

Why are different estimates of the effective reproductive number so different? A case study on COVID-19 in Germany

Elisabeth K. Brockhaus¹, Daniel Wolfram^{1, 2}, Tanja Stadler³, Michael Osthege^{4, 5}, Tanmay Mitra⁶, Jonas M. Littek¹, Ekaterina Krymova⁷, Anna J. Klesen¹, Jana S. Huisman^{3, 8}, Stefan Heyder⁹, Laura M. Helleckes^{4, 5}, Matthias an der Heiden¹⁰, Sebastian Funk^{11, 12}, Sam Abbott^{11, 12}, and Johannes Bracher^{1, 2}

¹Chair of Statistical Methods and Econometrics, Karlsruhe Institute of Technology (KIT), Karlsruhe, Germany

²Computational Statistics Group, Heidelberg Institute for Theoretical Studies (HITS), Heidelberg, Germany

³Department of Biosystems Science and Engineering, ETH Zurich, Zurich, Switzerland

⁴Institute of Biotechnology: IBG-1, Forschungszentrum Jülich GmbH, Jülich, Germany

⁵Institute of Biotechnology, RWTH Aachen University, Aachen, Germany

⁶Department of Systems Immunology and Braunschweig Integrated Centre of Systems Biology (BRICS), Helmholtz Centre for Infection Research, Braunschweig, Germany

⁷Swiss Data Science Center, EPF Lausanne and ETH Zurich, Zurich, Switzerland

⁸Physics of Living Systems, Department of Physics, Massachusetts Institute of Technology, Cambridge, MA, USA

⁹Institute of Mathematics, Technische Universität Ilmenau, Ilmenau, Germany

¹⁰Robert Koch Institute, Berlin, Germany

¹¹Department of Infectious Disease Epidemiology, London School of Hygiene & Tropical Medicine, London, UK

¹²Centre for Mathematical Modelling of Infectious Diseases, London School of Hygiene & Tropical Medicine, London, UK

April 27, 2023

Abstract

The effective reproductive number R_t has taken a central role in the scientific, political, and public discussion during the COVID-19 pandemic, with numerous real-time estimates of this quantity routinely published. Disagreement between estimates can be substantial and may lead to confusion among decision-makers and the general public. In this work, we compare different estimates of the national-level effective reproductive number of COVID-19 in Germany in 2020 and 2021. We consider the agreement between estimates from the same method but published at different time points (within-method agreement) as well as retrospective agreement across different approaches (between-method agreement). Concerning the former, estimates from some methods are very stable over time and hardly subject to revisions, while others display considerable fluctuations. To evaluate between-method agreement, we reproduce the estimates generated by different groups using a variety of statistical approaches, standardizing analytical choices to assess how they contribute to the observed disagreement. These analytical choices include the data source, data pre-processing, assumed generation time distribution, statistical tuning parameters, and various delay distributions. We find that in practice, these auxiliary choices in the estimation of R_t may affect results at least as strongly as the selection of the statistical approach. They should thus be communicated transparently along with the estimates.

* Correspondence to: E. K. Brockhaus (elisabeth.brockhaus@student.kit.edu), J.Bracher (johannes.bracher@kit.edu)

1 Introduction

The definition of the effective reproductive number R_t as the “the expected number of new infections caused by an infectious individual in a population where some individuals may no longer be susceptible” (Gostic et al., 2020) has become widely known even outside of the scientific community during the COVID-19 pandemic. Values above 1 imply epidemic growth, while values below 1 correspond to a decline. Public health agencies and academic groups from around the world have been publishing R_t values in a daily

45 rhythm since the beginning of the pandemic. In the political debate on the tightening or loosening of
 46 intervention measures, these numbers have been routinely cited. Likewise, numerous scientific works on the
 47 efficacy of control measures have attempted to link the development of R_t to specific policy choices (e.g.,
 48 [Haug et al. 2020](#); [Brauner et al. 2021](#); [Knock et al. 2021](#)).

49 A major difference between R_t and other epidemiological indicators is that it is not directly observable in
 50 practice. While numbers of confirmed cases or occupied hospital beds come with their own problems, they
 51 are *data*, i.e., observed values. The effective reproductive number, on the other hand, requires *estimation*
 52 unless the complete transmission chain is observed, which is unrealistic in most settings. Estimation is based
 53 on statistical models which combine data and epidemiological assumptions, leading to a considerable number
 54 of analytical choices to be made. Usually, various defensible options exist, which will influence the results.
 55 Estimates produced by different groups of researchers can therefore differ, as is illustrated in Figure 1. The
 56 top panel shows estimates of the effective reproductive number of COVID-19 in Germany from January 1,
 57 2021, to June 10, 2021, as published by eight different research teams on July 10, 2021. When taken at face
 58 value, these numbers often imply disagreement even on whether R_t was above or below 1. The widths of
 59 95% uncertainty intervals, shown in the bottom panel, vary considerably, and for some pairs of methods,
 60 they hardly overlap. In this article, we are concerned with how these discrepancies come about and how
 61 they are shaped by different analytical choices.

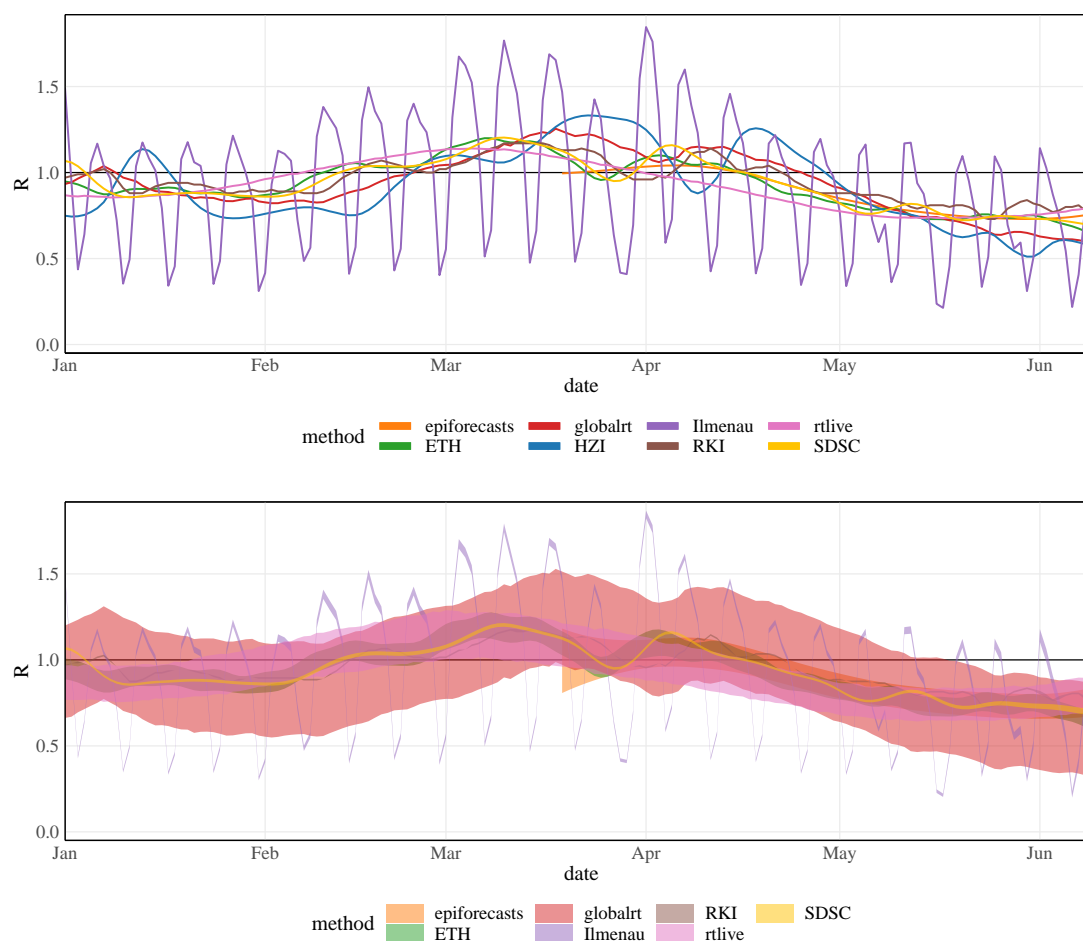


Figure 1: Estimates for the effective reproductive number of COVID-19 in Germany published by eight different research teams on July 10, 2021 (July 11, 2021, for HZI). Top: point estimates (only available for the last 15 weeks for epiforecasts); bottom: 95% uncertainty intervals (not available for HZI).

62 The pronounced differences between estimates of the effective reproductive number have been pointed
 63 out recently by [Wagenmakers et al. \(2022\)](#). In an illustration of different estimates of R_t in the United

64 Kingdom from early October 2020, they found the variability between different estimates to exceed the
65 width of the respective uncertainty intervals. Occasionally, the disagreement between estimates has also
66 spurred confusion in the public debate. For example, on October 27, 2020, Bavarian governor Markus Söder
67 cited an effective reproductive number of 0.57 for his state, which led representatives of the parliamentary
68 opposition to demand a loosening of restrictions ([RedaktionsNetzwerk Deutschland, 2020](#)). This number,
69 however, differed substantially from the value of 0.9 reported for Bavaria on the same day by Robert Koch
70 Institute ([RKI 2020a](#)), the German federal public health agency. As clarified subsequently by the Bavarian
71 State Office for Health and Food Safety, Söder had cited an estimate from Helmholtz Centre for Infection
72 Research (HZI, [Khailaie and Mitra et al. 2020](#)). The statement detailed that the Bavarian authorities
73 monitored estimates from RKI and HZI in parallel, but did not always state the respective source in public
74 communications. The situation is further complicated by the fact that estimates referring to the same day
75 and based on the same method often evolve over time, which has likewise been subject to public debate.
76 As an example, in Fall 2020 it was pointed out that the estimates by RKI were often corrected upwards
77 retrospectively ([Lauck, 2020](#)).

78 Given these challenges, a systematic comparative evaluation of R_t estimates is desirable. This, however, is
79 hampered by several conceptual difficulties. Firstly, there is leeway in the technical definition of the effective
80 reproductive number ([Funk et al., 2022](#)), and different approaches may not actually refer to the exact
81 same estimand. Secondly, the effective reproductive number remains a latent quantity even in hindsight.
82 Systematic comparison of estimates and true values is thus only feasible on synthetic data (e.g., [Gostic et al.](#)
83 [2020](#), [O’Driscoll et al. \(2021\)](#)). Simulation results, however, will necessarily depend on which model is used
84 to generate data, and it is unclear to what degree they translate to the real world. It has been argued
85 that R_t estimates can be evaluated based on derived short-term forecasts ([Teh et al., 2022](#)); this, however,
86 is challenging as e.g., errors in the estimated R_t and the assumed generation time distribution may cancel
87 out so that even bad R_t estimates can yield acceptable forecasts. In this work, we take a complementary
88 approach to simulation and forecasting studies by describing discrepancies between real-world R_t estimates
89 and relating them to underlying analytical choices. A somewhat similar approach has previously been taken
90 by [Pasetto et al. \(2021\)](#), who compared R_t estimates based on an SEIR model and the method by [Cori et al.](#)
91 ([2013](#)). We will analyze R_t estimates for COVID-19 in Germany to study the following aspects:

- 92 • *Within-method temporal coherence:* We assess to which degree estimates based on the same method
93 and referring to the same date, but published at different times, vary. In particular, we analyze the
94 agreement of consolidated point estimates with the uncertainty intervals published near-real-time.
- 95 • *Between-method agreement of retrospective estimates:* We retrospectively compare estimates across
96 different estimation methods. Reproducing the results published by different groups and harmonizing
97 analytical decisions, we gain insights into how they contribute to the observed discrepancies.

98 With our analysis we intend to enable readers to critically assess published R_t estimates and make informed
99 decisions when implementing an estimation scheme themselves. The remainder of the paper is structured
100 as follows. Section 2 provides an overview of the different choices an analyst needs to make to estimate
101 R_t . Moreover, different estimation approaches applied in real-time to German surveillance data during the
102 COVID-19 pandemic are described. In Section 3 we explore within-method temporal coherence, before
103 turning to the between-method agreement in Section 4. Section 5 concludes with a discussion.

104 2 Estimating R_t : The agony of choice

105 Estimating R_t requires numerous decisions by the analyst, ranging from the definition of R_t and the statistical
106 approach to epidemiological parameterizations and the choice of the data set (see also [Vegvari et al. 2022](#)).
107 In this section, we review these dimensions and contrast the decisions underlying various routinely published
108 estimates of R_t of COVID-19 in Germany. Table 1 provides an overview of the research groups whose
109 estimates we consider. Most systems were launched throughout the year 2020 (starting with epiforecasts
110 in early March), and some have in the meantime been retired. Table 2 provides an abridged summary of
111 the model characteristics. For all methods, estimates (and in most cases, analysis codes) were shared in
112 machine-readable format and under open licences in dedicated repositories, see Supplement S1.

Table 1: Overview of the groups who regularly published R_t estimates for Germany during the COVID-19 pandemic. Descriptions of the respective methodology are provided in Section 2.2. Note that web domains provided in footnotes may be discontinued at some point; the links to the repositories provided in Supplement S1 are likely to be more stable.

