introduction of interventions, can alter the SI throughout an
389 epidemic.²⁹ Whilst changing case definitions, new testing practices, and increased healthcare-seeking
390 behaviour, can all affect case ascertainment.¹⁵ Parameters chosen by users can also influence estimation
391 accuracy, for instance, the time window length for temporal smoothing and the prior for R_t .²⁷

392

393 To make the method simple to implement for current and future users of EpiEstim, this extension has been
394 fully integrated with the ‘estimate_R()’ function in the original R package on github.³⁰ Just one additional
395 parameter is required – the number of days data are aggregated over (with some other optional parameters).
396 More details regarding the applications of this method can be found in the package vignette and associated
397 examples.³⁰

398

399 Conclusion

400

401 We extended the widely used R_t estimation approach proposed by Cori et al.,¹⁵ and implemented in the R
402 package EpiEstim, to incorporate a new feature which allows R_t to be easily estimated from any temporal
403 aggregation of incidence data. We have demonstrated that the method performs well using both simulated

404 and real-world data, recovering or even improving upon the estimates that would have been made from
405 reported daily data. This extension is easy to use and computationally efficient, which will enable
406 epidemiologists and other public health professionals to apply EpiEstim to a wider range of diseases and
407 epidemic contexts.

408

409

- 410 1. Baker RE, Mahmud AS, Miller IF, Rajeev M, Rasambainarivo F, Rice BL, et al. Infectious disease in an era
411 of global change. *Nat Rev Microbiol*. 2022 Apr;20(4):193–205.
- 412 2. National flu and COVID-19 surveillance reports: 2021 to 2022 season [Internet]. GOV.UK. [cited 2022 Jun
413 27]. Available from: [https://www.gov.uk/government/statistics/national-flu-and-covid-19-surveillance-](https://www.gov.uk/government/statistics/national-flu-and-covid-19-surveillance-reports-2021-to-2022-season)
414 [reports-2021-to-2022-season](https://www.gov.uk/government/statistics/national-flu-and-covid-19-surveillance-reports-2021-to-2022-season)
- 415 3. Pacheco O, Beltrán M, Nelson CA, Valencia D, Tolosa N, Farr SL, et al. Zika Virus Disease in Colombia —
416 Preliminary Report. *New England Journal of Medicine*. 2020 Aug 6;383(6):e44.
- 417 4. Notifiable diseases: weekly reports for 2022 [Internet]. GOV.UK. [cited 2022 Jun 27]. Available from:
418 <https://www.gov.uk/government/publications/notifiable-diseases-weekly-reports-for-2022>
- 419 5. Notifiable Infectious Disease Tables | CDC [Internet]. 2021 [cited 2022 Jul 2]. Available from:
420 <https://www.cdc.gov/nndss/data-statistics/infectious-tables/index.html>
- 421 6. Mishra S, Scott J, Zhu H, Ferguson NM, Bhatt S, Flaxman S, et al. A COVID-19 Model for Local Authorities
422 of the United Kingdom [Internet]. medRxiv; 2020 [cited 2022 Jul 1]. p. 2020.11.24.20236661. Available
423 from: <https://www.medrxiv.org/content/10.1101/2020.11.24.20236661v1>
- 424 7. Role of Data Aggregation in Biosurveillance Detection Strategies with Applications from ESSENCE
425 [Internet]. [cited 2022 Jul 2]. Available from:
426 <https://www.cdc.gov/mmwr/preview/mmwrhtml/su5301a16.htm>
- 427 8. Ferguson NM, Cucunubá ZM, Dorigatti I, Nedjati-Gilani GL, Donnelly CA, Basáñez MG, et al. Countering
428 the zika epidemic in latin america. *Science*. 2016;353(6297):353–4.
- 429 9. Charniga K, Cucunubá ZM, Mercado M, Prieto F, Ospina M, Nouvellet P, et al. Spatial and temporal
430 invasion dynamics of the 2014–2017 Zika and chikungunya epidemics in Colombia. *PLOS Computational*
431 *Biology*. 2021 Jul 2;17(7):e1009174.
- 432 10. Cowling BJ, Fang VJ, Riley S, Peiris JSM, Leung GM. Estimation of the serial interval of influenza.
433 *Epidemiology*. 2009 May;20(3):344–7.
- 434 11. White LF, Wallinga J, Finelli L, Reed C, Riley S, Lipsitch M, et al. Estimation of the reproductive number
435 and the serial interval in early phase of the 2009 influenza A/H1N1 pandemic in the USA. *Influenza and*
436 *Other Respiratory Viruses*. 2009;3(6):267–76.
- 437 12. Bi Q, Wu Y, Mei S, Ye C, Zou X, Zhang Z, et al. Epidemiology and transmission of COVID-19 in 391 cases
438 and 1286 of their close contacts in Shenzhen, China: a retrospective cohort study. *The Lancet Infectious*
439 *Diseases*. 2020 Aug 1;20(8):911–9.
- 440 13. Rai B, Shukla A, Dwivedi LK. Estimates of serial interval for COVID-19: A systematic review and meta-
441 analysis. *Clin Epidemiol Glob Health*. 2021;9:157–61.