Institution	Abbrev.	Reference	Active period
ETH Zurich ¹	ETH	Huisman et al. (2022a)	since 2020-05-01
Robert-Koch Institute ²	RKI	an der Heiden and Hamouda (2020)	since 2020-06-04
Technische Universität Ilmenau ³	Ilmenau	Hotz et al. (2020)	since 2020-04-22
Swiss Data Science Center and Institut de Santé Globale, Université de Genève ⁴	SDSC	Krymova et al. (2022)	since 2020-10-01
epiforecasts group / LSHTM ⁵	epiforecasts	Abbott et al. (2020b)	2020-03-02 - 2022-03-31
Forschungszentrum Jülich ⁶	rtlive	System et al. (2020)	2020-09-24 - 2021-07-31
globalrt ⁷	globalrt	Arroyo-Marioli et al. (2021)	2021-02-15 - 2023-01-06
Helmholtz Centre for Infection Research ⁸	HZI	Khailaie and Mitra et al. (2021)	since 2020-04-29

2.1 Definition of R_t

There are at least two ways of formalizing the concept of the (time-varying) effective reproductive number (Gostic et al., 2020). The *case reproductive number*, R_t^{case} , quantifies how many new infections individuals who became infected at time t will cause on average. It is thus forward-looking and compares these individuals to the *following* generation of infected. The *instantaneous reproductive number*, R_t^{inst} , on the other hand, is backward-looking and compares them to the *previous* generation. Specifically, it is given by the expected number of infections occurring at t , divided by the number of previously infected individuals, each weighted by their relative infectiousness at time t . A simple discrete-time display of the recursive relationship between infections X_t occurring on days $t = 1, 2, \dots$ can help to understand this distinction (White et al., 2021). For the instantaneous reproductive number, the recursion, also called the *renewal equation*, is given by

$$\mathbb{E}(X_t \mid X_{t-1}, \dots, X_1) = R_t^{\text{inst}} \times \sum_{i=1}^{t-1} w_i X_{t-i}, \quad (1)$$

where w_i is the probability that the generation time (i.e., the time between primary and secondary infection) equals i time units. Here, the index t in R_t refers to the time of secondary infection. For the case reproductive number the recursion is

$$\mathbb{E}(X_t \mid X_{t-1}, \dots, X_1) = \sum_{i=1}^{t-1} R_{t-i}^{\text{case}} w_i X_{t-i}, \quad (2)$$

the index $t - i$ in R_{t-i} thus referring to the time of primary infection. We note that R_{t-i}^{case} can be seen as a convolution of R_t^{inst} and the generation time distribution. In the absence of sudden changes, shifting R_t^{case} back by the mean generation interval m usually leads to good agreement with R_t^{inst} (i.e., R_{t-m}^{case} and R_t^{inst} can be expected to be similar; Gostic et al. 2020).

2.2 Modelling and estimation approaches

Numerous statistical approaches exist to estimate R_t from data. We do not provide a comprehensive review, but focus on methods various research teams have employed in real time to estimate R_t of COVID-19 in Germany (see Table 1). Descriptions are kept concise and we point to the respective references for details.

¹<https://ibz-shiny.ethz.ch/covid-19-re-international/>

²https://www.rki.de/DE/Content/InfAZ/N/Neuartiges_Coronavirus/Situationsberichte/COVID-19-Trends/COVID-19-Trends.html

³<https://stochastik-tu-ilmenau.github.io/COVID-19/germany>

⁴<https://renkulab.shinyapps.io/COVID-19-Epidemic-Forecasting/>

⁵<https://epiforecasts.io/covid/>, previously at <https://cmmid.github.io/topics/covid19/global-time-varying-transmission.html>

⁶<https://rtlive.de>

⁷<http://www.globalrt.live/>

⁸<https://gitlab.com/simm/covid19/secir/-/wikis/Report>

134 **Variations and extensions of the Cori method.** Four groups made use of the method by [Cori et al.](#)
135 [\(2013\)](#), but with different parameterizations and data pre-processing. This method uses formulation
136 [\(1\)](#) combined with a Poisson distribution for new cases. Estimation of R_t is then carried out for sliding
137 windows of a width chosen by the analyst. In the widely used R package *EpiEstim* ([Cori et al., 2020](#))
138 inference is based on a Bayesian approach.

139 **RKI** ([an der Heiden and Hamouda, 2020](#); [RKI, 2020b](#)). First, sampling-based nowcasting
140 is applied in order to impute missing symptom onset dates in incidence data and to correct recent
141 values for reporting delays. Next, the method by [Cori et al. \(2013\)](#) is applied to each sampled time
142 series, using a fixed generation time and frequentist inference. Uncertainty intervals result from
143 the spread of the R_t estimates across different nowcasting samples. The estimation of uncertainty
144 from the Cori method is not taken into account. The window size is set to either 4 or 7 days. We
145 focus on the latter which has been used more widely.

146 **ETH** ([Huisman et al., 2022a](#)). Local polynomial regression (LOESS) is applied to the time
147 series of reported cases to account for weekday effects. The smoothed time series is deconvoluted
148 using various types of delay distribution to reconstruct the time series of infections. R_t is
149 then estimated using the *EpiEstim* package and a window size of 3 days. Uncertainty intervals are
150 obtained by combining the credible intervals and a block bootstrapping approach. The bootstrapping
151 step was only added on January 26, 2021, and led to a widening of intervals (leaving point
152 estimates unaffected). The ETH team published four estimates in parallel (based on confirmed
153 cases as used here, as well as on new hospitalizations, death, and test positivity percentages). We
154 focus on the R_t estimates referred to as “sliding window” (the default in the ETH dashboard).

155 **SDSC** ([Krymova et al., 2022](#)). The case time series is smoothed via a LOESS-based seasonal-
156 trend decomposition prior to estimation using the *EpiEstim* package. The window size is set to 4
157 days. The proposed extension is focused on the point estimates from the Cori method and does
158 not involve the computation of uncertainty intervals. The provided intervals correspond to those
159 returned by the *EpiEstim* package.

160 **Ilmenau** ([Hotz et al., 2020](#)). The effective reproductive number is estimated in a frequentist
161 fashion using equation [\(1\)](#) and a window size of one day. Wald-type confidence intervals are based
162 on newly derived asymptotic standard errors of the employed estimator.

163 **epiforecasts** ([Abbott et al., 2020b](#)). The estimation of R_t is based on a Bayesian latent variable ap-
164 proach, implemented in the R package *EpiNow2* ([Abbott et al., 2020a](#)). The infection dynamics are
165 modeled as in equation [\(1\)](#) and linked to the observed case time series via convolutions with the
166 assumed incubation time and reporting delay distributions. The observation model is given by a
167 negative-binomial distribution. A zero-mean Gaussian process with a Matérn kernel is used for the
168 first-order temporal differences of the effective reproductive number with the magnitude and length-
169 scale estimated jointly with other parameters. Like for ETH, estimates based on hospitalizations and
170 deaths were available, too, but we focus on estimates based on case incidences.

171 **rtlive** ([Systrom et al., 2020](#); [Osthege et al., 2021](#)) Estimates are based on relationship [\(2\)](#), which is
172 combined with a delay process from infection to detection and a re-scaling of case numbers with inverse
173 testing volumes. Inference is conducted in a Bayesian fashion. Similarly to the epiforecasts approach,
174 a negative binomial observation model is used and R_t is assigned a random walk prior.

175 **globalrt** ([Arroyo-Marioli et al., 2021](#)). This approach exploits a relationship between the epidemic
176 growth rate and the effective reproductive number which holds under the SIR (susceptible-infected-
177 removed) model. The effective reproductive number is assumed to follow a random walk and estimation
178 from observed growth rates is done via a Kalman filter or smoother. We here focus on the smooth-
179 ing version, which corresponds to a case reproductive number, as this was displayed in the public
180 dashboard. The generation time distribution is assumed to be exponential as in the SIR model.

181 **HZI** ([Khailaie and Mitra et al., 2021](#); [Knabl, Mitra and Kimpel et al., 2021](#)) A deterministic SE-
182 CIR (susceptible - exposed - carrier - infected - recovered) model with time-varying parameters is fitted
183 to cumulative case and death numbers, with certain parameters fixed to or varied around literature
184 estimates. Estimates of R_t are computed from the model parameters, which are estimated for sliding

185 10-day windows. We use estimates which in addition were smoothed using a 7-day moving average, as
 186 shown in the HZI dashboard.

187 2.3 Epidemiological assumptions and parameterization

188 All described approaches require some parameterization, i.e., specification of epidemiological assumptions.
 189 In particular, the distributions of the following durations and delays need to be chosen.

- 190 • The *generation time* (GT), i.e., time between primary and secondary infection. The impact of the
 191 chosen generation time distribution on R_t estimates is well-studied (Wallinga and Lipsitch, 2007). The
 192 longer the assumed mean generation time, the greater the amplitude of estimates away from 1 (i.e.,
 193 estimates are increased if $R_t > 1$ and decreased if $\hat{R}_t < 1$ for a prolonged period of time). The variance
 194 of the GT has a more subtle effect. If R_t is time-constant, R_t estimates are further from 1 the smaller
 195 the variance (Wallinga and Lipsitch, 2007); for time-varying R_t , the assumed variance also influences
 196 the smoothness of the estimated trajectories.
- 197 • The *incubation period* (IP), i.e., time from infection to symptom onset.
- 198 • The *reporting delay* (RD) between symptom onset and reporting. Changing the mean incubation time
 199 and reporting delay shifts R_t estimates in time. The impact of the variance is not well-studied, but it
 200 likely affects the smoothness of estimates.

201 Table 2 summarizes the distributions used by the different groups. The means and standard deviations
 202 of the generation time distribution are moreover displayed in Figure 2. To illustrate that the variability in
 203 assumed values is not limited to the German context we added values used by various European public health
 204 agencies (see Supplement S2 for sources). Note that the values for HZI are not explicitly provided in the
 205 manuscript by Khailaie and Mitra et al. (2021), but have been computed by us based on model parameters
 206 reported there (see Supplement S3.1). The globalrt dashboard allowed users to select a mean generation
 207 time between five and ten days; we here use the default setting of seven days.

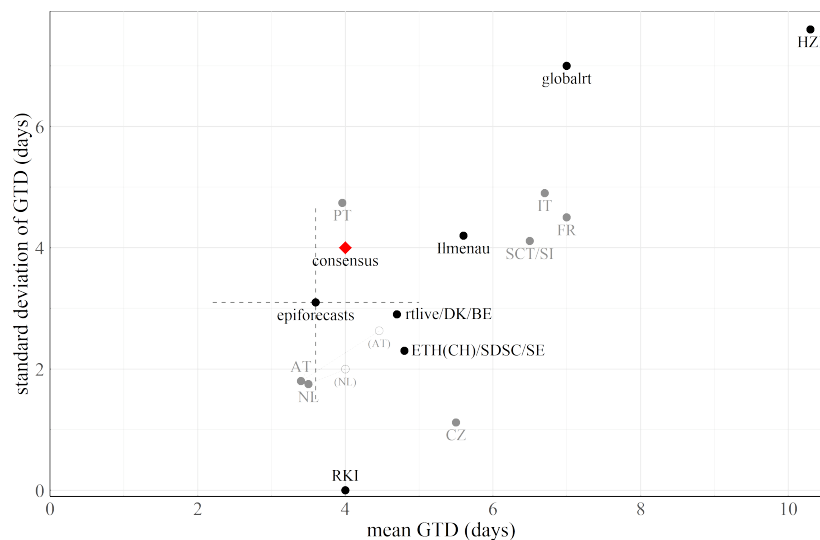


Figure 2: Scatter plot of mean generation time and corresponding standard deviation used by different research groups. The red rhombus represents a “consensus value” chosen for further analysis, see Section 4.1. epiforecasts accounted for uncertainty in the generation time distribution by assuming independent normal priors for the mean and standard deviation; we illustrate the respective 95% uncertainty intervals by a cross. For context, we also show values used by public health agencies of other European countries. In the Netherlands (due to the transition to the Omicron variant) and Austria (due to a data update) the parameterization was revised. For details and references see Supplement S2.

Table 2: Methodological characteristics and parameterizations of the compared estimation approaches. The table follows the structure of Sections 2.1–2.5. The *consensus* model is introduced in Section 4.1 By *conditional distribution of X_t* we refer to the distribution of new cases X_t in formulation (1) or (2). The concept of “revision due to smoothing” is discussed in Section 3.3. Abbreviations: GT = generation time; IP = incubation period; RD = reporting delay.

Panel A: Methods based on [Cori et al. \(2013\)](#) and a consensus parameterization used in Section 4.

	ETH	RKI	Ilmenau	SDSC ¹	consensus
type of R_t	instantaneous	instantaneous	instantaneous	instantaneous	instantaneous
underlying epidemic model	Cori et al.	Cori et al.	Cori et al.	Cori et al.	Cori et al.
regularization/prior on R_t	sliding window	sliding window	sliding window	sliding window	sliding window
cond. distr. of X_t	Poisson	Poisson	Poisson	Poisson	Poisson
inference	Bayesian	max. lik.	max. lik.	Bayesian	Bayesian
preprocessing	smooth. + deconv.	nowcast	–	smoothing	–
window size	3	7, 4	1	4	7
rev. due to smoothing	yes	no	no	yes	no
GT distribution type	gamma	constant	ad hoc	gamma	exponential
mean GT (sd)	4.8 (2.3)	4.0	5.6 (4.2)	4.8 (2.3)	4 (4)
source of GT	Nishiura et al.	–	–	Nishiura et al.	Figure 2
mean IP (sd)	5.3 (3.2)	1.0	5.0	–	0
mean RD (sd)	5.5 (3.8)	3.4	2.0	7.0	7
incidence data source	RKI, by onset date	RKI, by onset date	RKI, by test date	JHU	RKI, by test date

Panel B: Other Methods

	epiforecasts	rtlive	globalrt	HZI ¹
type of R_t	instantaneous	case	case	instantaneous
underlying epidemic model	Abbott et al.	Systrom et al.	Arroyo-Marioli et al.	Khailaie & Mitra et al.
regularization/prior on R_t	Gaussian process	random walk	random walk	sliding window
cond. distr. of X_t	negative binomial	negative binomial	Gaussian ²	deterministic
inference	Bayesian	Bayesian	Kalman smoother	literature est., least squares
preprocessing	–	–	–	–
window size	–	–	–	10 & 7
rev. due to smoothing	yes	yes	yes	no
generation time distr.	gamma	log-normal	exponential	mixt./conv. of exponentials
mean GT (sd)	3.6 (3.1) ³	4.7 (2.9)	7 (7)	10.3 (7.6)
source of GT	Ganyani et al.	Nishiura et al.	–	–
mean IP (sd)	5.4 (2.2) ³	5.0	–	5.2
mean RD (sd)	5.9 (14.6) ³	7.1 (5.9)	–	3.7
incidence data source	WHO	RKI, by test date	JHU	RKI, by test date

¹ Some statements were derived for the present study or retrieved from analysis codes rather than the referenced paper; for details on HZI see Appendix S3.

² The globalrt model operates on the scale of daily growth rates rather than incidences but implies a conditional Gaussian distribution for the latter.

³ The epiforecasts team was the only one to account for uncertainty in the GT, IP and RD distributions; see also Figure 2.

208 2.4 Methods-specific tuning parameters and prior distributions

209 The standardized display of analytical choices in Table 2 neglects that in each modeling approach, some
 210 additional decisions arise. Bayesian estimation as employed by several teams requires choosing prior dis-
 211 tributions. The HZI approach takes into account numerous epidemiological characteristics other than the
 212 generation time, which are informed by literature estimates. The SDSC and ETH approaches involve data
 213 smoothing and deconvolution, which require fixing various tuning parameters. These aspects cannot be
 214 standardized across methods, and we refrain from analyzing them in detail. Instead, we pragmatically leave
 215 them at the values specified by the respective teams wherever needed.

2.5 Input data sources

While R_t can also be estimated from death or hospitalization counts (Sherratt et al., 2021; Huisman et al., 2022a), we focus on estimates based on COVID-19 case numbers. In Germany and during the considered time period (April 2020 – July 2021), such data were regularly released by RKI (2020a), the World Health Organization (WHO; 2022), and the Center for Systems Science and Engineering at Johns Hopkins University (JHU; Dong et al. 2020). The WHO and JHU data were aggregated by the time cases first appeared in the respective data set. The RKI data were in a line list format containing a reference date called the *Meldedatum* (“reporting date”) and for a subset of cases the symptom onset date. The *Meldedatum* denoted when a local health authority digitally registered a case and usually corresponded to the date of the positive test. The Ilmenau, HZI, and rtlive groups aggregated the RKI data by this date. RKI and ETH used the date of symptom onset where available. While RKI completed missing onset dates via multiple imputation, ETH used the reporting date when the symptom onset date was not available and adjusted the reporting delay in the deconvolution accordingly. rtlive additionally used (not publicly available) data on testing volumes.

Figure 3 shows the different time series for January through June 2021. The series denoted “RKI, positive test” is aggregated by the date of the positive test using the implementation from rtlive. “RKI, symptom onset” is the time series by symptom onset date as reconstructed by RKI. The series by symptom onset is shifted to the left compared to the others; the WHO data are somewhat shifted to the right, while the JHU and RKI data by test date are largely aligned. All series display within-week seasonality, with a smaller amplitude for the RKI data by onset date. The JHU data occasionally display spikes absent in the other series. A last relevant aspect, going beyond Figure 3, is the temporal stability of the data. While the WHO and JHU data were only rarely subject to revisions, the last 3–5 entries in the RKI case data were typically still updated retrospectively.

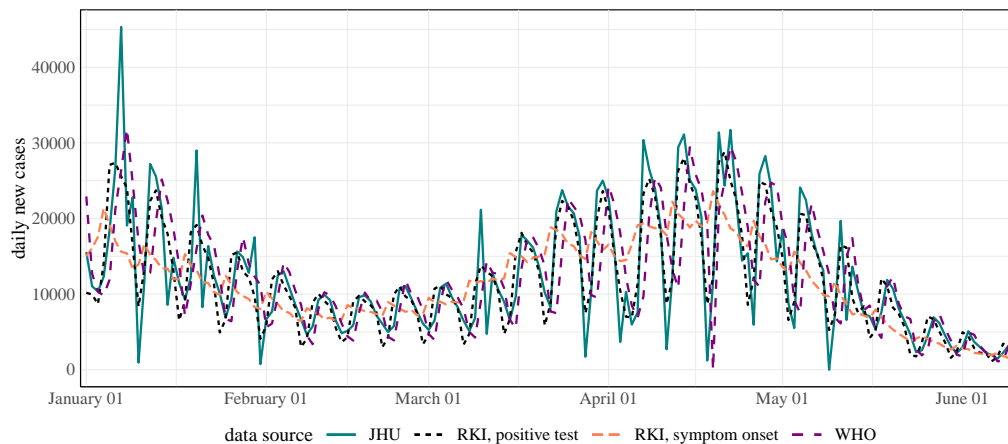


Figure 3: Case incidence time series used by different research groups. To enhance visibility we only display the period January through June 2021 (data version: November 23, 2021).

3 Within-model temporal coherence of real-time estimates

We now move to the analysis of R_t estimates based on the methods and parameterizations described before. Estimates were typically updated each day in an automated fashion. Oftentimes these updates also concerned estimates for the past, which are revised in light of new data. Consequently, for each *target date*, i.e., the date to which an estimate refers, a multitude of estimates issued on different *publication dates* are available. This raises the question of *temporal coherence* of estimates. By this, we mean that estimates issued at various times should not differ more than implied by the respective uncertainty intervals. Temporal coherence is a necessary, though not sufficient, prerequisite for reliable estimation. After all, if subsequent estimates from a method are incompatible, agreement with the underlying truth is necessarily limited. Our analyses are based on real-time estimates obtained from the repositories referenced in Supplement S1. We do not explicitly take into account possible modifications of methods during the considered time period; strictly speaking, we

249 thus assess the coherence of estimation systems, which may evolve over time, rather than uniquely defined
250 methods with fixed parameterizations.

251 3.1 Illustrating the evolution of R_t estimates over time

252 Figure 4 illustrates how real-time R_t estimates from different methods evolved over time. For each method
253 and a 70-day period, it overlays real-time estimates and an estimate made six months later (black line)
254 when all data and results can be expected to have stabilized. Where available, 95% uncertainty intervals¹
255 are shown as shaded areas. We display estimates published on Thursdays where available and published on
256 neighboring days otherwise. Dates of publication are indicated by vertical lines. Note that most teams do
257 not provide estimates up to the publication date, i.e., the R_t trajectories do not reach the vertical line in the
258 respective color. Moreover, some teams (epiforecasts, Ilmenau, SDSC) marked estimates for recent dates as
259 “based on partial data” or “forecast”, which we indicate by dashed and dotted lines, respectively.

260 Some patterns can be discerned in how and how strongly estimates are revised. While the HZI estimates
261 hardly changed, for RKI and Ilmenau recent values tended to be corrected upwards. The ETH estimates,
262 on the other hand, were mostly corrected downwards for the displayed period. rtlive, epiforecasts and (to a
263 lesser degree) globalrt estimates tended to be corrected upwards when R_t was increasing and downwards in
264 periods when R_t was decreasing. For SDSC, there were some pronounced corrections, but without a clear
265 pattern. Moreover, the approaches differed in the width of the uncertainty intervals. While those from SDSC
266 and Ilmenau were very narrow, those of rtlive and globalrt were so wide that they almost always included
267 the threshold value of 1. For most methods, uncertainty increased for recent dates, leading to funnel-shaped
268 bands. This was particularly prominent for epiforecasts, whereas the SDSC intervals were of almost constant
269 width. As mentioned in Section 2.2, the ETH method was revised in early 2021; this change explains why
270 the consolidated intervals are wider than those from Fall 2020. The HZI estimates were published in the
271 form of samples, but it is unclear whether they can be seen as an uncertainty quantification. As estimates
272 were displayed without uncertainty bands on the HZI website, we likewise omit them.

273 3.2 Systematic assessment of temporal coherence

274 To substantiate these observations, we assess the temporal coherence of estimates quantitatively. Unlike in
275 Figure 4 we do not use estimates made at a single later time point as the consolidated ones. Instead, for each
276 target date we use estimates generated 70 days later. This ensures that the time during which the estimates
277 could be revised is the same for all target dates. Based on this definition we computed the following.

- 278 • The fraction of instances in which the 95% uncertainty intervals issued in real-time covered the respec-
279 tive consolidated point estimate.
- 280 • The average width of 95% uncertainty intervals.
- 281 • The mean absolute difference (MAD) between real-time and consolidated point estimates. This reflects
282 the volatility of real-time estimates relative to the consolidated ones.
- 283 • The mean signed difference (MSD) of real-time and consolidated estimates. This reflects if revisions are
284 systematically in one direction. We orient this such that positive values indicate upwards corrections.

285 Figure 5 summarizes the results for estimates published between October 1, 2020, and July 22, 2021. Not all
286 models were operated during the entirety of this period, but we consider it a reasonable overlap (see Table
287 1 and Supplementary Figure S15 on when methods were operated). This period includes two full waves of
288 infections (Figure S16) so that effects caused by rising or falling case numbers should largely cancel out.
289 Results are shown as a function of the number of days between the publication date and the target date. E.g.,
290 “10d back” means that the estimate refers to the date ten days before the time of estimation. We here stuck
291 to the labeling of estimates by the respective research teams. As they assumed different incubation periods
292 and reporting delays (Table 2), estimates from different methods are not necessarily aligned. Notably, the
293 estimates (and thus curves in Figure 5) by epiforecasts, rtlive, and ETH are shifted to the left relative to
294 the others, as longer incubation periods and reporting delays were assumed. In Supplementary Figure S17

¹epiforecasts reported 90% rather than 95% uncertainty intervals, along with a standard deviation. As the 90% intervals agreed well with the Gaussian approximation $\text{mean} \pm 1.645 \times \text{sd}$, we approximated the 95% intervals as $\text{mean} \pm 1.96 \times \text{sd}$.

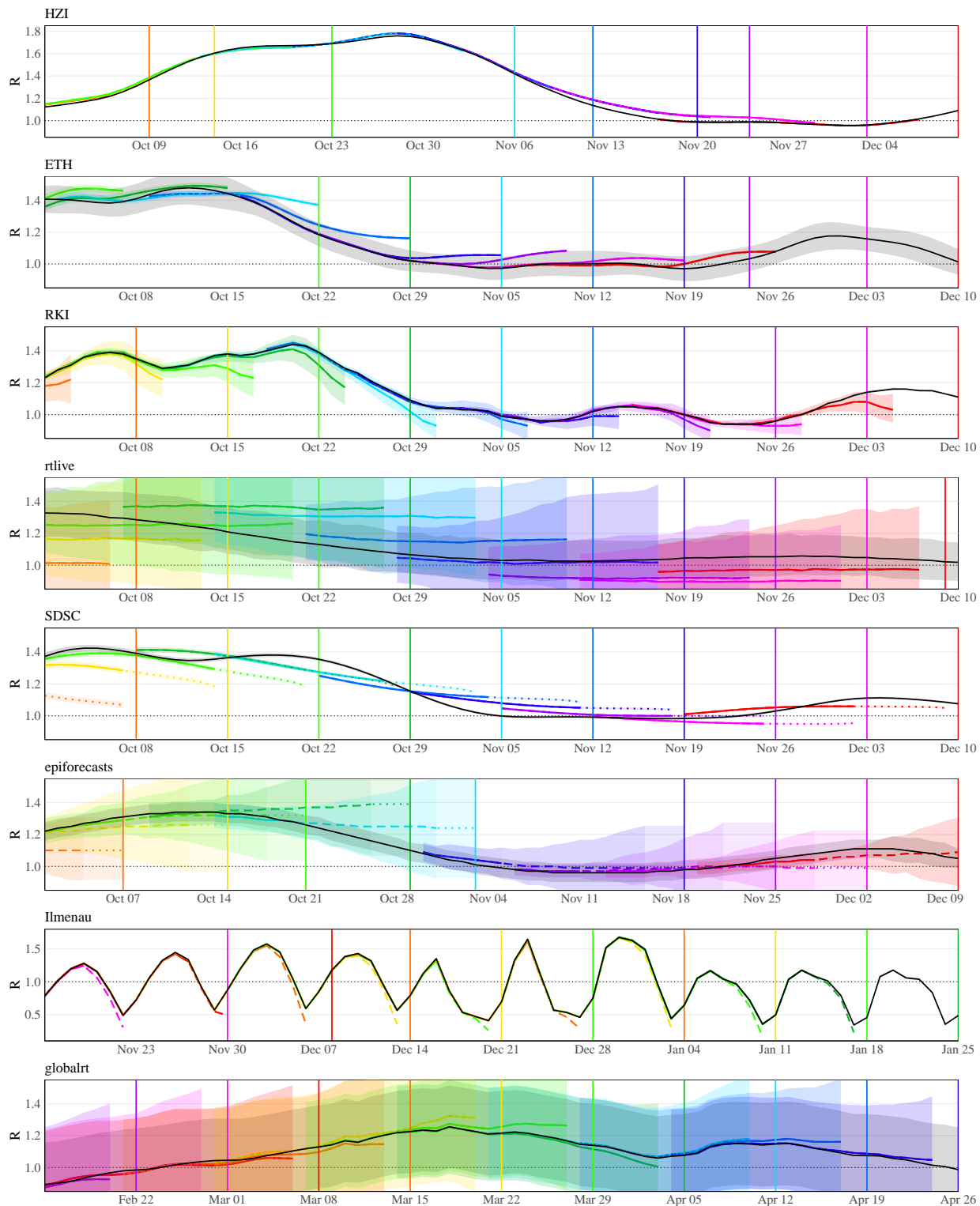


Figure 4: R_t estimates published between October 1, 2020, and December 10, 2020, and a consolidated estimate published 6 months later (epiforecasts: 15 weeks later). Note the different scales of the y-axes for HZI and Ilmenau, and the different time periods for Ilmenau and globalrt (which were not operated during the period shown for the other models). The consolidated ETH intervals are wider than those issued in real time due to a revision of methodology. The line type represents the label assigned to the estimate by the respective team: solid: “estimate”, dashed: “estimate based on partial data”, dotted: “forecast”. Shaded areas show 95% uncertainty intervals.