- 442 14. Nash RK, Nouvellet P, Cori A. Real-time estimation of the epidemic reproduction number: Scoping
443 review of the applications and challenges. *PLOS Digital Health*. 2022 Jun 27;1(6):e0000052.
- 444 15. Cori A, Ferguson NM, Fraser C, Cauchemez S. A New Framework and Software to Estimate Time-Varying
445 Reproduction Numbers During Epidemics. *Am J Epidemiol*. 2013 Nov 1;178(9):1505–12.
- 446 16. Cori [aut A, cre, Cauchemez S, Ferguson NM, Fraser C, Dahlgvist E, et al. EpiEstim: Estimate Time
447 Varying Reproduction Numbers from Epidemic Curves [Internet]. 2021 [cited 2022 Jun 9]. Available
448 from: <https://CRAN.R-project.org/package=EpiEstim>
- 449 17. Wallinga J, Lipsitch M. How generation intervals shape the relationship between growth rates and
450 reproductive numbers. *Proceedings of the Royal Society B: Biological Sciences*. 2007 Feb
451 22;274(1609):599–604.
- 452 18. Riley P, Cost AA, Riley S. Intra-Weekly Variations of Influenza-Like Illness in Military Populations. *Military
453 Medicine*. 2016 Apr 1;181(4):364–8.
- 454 19. Cases in the UK | Coronavirus in the UK [Internet]. [cited 2022 Jan 9]. Available from:
455 <https://coronavirus.data.gov.uk/details/cases>
- 456 20. Jombart T, Nouvellet P, Bhatia S, Kamvar ZN, Taylor T, Ghazzi S. projections: Project Future Case
457 Incidence [Internet]. 2021 [cited 2022 Jun 9]. Available from: [https://CRAN.R-
458 project.org/package=projections](https://CRAN.R-project.org/package=projections)
- 459 21. CSSEGISandData. COVID-19 Data Repository by the Center for Systems Science and Engineering (CSSE) at
460 Johns Hopkins University [Internet]. 2022 [cited 2022 Jul 5]. Available from:
461 <https://github.com/CSSEGISandData/COVID-19>
- 462 22. Flaxman S, Mishra S, Gandy A, Unwin HJT, Mellan TA, Coupland H, et al. Estimating the effects of non-
463 pharmaceutical interventions on COVID-19 in Europe. *Nature*. 2020;584(7820):257–61.
- 464 23. Nouvellet P, Bhatia S, Cori A, Ainslie KEC, Baguelin M, Bhatt S, et al. Reduction in mobility and COVID-19
465 transmission. *Nat Commun*. 2021 Feb 17;12(1):1090.
- 466 24. Yamauchi T, Takeuchi S, Yamano Y, Kuroda Y, Nakadate T. Estimation of the effective reproduction
467 number of influenza based on weekly reports in Miyazaki Prefecture. *Scientific reports*. 2019;9(1):1–9.
- 468 25. Nishiura H, Chowell G. Early transmission dynamics of Ebola virus disease (EVD), West Africa, March to
469 August 2014. *Eurosurveillance*. 2014 Sep 11;19(36):20894.
- 470 26. Abbott S, Hellewell J, Thompson RN, Sherratt K, Gibbs HP, Bosse NI, et al. Estimating the time-varying
471 reproduction number of SARS-CoV-2 using national and subnational case counts. *Wellcome Open Res*.
472 2020 Jun 1;5:112.
- 473 27. Gostic KM, McGough L, Baskerville EB, Abbott S, Joshi K, Tedijanto C, et al. Practical considerations for
474 measuring the effective reproductive number, Rt. *PLOS Computational Biology*. 2020 Dec
475 10;16(12):e1008409.
- 476 28. Britton T, Scalia Tomba G. Estimation in emerging epidemics: biases and remedies. *Journal of The Royal
477 Society Interface*. 2019 Jan 31;16(150):20180670.
- 478 29. Ali ST, Wang L, Lau EHY, Xu XK, Du Z, Wu Y, et al. Serial interval of SARS-CoV-2 was shortened over time
479 by nonpharmaceutical interventions. *Science*. 2020 Aug 28;369(6507):1106–9.

- 480 30. mrc-ide/EpiEstim: A tool to estimate time varying instantaneous reproduction number during epidemics
481 [Internet]. [cited 2022 Aug 9]. Available from: <https://github.com/mrc-ide/EpiEstim>
482