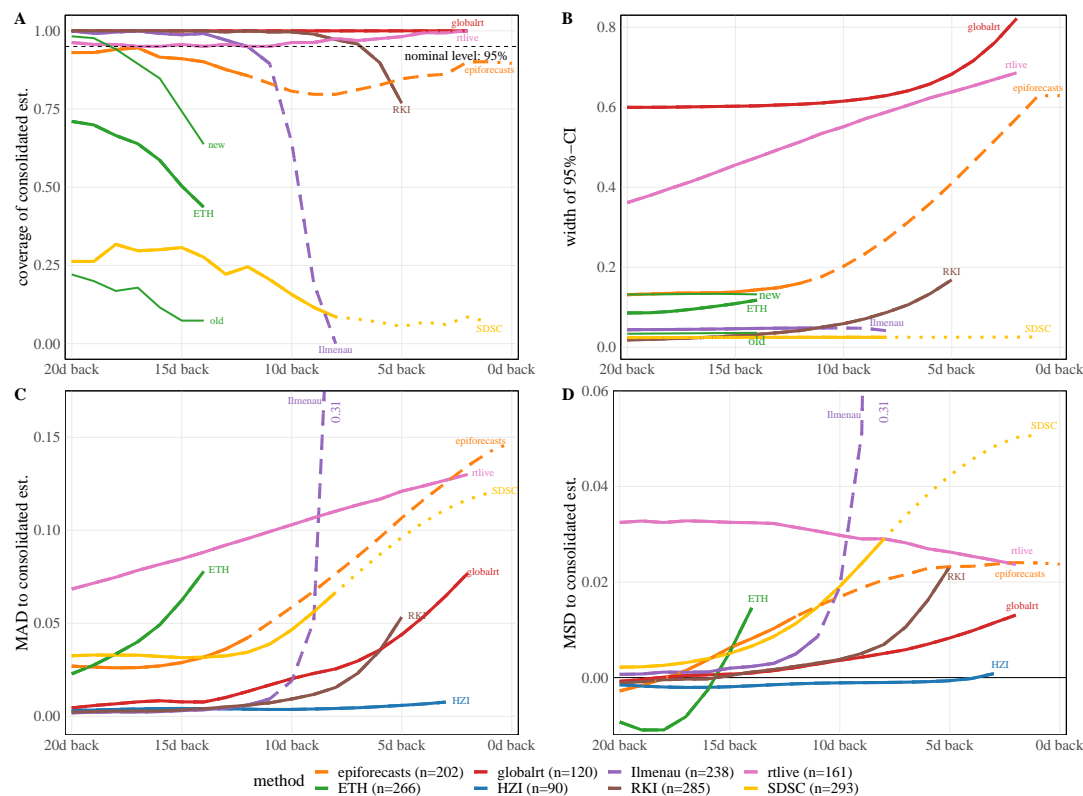


Figure 5: Temporal coherence of R_t estimates. Panels: **A** Proportion of 95%-uncertainty intervals issued in real-time which contained the consolidated estimate. **B** Mean width of 95%-uncertainty intervals. **C** Mean absolute difference of the real-time and consolidated estimates. **D** Same as C, but signed rather than absolute differences. All indicators are shown as a function of the time between the target date (as stated by the teams) and the publication date. Averages refer to the period October 1, 2020 – July 22, 2021 (see Figure S15 for exact periods during which methods were operated). The consolidated estimate corresponds to the one published 70 days after the respective target date. For ETH two additional lines are included in the top row differentiating between intervals obtained from the old procedure before January 26, 2021 ($n = 95$), and from the new bootstrap approach afterward ($n = 171$; see model description in Section 2.2).

we provide a display where curves are aligned to improve comparability. The respective shifts have been determined in a data-driven way, see Section 4.1 and Supplement S4.

Panel A shows the coverage fractions of the 95% uncertainty intervals as defined above. These were in the order of 95% for rtlive and consistently 100% for globalart. For epiforecasts, coverage was close to nominal less than 4 and more than 14 days back, while there was a moderate dip in between (this concerned mostly estimates marked as “based on partial data”). RKI and Ilmenau achieved close to complete coverage for dates further back in the past, starting from 9 and 14 days back, respectively. For more recent values, however, coverage dropped. This was particularly pronounced for Ilmenau, with coverage falling to 0% at 8 days back. ETH overall achieved coverage values of 40% to 75% during the period examined in this paper. As can be seen from the additional lines labeled “old” and “new”, coverage was considerably higher for estimates published after January 26, 2021, when the computation of intervals was revised (with the explicit goal to account for more sources of uncertainty, see Huisman et al. 2022b). The coverage of the SDSC (default *EpiEstim*) intervals was around 25% for values labeled as “observed” and dropped to roughly 10% for values labeled “predicted”. Panel B shows the average width of the 95% uncertainty intervals. The funnel-shaped character of the confidence intervals of globalart, rtlive, epiforecasts, and RKI is reflected in the upward shape of the respective curves. As already visible in Figure 4, the uncertainty intervals issued by globalart and rtlive were considerably wider than those from the other groups. SDSC and Ilmenau issued the most narrow intervals. Prior to the change in methodology in January 2021, the ETH intervals were

313 similarly narrow but became wider afterward.

314 Panels C and D display the mean absolute and mean signed difference between real-time and consolidated
315 estimates, respectively. For all methods, the mean absolute difference was the largest for recent values. A
316 particularly striking picture is seen for the Ilmenau estimates, where the average correction of estimates 8
317 days back was 0.31. For some methods the MAD approached zero after a few days (HZI, RKI, Ilmenau,
318 globalrt), indicating that the estimates stabilized. For the remaining models, the average corrections were
319 clearly non-zero even 20 days back, with epiforecasts and SDSC showing a flat pattern from around 12 days
320 back. Panel D shows that for most methods estimates tended to be corrected upwards, especially recent
321 ones. As already visible from Figure 4, this includes RKI and Ilmenau. For ETH the picture is somewhat
322 difficult to interpret, as the sign of the average correction flips at 16 days back. It should be noted that
323 for most models the mean signed differences were much lower than the mean average differences, indicating
324 that corrections in both directions occurred.

325 To assess the sensitivity of these results to the definition of the consolidated estimates, we compared
326 estimates published 50 and 70 days after the target date. As shown in Supplementary Figure S18, these
327 agree closely. The exact definition of the consolidated estimates is thus not crucial for our results.

328 3.3 Interpretation of observed patterns

329 We now provide some interpretation of the identified patterns, pointing out possible connections to modeling
330 choices. Retrospective revisions of R_t estimates can stem from two main mechanisms. Firstly, past incidence
331 values can be revised in the input data, which will lead the same estimation method to produce different
332 results when re-run. The RKI data were subject to such revisions, which affected in particular the last few
333 days. The JHU and WHO data, on the other hand, were rarely revised. Data revisions were typically upward
334 as delayed reports were added. It seems likely that the strong upward corrections in the Ilmenau estimates
335 stem from this aspect as reporting delays were not accounted for explicitly. The RKI method included a
336 nowcasting step to account for delays, but the correction seems to have been slightly too weak. The rtlive
337 model accounted for revisions by an empirically determined reporting delay distribution. However, it also
338 relied on testing volume data which was more prone to data revision.

339 In the Cori and the HZI methods, the length of the estimation window moderates how strongly results
340 can change due to data revisions. The Ilmenau model, which used a one-day window, was strongly affected
341 as estimation hinged purely on the rather unstable last data point. The HZI model, on the other hand,
342 used a ten-day window for estimation and additionally smoothed the consolidated estimates via a trailing
343 seven-day moving average. The consolidated estimates were thus based on a 16-day window (with some
344 weighting). As the revisions of the RKI data only concerned a small part of this window, the resulting
345 revisions of estimates were negligible. We illustrate this in Supplementary Figure S13, which shows that
346 without the additional smoothing step slightly more pronounced revisions of estimates occurred.

347 The second reason why estimates may change is smoothing during the estimation process. This can
348 enter either via data pre-processing (ETH, SDSC) or model assumptions on the R_t trajectory (Gaussian
349 process assumption in epiforecasts, random walk in rtlive and globalrt). Via smoothing, a new data point
350 can influence how the model treats previous data, and thus impact the results for preceding target dates.
351 We note that smoothing is a planned feature of the approaches in question. Indeed, estimates up to the
352 day of estimation as available from epiforecasts would not be feasible without a generative assumption
353 implying some smoothness. The trade-off is that near-real-time estimates are increasingly extrapolations
354 of the previous R_t trajectory, and likely to change once more data become available. This explains why
355 estimates from epiforecasts, globalrt, and rtlive were often corrected upwards when R_t was on the rise and
356 downwards when it was on the fall. For methods based on trailing estimation windows (RKI, Ilmenau, HZI)
357 revisions cannot arise from this aspect, even though window sizes larger than one also imply some smoothing.

358 Lastly, how well uncertainty intervals cover consolidated estimates depends on how wide the former are.
359 By issuing wide intervals, globalrt and rtlive achieved high coverage despite substantial revisions. While we
360 defer a discussion of the overall interval widths to Section 4, we provide some remarks on the widening of
361 intervals for target dates close to the publication date. This funnel-like pattern was particularly pronounced
362 for epiforecasts, globalrt, and rtlive. These methods provided estimates closest up to the publication date,
363 which as mentioned before, got less and less constrained by data. In the Bayesian framework, this translated
364 naturally to wider uncertainty intervals. In the case of rtlive, this was reinforced by hard-coded assumptions
365 on the variability of the random walk. In the RKI approach, the uncertainty from the nowcasting step was

366 forwarded to the R_t estimation, leading to similarly expanding intervals. For both epiforecasts and RKI, this
367 widening was not quite pronounced enough, however, and interval coverage fell below the minimum desired
368 level of 95%. The Ilmenau, ETH, and SDSC (default *EpiEstim*) approaches showed little to no widening of
369 intervals. The likely reason is that the uncertainty about the recent data points was not forwarded to the
370 R_t estimation from earlier preprocessing steps (see the discussion section of [Hotz et al. 2020](#) on additional
371 sources of uncertainty). In all three cases, this led to a drop in 95% interval coverage below 50%.

372 4 Between-method agreement of retrospective estimates

373 We now turn to the agreement across estimates by different research groups, which as shown in Figure 1 can
374 differ substantially. Our approach is to standardize analytical choices in order to assess their contribution
375 to the overall disagreement. This is inspired by the *vibration of effects* framework ([Klau et al., 2020](#)), which
376 for observational studies serves to assess the sensitivity of effect estimates to aspects like model choice and
377 measurement errors. While e.g., the impact of the assumed generation time distribution on estimates is well-
378 understood at a theoretical level (see Section 2.3), we aim to answer an empirical question: What differences
379 arise in practice when different researchers independently take the necessary analytical decisions?

380 4.1 Sequential standardization and individual variation of analytical choices

381 As visible from Table 2, the available R_t estimates are not only the results of different statistical methods
382 but also of different parameterizations and input data. Isolating the contributions of these aspects requires
383 standardizing the remaining dimensions as far as possible. In what follows we describe a “consensus setting”
384 which we implement for each of the represented methods (see also Panel A of Table 2, last column).

- 385 • *Incidence data:* We use RKI data, which are the most common choice among teams. We use data by
386 test date as aggregated by rtlive for all methods requiring a simple time series. Models making use of
387 information on symptom onset dates (RKI, ETH) or test positivity percentages (rtlive) can keep using
388 these as we consider this an integral part of their method.
- 389 • *Epidemic model:* We employ the [Cori et al. \(2013\)](#) method, a common building block in the considered
390 approaches, in its basic form without any pre-processing steps.
- 391 • *Window size:* When applying the [Cori et al. \(2013\)](#) method we use a window size of 7 days. This is a
392 common choice as it reduces fluctuations arising from within-week reporting patterns.
- 393 • *Generation time distribution:* We assume an exponential distribution with rate 1/4, i.e., mean and
394 standard deviation equal to 4 days. While an exponential distribution may not be the most common
395 choice to match the epidemiology of COVID-19, this enables us to include the globalrt model, which
396 can only accommodate an exponential GTD.
- 397 • *Incubation period and reporting delay:* These aspects are challenging to standardize across methods,
398 as variation in delays is an integral part of some methods (e.g., epiforecasts) but incompatible with
399 others (e.g., Ilmenau, SDSC). Temporal misalignment resulting from these aspects is therefore handled
400 pragmatically by shifting estimates in time. As the consensus setting, we assume that the reporting
401 delay and incubation period sum up to seven days and shift estimates accordingly.
- 402 • *Definition of R_t :* By using the [Cori et al. \(2013\)](#) method, we estimate instantaneous reproductive
403 numbers. Based on the notion that R_{t-i}^{inst} typically lags behind R_{t-i}^{case} by one mean generation time (see
404 Section 2.1), we again resort to shifting estimates of case reproductive numbers in time.

405 In practice, it proved challenging to determine exactly how estimates needed to be shifted to account for
406 differing assumptions on incubation periods, reporting delays, and type of R_t . Following [Alvarez et al.](#)
407 (2021), we therefore adjust temporal shifts for each method in a data-driven way by minimizing the mean
408 absolute difference to the consensus estimates (see Supplement S4 for details and additional analyses based
409 on reported delay distributions).

410 We moreover note that while all other approaches can be reproduced with standardized settings, some
411 compromises are necessary for the HZI model. The input data already correspond to the consensus choice and

412 the various delay distributions are handled by shifting estimates (as for all other models). The generation time
413 distribution, however, cannot be set directly to the consensus setting, as it is not an independent parameter
414 in the HZI model. Instead, it arises from the interplay of numerous other parameters. We therefore opt
415 to transform the published estimates using a relationship linking the generation time distribution and R_t
416 estimates from [Wallinga and Lipsitch \(2007, see Supplement S3.2\)](#).

417 It is not practically feasible to assess all combinations of standardizing or not standardizing the different
418 analytical choices. We, therefore, vary them in two specific fashions. In the first procedure, we start from
419 the original settings used by different teams. Then, in the above order, we standardize all analytical choices
420 apart from the statistical estimation approach (including possible pre-processing steps). We refer to this as
421 *sequential standardization*. The second procedure starts from the consensus model (i.e., a simple application
422 of the [Cori et al. 2013](#) approach) and subsequently varies the different analytical choices one by one. We refer
423 to this as *individual variation*. An advantage of individual variation is that it does not require specifying
424 an order in which the various dimensions are aligned. The sequential approach, on the other hand, helps to
425 illustrate the compounding of the various effects.

426 As some of the considered approaches are computationally costly (in particular the Bayesian hierarchical
427 models by epiforecasts and rtlive) it is not feasible to re-run the estimations under different parameterizations
428 for all considered estimation dates. We, therefore, refrain from mimicking a real-time setting and assess
429 between-method agreement retrospectively for a single estimation date. Specifically, we consider estimates
430 for the period April 1, 2020, until June 10, 2021, based on data as available on July 10, 2021.

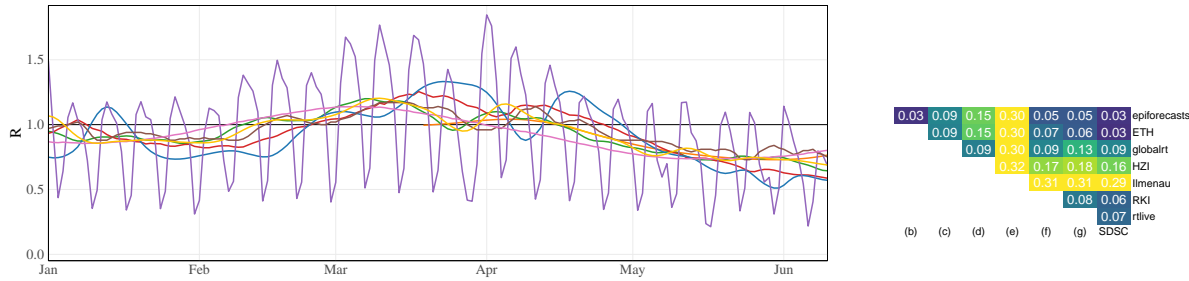
431 4.2 Results for point estimates

432 Figure 6 shows how the agreement between methods improves step by step in the sequential standardization
433 of analytical choices. The visual impression of closer and closer alignment from the left column is confirmed
434 by the matrices of mean absolute differences in the right column. These range from 0.03 (epiforecasts vs.
435 ETH vs. SDSC) up to 0.32 (Ilmenau vs. HZI) for the original versions of the estimates, with an average
436 pairwise value of 0.15. This is a substantial difference given that the estimates are mainly between 0.75
437 and 1.25. Once all analytical choices other than the estimation method and data pre-processing are aligned,
438 mean absolute differences range from 0.01 (ETH vs SDSC) to 0.07 (epiforecasts vs rtlive). Particularly
439 strong improvements result from standardizing the window size where applicable and the generation time
440 distribution. Aligning the window size removes the periodic fluctuations in the Ilmenau estimates, which are
441 based on a very short window of just one day. Standardizing the generation time distribution has a strong
442 impact on the HZI estimates, which are based on a long mean generation time of 10.3 days.

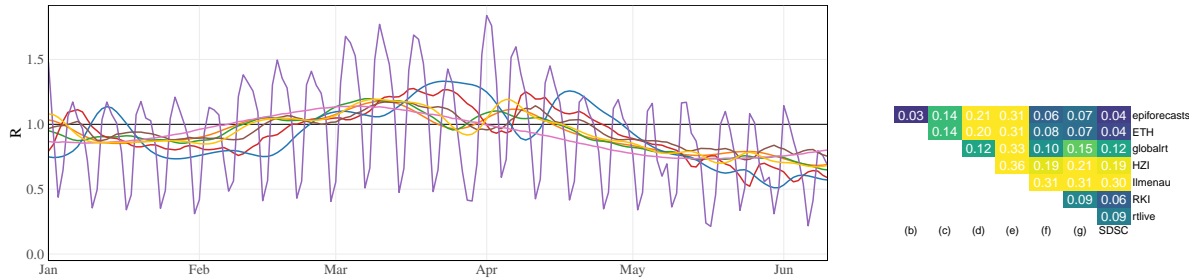
443 As can be seen from the improvement between steps 3 and 4, temporal shifting of estimates is necessary
444 to achieve good alignment. This shift, which is determined in a data-driven way, accounts for differences
445 arising from the assumed incubation periods and reporting delays as well as the choice between case and
446 instantaneous reproductive numbers. In almost all cases the shifts agree well with what would be expected
447 based on the respective model descriptions, see Supplementary Table S3. An alternative display where shifts
448 are determined based on these descriptions is available in Supplementary Figure S14.

449 Results for the individual variation approach are shown in Figure 7. Here, we also vary the data pre-
450 processing step separately; this corresponds to nowcasting for RKI, smoothing for SDSC, and a combination
451 of nowcasting, smoothing, and deconvolution for ETH. Pre-processing as well as the choice of data source
452 impact the smoothness of the estimates, but in terms of mean absolute deviations play a limited role. The
453 window size and generation time distribution have a more substantial impact on the results. The resulting
454 mean absolute differences are in fact more pronounced than when varying the estimation approach (bottom
455 panel). As implied by theory, the estimates are fanned out away from $R_t = 1$ when longer mean generation
456 times are used. In particular, the HZI parameterization with a mean generation time of 10.3 days stands
457 out. Concerning the window length in the Cori approach, choices that are not multiples of 7 lead to periodic
458 fluctuations in the estimates. We note, however, that the ETH and SDSC teams, who use widths of 3 and
459 4 days, employ data pre-processing steps to suppress this behavior.

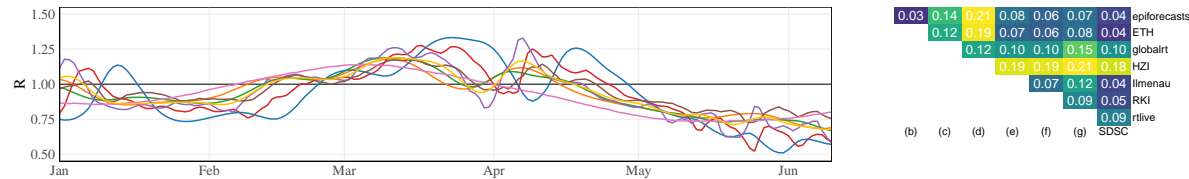
Step 0: Estimates as published on July 10, 2021 (Figure 1).



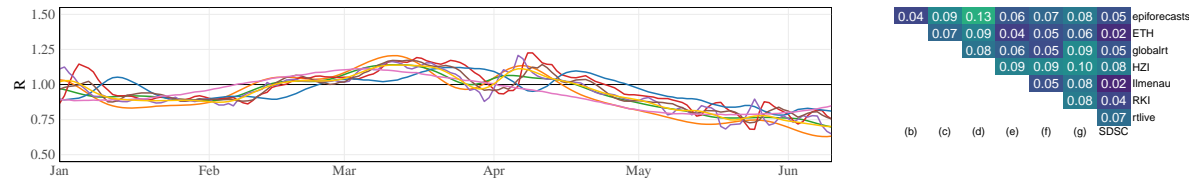
Step 1: Standardize input data to RKI by *Meldedatum*.



Step 2: Standardize window size in [Cori et al. \(2013\)](#) method to 7 days.



Step 3: Standardize GTD to a gamma distribution with mean 4 and standard deviation 4.



Step 4: Data-driven temporal alignment.

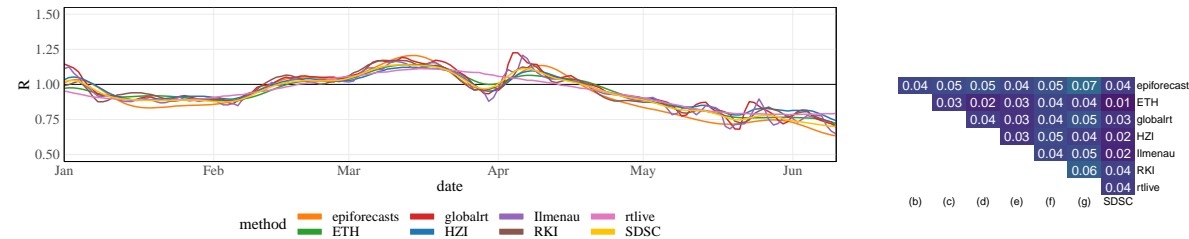


Figure 6: Step-by-step alignment of analytical choices to the consensus specifications. The left column shows the resulting R_t estimates for a subset of the considered time period. The right column shows the mean absolute differences between point estimates obtained from the different approaches. In the bottom panel all considered aspects other than the estimation method (incl. data pre-processing) are aligned.

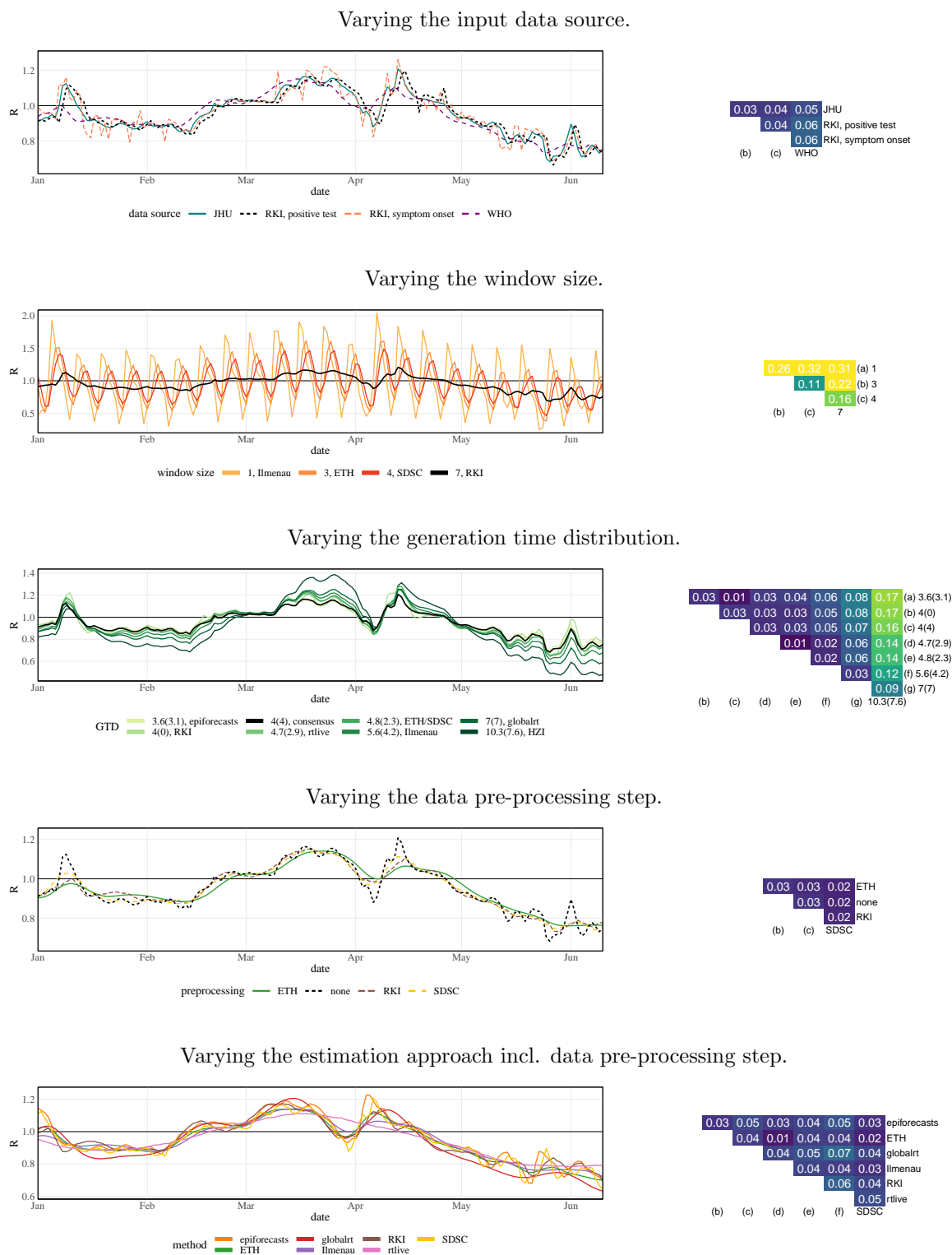


Figure 7: Individual variation of analytical choices in the consensus model. Left column: R_t estimates for a subset of the considered time period. Right column: mean absolute differences between point estimates. The values over which the respective quantities are varied correspond to those chosen by the different teams. For the generation time distribution, we adopt the notation mean (standard deviation). Note that the bottom panel is identical to the one of Figure 6

460 4.3 Some remarks on uncertainty intervals

461 An analog display of the bottom panels of Figures 6 and 7 showing 95% uncertainty intervals can be found
462 in Figure 8. Here, all analytical choices have been standardized apart from the estimation method and
463 data pre-processing. While similarly to the point forecasts, the intervals are more aligned in terms of their
464 temporal course, considerable differences in their widths remain. Rather narrow intervals are produced by
465 the Ilmenau, SDSC, RKI, and ETH approaches (based on the updated version of the method). The intervals
466 obtained from the epiforecasts, globalrt, and rtlive methods are wider. This divide coincides with variations
467 of the [Cori et al. \(2013\)](#) method on the one hand, and more complex hierarchical approaches on the other.

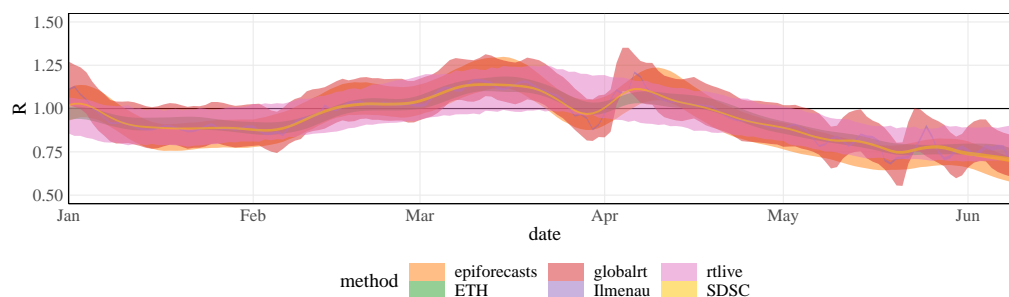


Figure 8: 95% uncertainty intervals corresponding to Figure 6, Step 4.

468 One particularity of the [Cori et al. \(2013\)](#) approach compared to the three others is that it combines
469 equation (1) with a conditional Poisson distribution of observed case counts. The epiforecasts and rtlive
470 approaches assume a negative binomial distribution and globalrt implicitly assumes a Gaussian distribution.
471 These distributions, unlike the Poisson distribution, have a free parameter steering the degree of dispersion.
472 It is known that in generalized regression, assuming a Poisson distribution can lead to an underestimation of
473 standard errors when the data are actually over-dispersed ([Dean and Lundy, 2016](#)). To assess whether this
474 aspect plays a role in the observed patterns we re-ran the [Cori et al. \(2013\)](#) method swapping the Poisson for
475 a negative binomial distribution. As can be seen from Figure 9, this results in considerably wider uncertainty
476 intervals, comparable to those from globalrt.

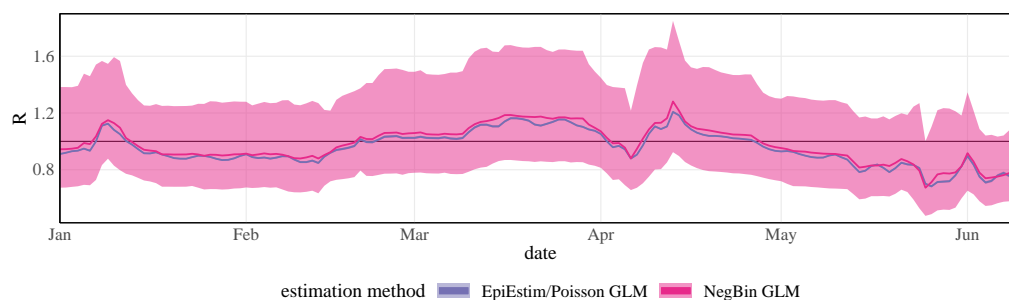


Figure 9: Comparison of 95% uncertainty intervals of the Cori method (consensus settings) with a Poisson (light) and negative binomial distribution (dark). The uncertainty intervals under the Poisson distribution are hardly discernible from the line representing the point estimate.

477 We note that the negative binomial version of the [Cori et al. \(2013\)](#) method needed to be newly im-
478 plemented. For technical ease and to avoid having to specify prior distributions we performed frequentist
479 estimation via the function `glm.nb` from the R package `MASS`. The overdispersion parameter of the negative
480 binomial distribution was estimated jointly with R_t (under the assumption of a constant value over the
481 7-day estimation window); see Supplement S6 for details. As an analog implementation of the Poisson ver-
482 sion yielded almost identical results to *EpiEstim* we consider the use of a frequentist rather than Bayesian
483 implementation unproblematic.

484 Another potentially relevant difference between the [Cori et al. \(2013\)](#) approach and the three others
485 involves the assumptions on the process governing R_t . While the former assumes R_t to be constant on a
486 certain time window, the others assume truly generative models, specifically random walks or a Gaussian
487 process. Both serve to stabilize estimates. However, it is difficult to assess their impact, as replacing these
488 assumptions would be a fundamental change to the respective models.

489 5 Discussion

490 In this paper, we assessed temporal coherence and between-method agreement of R_t estimates for COVID-19
491 in Germany. We found that for most considered methods, the real-time estimates for dates close to the pub-
492 lication date were subject to substantial revisions. In many cases, these were more pronounced than implied
493 by the accompanying uncertainty intervals. Some methods were able to avoid temporal incoherence but at
494 the cost of wide uncertainty intervals. In our retrospective assessment of the between-method agreement, we
495 found that while the choice of estimation method led to some discrepancies, surrounding analytical choices,
496 e.g., on the generation time distribution, were at least as influential.

497 Our assessment of temporal self-coherence highlights the importance of continuously tracking the real-
498 time behavior of R_t estimates. If these are overly fluctuating or subject to systematic corrections, this may
499 lead to a loss in user trust. However, the stability of estimates is not the only relevant goal, and there is a
500 trade-off with the timeliness of estimates. R_t estimates are quickly outdated, and results for recent days are
501 the most relevant for public health purposes. These are unavoidably subject to increased uncertainty. This
502 needs to be acknowledged by users, and uncertainty needs to be quantified and communicated appropriately.
503 We believe that analyses of temporal coherence as presented in our work can be a useful tool to this end.

504 In our between-methods comparison of estimates, we found that in particular the assumed generation
505 time distribution and the choice of window sizes drove differences between estimates published by different
506 research teams. These decisions and their potential impact should thus be communicated transparently.
507 The approach taken by `globalrt`, where users can vary the mean generation time, is promising, though some
508 contextualization on which values are well-supported by the state of research may be helpful. Temporal
509 shifts arising from different assumptions on incubation periods and reporting delays proved relevant, too,
510 as they shift R_t estimates in time. This is of particular importance when linking the latter to intervention
511 measures. The respective delay distributions should thus be chosen with care.

512 It has been argued that to reduce the dependence on specific assumptions, different estimates could be
513 combined into a consensus R_t value or range. While in the United Kingdom meta-analysis techniques have
514 been applied to this end ([Maishman et al., 2022](#)), this is not without pitfalls. Unlike in classical meta-
515 analysis, different estimates are typically obtained from the same data and thus inherently dependent. As
516 pointed out by [Nicholson et al. \(2022\)](#), this leads to estimators with unclear statistical properties. Moreover,
517 when merging estimates based on different assumptions, the estimand becomes unclear, as do the assumptions
518 underlying the consensus estimate. To combine estimates of the basic reproductive number R_0 , an appealing
519 approach where information is pooled separately for the generation time distribution and the epidemic growth
520 rate has been suggested by [Park et al. \(2020\)](#). This could likely be translated to R_t estimation.

521 In the present work, we focused exclusively on estimates based on national-level case incidence data. We
522 did not take into account regional or age stratification, which can be incorporated e.g., in compartmental
523 epidemiological models to estimate R_t ([Knock et al., 2021](#)). Reproductive numbers can also be estimated
524 from other data streams including hospitalizations, deaths ([Sherratt et al., 2021](#); [Huisman et al., 2022a](#)),
525 wastewater surveillance ([Huisman et al., 2022b](#)) and PCR cycle threshold data ([Hay et al., 2021](#)). While these
526 may resolve some of the issues of case incidences, e.g., their sensitivity to testing strategies, the dependence
527 of estimates on analytical choices remains largely the same. Nonetheless, considering estimates based on
528 various data streams may yield a more comprehensive picture. More generally, we underscore that the R_t
529 value should not be interpreted in isolation, but in conjunction with other epidemiological indicators like
530 the overall case and hospitalization numbers or genetic data on the prevalence of different variants.

531 Data and code availability

532 Data and code to reproduce the presented results can be found at [https://github.com/ElisabethBrockhaus/
533 Rt_estimate_reconstruction](https://github.com/ElisabethBrockhaus/Rt_estimate_reconstruction) and https://github.com/KITmetricslab/reproductive_numbers.

534 Acknowledgements

535 The authors would like to thank Simas Kucinskas (globalrt) for making estimates available and providing
536 additional information. Johannes Bracher was supported by the Helmholtz Foundation via the project
537 *SIMCARD*. Johannes Bracher’s work was moreover partly funded by the Deutsche Forschungsgemeinschaft
538 (DFG, German Research Foundation) – project number 512483310. Sam Abbott and Sebastian Funk were
539 supported by The Wellcome Trust (210758/Z/18/Z).

540 References

- 541 ABBOTT, S., J. HELLEWELL, K. SHERRATT, K. GOSTIC, J. HICKSON, H. S. BADR, M. DEWITT,
542 R. THOMPSON, EPIFORECASTS, AND S. FUNK (2020a): *EpiNow2: Estimate Real-Time Case Counts and*
543 *Time-Varying Epidemiological Parameters*.
- 544 ABBOTT, S., J. HELLEWELL, R. N. THOMPSON, K. SHERRATT, H. P. GIBBS, N. I. BOSSE, J. D.
545 MUNDAY, S. MEAKIN, E. L. DOUGHTY, J. Y. CHUN, ET AL. (2020b): “Estimating the time-varying
546 reproduction number of SARS-CoV-2 using national and subnational case counts,” *Wellcome Open Re-*
547 *search*, 5, 112.
- 548 ALVAREZ, L., M. COLOM, J.-D. MOREL, AND J.-M. MOREL (2021): “Computing the daily reproduction
549 number of COVID-19 by inverting the renewal equation using a variational technique,” *Proceedings of the*
550 *National Academy of Sciences*, 118, e2105112118.
- 551 AN DER HEIDEN, M. AND O. HAMOUDA (2020): “Schätzung der aktuellen Entwicklung der SARS-CoV-2-
552 Epidemie in Deutschland – Nowcasting,” *Epidemiologisches Bulletin*, 2020, 10–15.
- 553 ARROYO-MARIOLI, F., F. BULLANO, S. KUCINSKAS, AND C. RONDÓN-MORENO (2021): “Tracking R of
554 COVID-19: A new real-time estimation using the Kalman filter,” *PloS one*, 16, e0244474.
- 555 BRAUNER, J. M., S. MINDERMAN, M. SHARMA, D. JOHNSTON, J. SALVATIER, T. GAVENČIAK, A. B.
556 STEPHENSON, G. LEECH, G. ALTMAN, V. MIKULIK, A. J. NORMAN, J. T. MONRAD, T. BESIROGLU,
557 H. GE, M. A. HARTWICK, Y. W. TEH, L. CHINDELEVITCH, Y. GAL, AND J. KULVEIT (2021): “Inferring
558 the effectiveness of government interventions against COVID-19,” *Science*, 371, eabd9338.
- 559 CORI, A., N. FERGUSON, C. FRASER, E. DAHLQWIST, P. DEMARSH, T. JOMBART, Z. KAMVAR,
560 J. LESSLER, S. LI, J. POLONSKY, ET AL. (2020): “Package ‘EpiEstim’,” *CRAN: Vienna, Austria*.
- 561 CORI, A., N. M. FERGUSON, C. FRASER, AND S. CAUCHEMEZ (2013): “A new framework and software to
562 estimate time-varying reproduction numbers during epidemics,” *American Journal of Epidemiology*, 178,
563 1505–1512.
- 564 DEAN, C. B. AND E. R. LUNDY (2016): *Wiley StatsRef: Statistics Reference Online*, John Wiley & Sons,
565 Ltd, chap. Overdispersion, 1–9.
- 566 DONG, E., H. DU, AND L. GARDNER (2020): “An interactive web-based dashboard to track COVID-19 in
567 real time,” *The Lancet infectious diseases*, 20, 533–534.
- 568 FLAXMAN, S., S. MISHRA, A. GANDY, H. J. T. UNWIN, T. A. MELLAN, H. COUPLAND, C. WHITTAKER,
569 H. ZHU, T. BERAH, J. W. EATON, ET AL. (2020): “Estimating the effects of non-pharmaceutical
570 interventions on COVID-19 in Europe,” *Nature*, 584, 257–261.
- 571 FUNK, S., S. ABBOTT, AND J. BRACHER (2022): “Sebastian Funk, Sam Abbott and Johannes Bracher’s
572 Discussion Contribution to the Papers in Session 2 of The Royal Statistical Society’s Special Topic Meeting
573 on Covid-19 Transmission: 11 June 2021,” *Journal of the Royal Statistical Society Series A: Statistics in*
574 *Society*, 185, S103–S104.
- 575 GANYANI, T., C. KREMER, D. CHEN, A. TORNERI, C. FAES, J. WALLINGA, AND N. HENS (2020):
576 “Estimating the generation interval for coronavirus disease (COVID-19) based on symptom onset data,”
577 *Eurosurveillance*, 25, 2000257.

- 578 GOSTIC, K. M., L. MCGOUGH, E. B. BASKERVILLE, S. ABBOTT, K. JOSHI, C. TEDIJANTO, R. KAHN,
579 R. NIEHUS, J. A. HAY, P. M. DE SALAZAR, ET AL. (2020): “Practical considerations for measuring the
580 effective reproductive number, R_t ,” *PLoS Computational Biology*, 16, e1008409.
- 581 HAUG, N., L. GEYRHOFFER, A. LONDEI, E. DERVIC, A. DESVARS-LARRIVE, V. LORETO, B. PINIOR,
582 S. THURNER, AND P. KLIMEK (2020): “Ranking the effectiveness of worldwide COVID-19 government
583 interventions,” *Nature Human Behaviour*, 4, 1303–1312.
- 584 HAY, J. A., L. KENNEDY-SHAFFER, S. KANJILAL, N. J. LENNON, S. B. GABRIEL, M. LIPSITCH, AND
585 M. J. MINA (2021): “Estimating epidemiologic dynamics from cross-sectional viral load distributions,”
586 *Science*, 373, eabh0635.
- 587 HOTZ, T., M. GLOCK, S. HEYDER, S. SEMPER, A. BÖHLE, AND A. KRÄMER (2020): “Monitoring the
588 spread of COVID-19 by estimating reproduction numbers over time,” *arXiv preprint arXiv:2004.08557*.
- 589 HUISMAN, J. S., J. SCIRE, D. C. ANGST, J. LI, R. A. NEHER, M. H. MAATHUIS, S. BONHOEFFER,
590 AND T. STADLER (2022a): “Estimation and worldwide monitoring of the effective reproductive number
591 of SARS-CoV-2,” *eLife*, 11, e71345.
- 592 HUISMAN, J. S., J. SCIRE, L. CADUFF, X. FERNANDEZ-CASSI, P. GANESANANDAMOORTHY, A. KULL,
593 A. SCHEIDEGGER, E. STACHLER, A. B. BOEHM, B. HUGHES, A. KNUDSON, A. TOPOL, K. R. WIGGIN-
594 TON, M. K. WOLFE, T. KOHN, C. ORT, T. STADLER, AND T. R. JULIAN (2022b): “Wastewater-Based
595 Estimation of the Effective Reproductive Number of SARS-CoV-2,” *Environmental Health Perspectives*,
596 130, 057011.
- 597 KHAILAIE AND MITRA, A. BANDYOPADHYAY, M. SCHIPS, P. MASCHERONI, P. VANELLA, B. LANGE,
598 S. BINDER, AND M. MEYER-HERMANN (2020): “SECIR Report,” [https://gitlab.com/simm/covid19/
599 secir/-/wikis/Report](https://gitlab.com/simm/covid19/secir/-/wikis/Report), accessed: 2022-07-18.
- 600 KHAILAIE AND MITRA, A. BANDYOPADHYAY, M. SCHIPS, P. MASCHERONI, P. VANELLA, B. LANGE,
601 S. C. BINDER, AND M. MEYER-HERMANN (2021): “Development of the reproduction number from
602 coronavirus SARS-CoV-2 case data in Germany and implications for political measures,” *BMC medicine*,
603 19, 1–16.
- 604 KLAU, S., S. HOFFMANN, C. J. PATEL, J. P. IOANNIDIS, AND A.-L. BOULESTEIX (2020): “Examining
605 the robustness of observational associations to model, measurement and sampling uncertainty with the
606 vibration of effects framework,” *International Journal of Epidemiology*, 50, 266–278.
- 607 KNABL, MITRA AND KIMPEL, A. RÖSSLER, A. VOLLAND, A. WALSER, H. ULMER, L. PIPPERGER,
608 S. BINDER, L. RIEPLER, K. BATES, A. BANDYOPADHYAY, M. SCHIPS, M. RANJAN, B. FALKENSAM-
609 MER, W. BORENA, M. MEYER-HERMANN, AND V. VON LAER (2021): “High SARS-CoV-2 seroprevalence
610 in children and adults in the Austrian ski resort of Ischgl,” *Communications Medicine*, 1.
- 611 KNOCK, E. S., L. K. WHITTLES, J. A. LEES, P. N. PEREZ-GUZMAN, R. VERITY, R. G. FITZJOHN,
612 K. A. M. GAYTHORPE, N. IMAI, W. HINSLEY, L. C. OKELL, A. ROSELLO, N. KANTAS, C. E. WAL-
613 TERS, S. BHATIA, O. J. WATSON, C. WHITTAKER, L. CATTARINO, A. BOONYASIRI, B. A. DJAAFARA,
614 K. FRASER, H. FU, H. WANG, X. XI, C. A. DONNELLY, E. JAUNEIKAITE, D. J. LAYDON, P. J.
615 WHITE, A. C. GHANI, N. M. FERGUSON, A. CORI, AND M. BAGUELIN (2021): “Key epidemiological
616 drivers and impact of interventions in the 2020 SARS-CoV-2 epidemic in England,” *Science Translational
617 Medicine*, 13, eabg4262.
- 618 KRYMOVA, E., B. BÉJAR, D. THANOU, T. SUN, E. MANETTI, G. LEE, K. NAMIGAI, C. CHOIRAT,
619 A. FLAHAULT, AND G. OBOZINSKI (2022): “Trend estimation and short-term forecasting of COVID-19
620 cases and deaths worldwide,” *Proceedings of the National Academy of Sciences*, 119, e2112656119.
- 621 LAUCK, D. (2020): “Corona-Zahlen des RKI: Täglicher R-Wert stimmt oft nicht,” [https://www.tagesschau.
622 de/faktenfinder/r-wert-rki-101.html](https://www.tagesschau.de/faktenfinder/r-wert-rki-101.html), accessed: 2022-10-11.

- 623 MAISHMAN, T., S. SCHAAP, D. S. SILK, S. J. NEVITT, D. C. WOODS, AND V. E. BOWMAN (2022):
624 “Statistical methods used to combine the effective reproduction number, $R(t)$, and other related measures
625 of COVID-19 in the UK,” *Statistical Methods in Medical Research*, 31, 1757–1777, pMID: 35786070.
- 626 NICHOLSON, G., M. BLANGIARDO, M. BRIERS, P. J. DIGGLE, T. E. FJELDE, H. GE, R. J. B. GOUDIE,
627 R. JERSAKOVA, R. E. KING, B. C. L. LEHMANN, A.-M. MALLON, T. PADELLINI, Y. W. TEH,
628 C. HOLMES, AND S. RICHARDSON (2022): “Interoperability of Statistical Models in Pandemic Prepared-
629 ness: Principles and Reality,” *Statistical Science*, 37, 183 – 206.
- 630 NISHIURA, H., N. M. LINTON, AND A. R. AKHMETZHANOV (2020): “Serial interval of novel coronavirus
631 (COVID-19) infections,” *International Journal of Infectious Diseases*, 93, 284–286.
- 632 OSTHEGE, M., L. HELLECKES, A. ANDORRA, AND MAPEPER (2021): “rtlive-dash-de,” <https://github.com/michaelosthege/rtlive-dash-de>,
633 accessed: 2023-03-07.
- 634 O’DRISCOLL, M., C. HARRY, C. A. DONNELLY, A. CORI, AND I. DORIGATTI (2021): “A comparative
635 analysis of statistical methods to estimate the reproduction number in emerging epidemics, with implica-
636 tions for the current Coronavirus Disease 2019 (COVID-19) Pandemic,” *Clinical Infectious Diseases*, 73,
637 e215–e223.
- 638 PARK, S. W., B. M. BOLKER, D. CHAMPREDON, D. J. D. EARN, M. LI, J. S. WEITZ, B. T. GRENFELL,
639 AND J. DUSHOFF (2020): “Reconciling early-outbreak estimates of the basic reproductive number and its
640 uncertainty: framework and applications to the novel coronavirus (SARS-CoV-2) outbreak,” *Journal of*
641 *The Royal Society Interface*, 17, 20200144.
- 642 PASETTO, D., J. C. LEMAITRE, E. BERTUZZO, M. GATTO, AND A. RINALDO (2021): “Range of reproduc-
643 tion number estimates for COVID-19 spread,” *Biochemical and Biophysical Research Communications*,
644 538, 253–258, cOVID-19.
- 645 REDAKTIONSNETZWERK DEUTSCHLAND (2020): “Reproduktionszahl: Wieso RKI und HZI unterschiedliche
646 Werte melden,” [https://www.rnd.de/gesundheit/reproduktionszahl-darum-kommen-rki-und-hzi-bei-der-
647 berechnung-zu-unterschiedlichen-ergebnissen-UKYFVVLJH7MLQEWODUWVHW2BVY.html](https://www.rnd.de/gesundheit/reproduktionszahl-darum-kommen-rki-und-hzi-bei-der-berechnung-zu-unterschiedlichen-ergebnissen-UKYFVVLJH7MLQEWODUWVHW2BVY.html), pub-
648 lished: 2020-04-30.
- 649 RKI (2020a): “COVID-19 Datenhub,” [https://npgeo-corona-npgeo-de.hub.arcgis.com/datasets/
650 dd4580c810204019a7b8eb3e0b329dd6_0/explore](https://npgeo-corona-npgeo-de.hub.arcgis.com/datasets/dd4580c810204019a7b8eb3e0b329dd6_0/explore), accessed: 2021-11-23.
- 651 ——— (2020b): “Erläuterung der Schätzung der zeitlich variierenden Reproduktionszahl R_t ,” Robert
652 Koch Institut Berlin, [https://www.rki.de/DE/Content/InfAZ/N/Neuartiges_Coronavirus/Projekte_RKI/
653 R-Wert-Erlaeuterung.pdf](https://www.rki.de/DE/Content/InfAZ/N/Neuartiges_Coronavirus/Projekte_RKI/R-Wert-Erlaeuterung.pdf).
- 654 SHERRATT, K., S. ABBOTT, S. R. MEAKIN, J. HELLEWELL, J. D. MUNDAY, N. BOSSE, C. C.-. WORKING
655 GROUP, M. JIT, AND S. FUNK (2021): “Exploring surveillance data biases when estimating the repro-
656 duction number: with insights into subpopulation transmission of COVID-19 in England,” *Philosophical*
657 *Transactions of the Royal Society B*, 376, 20200283.
- 658 SYSTROM, K., M. KRIEGER, T. VLADECK, M. OSTHEGE, L. HELLECKES, A. ANDORRA, T. WIECKI,
659 D. FERRERO, MAPEPER, T. MICK, AND L. NACHTERGAELE (2020): “Rt.live and Rtlive-global,” GitHub
660 repositories, stable Zenodo release at <https://doi.org/10.5281/zenodo.7300132>, version v1.0.0.
- 661 TEH, Y. W., B. ELESSEDY, B. HE, M. HUTCHINSON, S. ZAIDI, A. BHOOPCHAND, U. PAQUET, N. TOMA-
662 SEV, J. READ, AND P. J. DIGGLE (2022): “Efficient Bayesian inference of Instantaneous Reproduction
663 Numbers at Fine Spatial Scales, with an Application to Mapping and Nowcasting the Covid-19 Epidemic
664 in British Local Authorities,” *Journal of the Royal Statistical Society Series A: Statistics in Society*, 185,
665 S65–S85.
- 666 VEGVARI, C., S. ABBOTT, F. BALL, E. BROOKS-POLLOCK, R. CHALLEN, B. S. COLLYER, C. DANGER-
667 FIELD, J. R. GOG, K. M. GOSTIC, J. M. HEFFERNAN, T. D. HOLLINGSWORTH, V. ISHAM, E. KENAH,
668 D. MOLLISON, J. PANOVSKA-GRIFFITHS, L. PELLIS, M. G. ROBERTS, G. S. TOMBA, R. N. THOMP-
669 SON, AND P. TRAPMAN (2022): “Commentary on the use of the reproduction number R during the
670 COVID-19 pandemic,” *Statistical Methods in Medical Research*, 31, 1675–1685, pMID: 34569883.

- 671 WAGENMAKERS, E.-J., A. SARAFILOGLOU, AND B. ACZEL (2022): “One statistical analysis must not rule
672 them all,” .
- 673 WALLINGA, J. AND M. LIPSITCH (2007): “How generation intervals shape the relationship between growth
674 rates and reproductive numbers,” *Proceedings of the Royal Society B: Biological Sciences*, 274, 599–604.
- 675 WHITE, L. F., C. B. MOSER, R. N. THOMPSON, AND M. PAGANO (2021): “Statistical estimation of the
676 reproductive number from case notification data,” *American Journal of Epidemiology*, 190, 611–620.
- 677 WORLD HEALTH ORGANIZATION (2022): “WHO Coronavirus (COVID-19) Dashboard,” [https://covid19.
678 who.int/](https://covid19.who.int/).

679 Supplementary material for Brockhaus et al.: Why are different esti-
680 mates of the effective reproductive number so different? A case study
681 on COVID-19 in Germany

682 Contact: E. K. Brockhaus (elisabeth.brockhaus@student.kit.edu), J. Bracher (johannes.bracher@kit.edu)

683 S1 Repositories from which real-time estimates were obtained

- 684 • epiforecasts: <https://github.com/epiforecasts/covid-rt-estimates>
- 685 • ETH: <https://github.com/covid-19-Re/dailyRe-Data>
- 686 • globalrt: <https://github.com/cronndonm/TrackingR>
- 687 • HZI: <https://gitlab.com/simm/covid19/secir/-/tree/master>
- 688 • Ilmenau: <https://github.com/Stochastik-TU-Ilmenau/COVID-19/tree/gh-pages>
- 689 • RKI: https://github.com/robert-koch-institut/SARS-CoV-2-Nowcasting_und_-R-Schaetzung
- 690 • rtlive: <https://zenodo.org/record/5683308>, <https://github.com/michaelosthege/rtlive-global>
- 691 • SDSC: <https://renkulab.io/gitlab/covid-19/covid-19-forecast/-/tree/master>

692 S2 Sources for generation time distributions shown in Figure 2

693 We here provide the sources for the generation time distributions used by European public health agencies
694 as displayed in Figure 2.

- 695 • Austria: Richter, Schmid and Stadlober: *Methodenbeschreibung für die Schätzung von epidemio-*
696 *logischen Parametern des COVID19 Ausbruchs, Österreich.* [https://www.ages.at/fileadmin/Corona/](https://www.ages.at/fileadmin/Corona/Epidemiologische-Parameter/Methoden_zur_Sch%C3%A4tzung_der_epi_Parameter.pdf)
697 [Epidemiologische-Parameter/Methoden_zur_Sch%C3%A4tzung_der_epi_Parameter.pdf](https://www.ages.at/fileadmin/Corona/Epidemiologische-Parameter/Methoden_zur_Sch%C3%A4tzung_der_epi_Parameter.pdf). The Cori et al.
698 (2013) method is used for estimation. The generation time distribution was initially set to a gamma
699 distribution with mean 4.46 days and standard deviation 2.63 days. Later this was revised to 3.37 and
700 1.83 days, respectively.
- 701 • Belgium: Sciensano: *COVID-19 Bulletin épidémiologique hebdomadaire (19 mai 2022)*, [https://covid-](https://covid-19.sciensano.be/sites/default/files/Covid19/COVID-19_Weekly%20report_20220519%20-%20FR.pdf)
702 [19.sciensano.be/sites/default/files/Covid19/COVID-19_Weekly%20report_20220519%20-%20FR.pdf](https://covid-19.sciensano.be/sites/default/files/Covid19/COVID-19_Weekly%20report_20220519%20-%20FR.pdf).
703 The Cori et al. (2013) method is used for estimation. The generation time distribution is set to a
704 gamma distribution with mean 4.7 days and standard deviation 2.9 days (source of parameterization:
705 personal correspondence).
- 706 • Czech Republic: Majék et al (2020): *Modelling the first wave of the COVID-19 epidemic in the*
707 *Czech Republic and the role of government interventions.* medRxiv, [https://doi.org/10.1101/2020.09.](https://doi.org/10.1101/2020.09.10.20192070)
708 [10.20192070](https://doi.org/10.1101/2020.09.10.20192070). The generation time distribution is a discrete uniform over $\{4, 5, 6, 7\}$, implying a mean
709 of 5.5 and standard deviation of 1.2.
- 710 • Denmark: Statens Seruminstitut: *COVID-19 i Danmark: Epidemiologisk trend og fokus: kontakttal,*
711 *11. juni 2020* (2020). <https://files.ssi.dk/COVID19-epi-trendogfokus-11062020>. The generation time
712 distribution from Nishiura et al. (2020) is used in the Cori et al. (2013) method, which corresponds to
713 a mean of 4.7 and a standard deviation of 2.9 days.
- 714 • France: Santé Publique France (2021): *COVID-19 – Point épidémiologique hebdomadaire no 71 du 08*
715 *juillet 2021.* [https://www.santepubliquefrance.fr/content/download/358653/document_file/COVID19-](https://www.santepubliquefrance.fr/content/download/358653/document_file/COVID19-PE_20210708_signets.pdf)
716 [PE_20210708_signets.pdf](https://www.santepubliquefrance.fr/content/download/358653/document_file/COVID19-PE_20210708_signets.pdf). The R_t estimates are obtained using the method by Cori et al. (2013) with
717 a window size of 7 days. The mean and standard deviation of the generation time distribution are not
718 reported, but using trial and error could be reconstructed as approximately 7 and 4.5 days, respectively.

- 719 • Italy: Guzzetta and Merler (2020): *Stime della trasmissibilità di SARS-CoV-2 in Italia*. Istituto
720 Superiore di Sanità / EpiCentro. <https://www.epicentro.iss.it/coronavirus/open-data/rt.pdf>. The
721 method by Cori et al. (2013) is applied; the generation time distribution is a gamma distribution with
722 mean 6.7 days and standard deviation 4.9 days.
- 723 • Netherlands: Rijksinstituut voor Volksgezondheid en Milieu (2021): *Covid-19 reproductiegetal*. <https://data.rivm.nl/meta/srv/eng/catalog.search;jsessionid=1B3A9B193CB3B1946836BCA3D1BF3A11>.
724 The Wallinga-Lipsitch method as implemented in the *EpiEstim* R package (Cori et al., 2020) is used.
725 For pre-Omicron variants, the mean and standard deviation of the generation time are set to 4 and 2
726 days, respectively. For Omicron, 3.5 and 1.75 days are used (source for standard deviations: personal
727 correspondence).
728
- 729 • Portugal: Instituto Nacional de Saúde (2022): *COVID-19 – curva epidemica e parâmetros de trans-*
730 *misidade, 18.05.2022*. [https://www.insa.min-saude.pt/category/areas-de-atuacao/epidemiologia/covid-](https://www.insa.min-saude.pt/category/areas-de-atuacao/epidemiologia/covid-19-curva-epidemica-e-parametros-de-transmissibilidade/)
731 [19-curva-epidemica-e-parametros-de-transmissibilidade/](https://www.insa.min-saude.pt/category/areas-de-atuacao/epidemiologia/covid-19-curva-epidemica-e-parametros-de-transmissibilidade/). The mean and standard deviation of the
732 generation time are set to 3.96 and 4.74 days, respectively (based on Du et al, [https://wwwnc.cdc.](https://wwwnc.cdc.gov/eid/article/26/6/20-0357_article)
733 [gov/eid/article/26/6/20-0357_article](https://wwwnc.cdc.gov/eid/article/26/6/20-0357_article)).
- 734 • Scotland: Scottish Government (2020): *Coronavirus (COVID-19): modeling the epidemic in Scotland*
735 *(Issue No. 24)*. [https://www.gov.scot/binaries/content/documents/govscot/publications/research-](https://www.gov.scot/binaries/content/documents/govscot/publications/research-and-analysis/2020/10/coronavirus-covid-19-modelling-epidemic-issue-no-24/documents/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24/govscot%3Adocument/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24.pdf)
736 [and-analysis/2020/10/coronavirus-covid-19-modelling-epidemic-issue-no-24/documents/coronavirus-covid-](https://www.gov.scot/binaries/content/documents/govscot/publications/research-and-analysis/2020/10/coronavirus-covid-19-modelling-epidemic-issue-no-24/documents/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24/govscot%3Adocument/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24.pdf)
737 [19-modelling-epidemic-scotland-issue-no-24/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-](https://www.gov.scot/binaries/content/documents/govscot/publications/research-and-analysis/2020/10/coronavirus-covid-19-modelling-epidemic-issue-no-24/documents/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24/govscot%3Adocument/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24.pdf)
738 [24/govscot%3Adocument/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24.pdf](https://www.gov.scot/binaries/content/documents/govscot/publications/research-and-analysis/2020/10/coronavirus-covid-19-modelling-epidemic-issue-no-24/documents/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24/govscot%3Adocument/coronavirus-covid-19-modelling-epidemic-scotland-issue-no-24.pdf). Estima-
739 tion is based on the model by Flaxman et al. (2020), which uses a gamma distribution with mean 6.5
740 and standard deviation 4.11 days (see their Supplementary Information, page 13).
- 741 • Sweden: Folkhälsomyndigheten (2022): *Skattning av det momentana reproduktionstalet, 18/05/2022*.
742 Mean and standard deviation of the generation time are set to 4.8 and 2.3 days, respectively. Unfor-
743 tunately, the respective document is no longer available online.
- 744 • Slovenia: Rok Blagus, Manevski and Pohar Perme (2020): *Estimation of the reproductive number and*
745 *the outbreak size of SARS-CoV-2 in Slovenia*. Slovenian Medical Journal, [http://dx.doi.org/10.6016/](http://dx.doi.org/10.6016/ZdravVestn.3068)
746 [ZdravVestn.3068](http://dx.doi.org/10.6016/ZdravVestn.3068). As for Scotland, the model and assumed generation time corresponds to the one
747 from Flaxman et al. (2020).

748 S3 Additional remarks on the HZI approach

749 S3.1 Determining the generation time distribution

750 The generation time distribution is not an independent parameter in the SECIR model applied by Khailaie
751 and Mitra et al. (2021) but results from the interplay of several other parameters. The generation time
752 distribution of the model has no closed form. As it arises from the transitions between different compartments
753 in a classic compartmental model, it corresponds to a mixture of convolutions of exponential distributions.
754 For the purposes of our study, we obtain the mean and standard deviation via simulation. Figure S10
755 shows compartments of the model, the transitions between which are governed by the following rates and
756 probabilities (see also Knabl, Mitra and Kimpel et al. 2021, particular Supplementary Material p12–13):

- 757 • $\alpha = 0.22$, i.e. a 78% probability of entering the state C_I (carrier who will move on to infected) after
758 exposure.
- 759 • $R_2 = 1/3.2$ and $R_3 = 1/2$ imply a mean incubation period of 5.2 days.
- 760 • $R_4 = 1/7$, i.e., a mean time to recovery of seven days for undetected infected.
- 761 • $\mu = 0.085$ is the probability of detection for any infected individual. This value has been taken from
762 the code repository² rather than the manuscript.

²https://gitlab.com/simm/covid19/secir/-/raw/master/codes/settings/param_random.csv

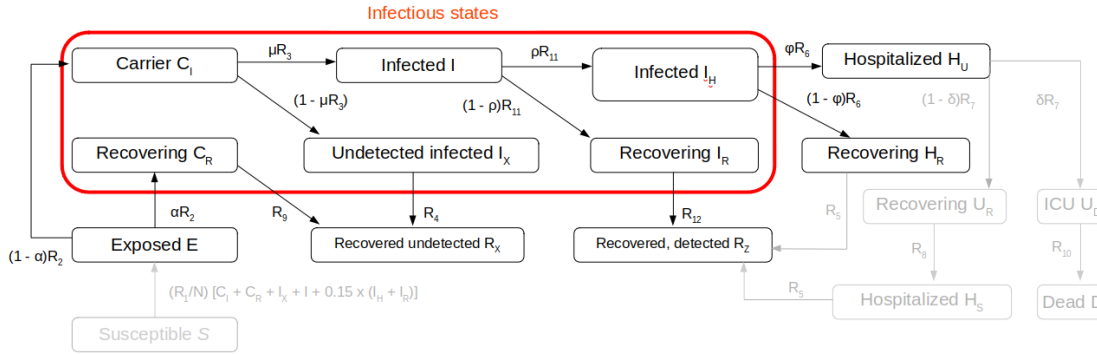


Figure S10: Compartmental structure of the SECIR model by HZI. Compartments relevant to the generation time distribution are shown in black, the remaining compartments in grey. This figure reproduces Figure 1 from [Khailaie and Mitra et al. \(2021\)](#).

- 763 • $R_{11} = 1/3.7$, i.e., on average it takes 3.7 days to move from I_H to H_U . R_6 is chosen such that $1/R_6$
764 and $1/R_{11}$ sum to 4.25 days, which is the assumed time to hospitalization since symptom onset. This
765 implies $R_6 = 1/0.55$.
- 766 • $\rho = 0.13$, i.e., the probability that a detected infected patient requires hospitalization (moves to I_H)
767 rather than recovering (moving to I_R). This value has been taken from the code repository.
- 768 • $\varphi = 0.47$, i.e. a hospitalization probability of 47% for individuals reaching the I_H state (but this value
769 is irrelevant for the generation time distribution).
- 770 • $R_{12} = 1/3.3$, i.e., the average time left to full recovery for individuals who have already arrived in I_R
771 is 3.3 days.
- 772 • The remaining parameters are not relevant for the computation of the generation time distribution
773 and are thus omitted here.

774 Moreover, [Khailaie and Mitra et al. \(2021\)](#) assume that individuals in the C_I, C_R, I_X , and I compartments
775 have the same infectiousness, while infectiousness in the I_H and I_R compartments is reduced by a factor of
776 0.15. We note that we here only use the assumed mean values of the different parameters and simplifyingly
777 neglect that they are randomly varied around these values in [Khailaie and Mitra et al. \(2021\)](#). To obtain
778 the generation time distribution numerically we then proceed as follows.

- 779 1. For a total of 5000 individuals we first sample the path the individual takes through the different
780 compartments from E onwards (e.g., $E \rightarrow C_R \rightarrow R_X$ or $E \rightarrow C_I \rightarrow I \rightarrow I_H \rightarrow H_U$). This involves
781 the probabilities α, μ, ρ and φ .
- 782 2. We then sample the duration of stay in each of the compartments from exponential distributions with
783 the respective transition rates.
- 784 3. For the time spent in infectious compartments ($C_I, C_R, I, I_X, I_H, I_R$) we sample times of secondary
785 infections from Poisson processes with suitably chosen rates (with diminished intensity for the I_H and
786 I_R compartments). In practice, this is done by first sampling the total number of events from a suitable
787 Poisson distribution and then sampling the respective event times from a uniform distribution over the
788 time spent in the compartment.
- 789 4. For each infection event we compute the total time since the entry of the infecting individual in the E
790 compartment. This corresponds to the realized generation time.

791 We then evaluate the empirical distribution of these generation times. The resulting histogram is shown in
792 Figure S11. The mean and standard deviations are given by 10.3 and 7.6, respectively.

793 Code to reproduce these results is available in https://github.com/ElisabethBrockhaus/Rt_estimate_reconstruction/blob/main/HZI.

Distribution of generation times

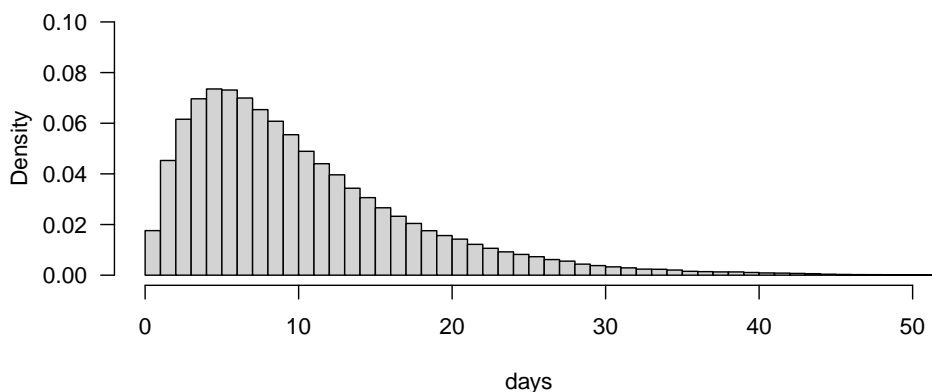


Figure S11: Histogram of sampled generation times from the HZI model.

795 S3.2 Transformation of estimates to standardize generation time distribution

The HZI results proved to be technically challenging to reproduce, and due to the setup of the model, the generation time distribution cannot be manipulated directly in the model. Rather than re-running the model with a different generation time distribution, we, therefore, opted to transform the available estimates to approximate how they would have looked under a different generation time distribution. We employ the following relationship from [Wallinga and Lipsitch \(2007, Equation 3.6\)](#):

$$R = \frac{r}{\sum_{i=1}^n y_i \{ \exp(-ra_{i-1}) - \exp(-ra_i) \} / (a_i - a_{i-1})}.$$

796 Here, R is the reproductive number, r is the growth rate, a_0, a_1, \dots, a_n are the category bounds of a
797 histogram, and y_1, y_2, \dots, y_n the respective relative frequencies. By plugging in the distribution from Figure
798 S11 and samples from the consensus distribution $\text{Exp}(1/4)$, we can obtain mappings from the growth rate to
799 estimated reproductive numbers under the two generation time distributions (we here use bins of width 0.01
800 for the histogram). Combining these two, we can map the reproductive numbers computed by HZI under the
801 generation time distribution from Figure S11 to reproductive numbers under the consensus generation time
802 distribution. This mapping is displayed in Figure S12. As one would expect, values of 0 and 1 are mapped
803 to themselves, while otherwise, the values under the consensus distribution (which has a considerably lower
804 mean) are closer to 1.

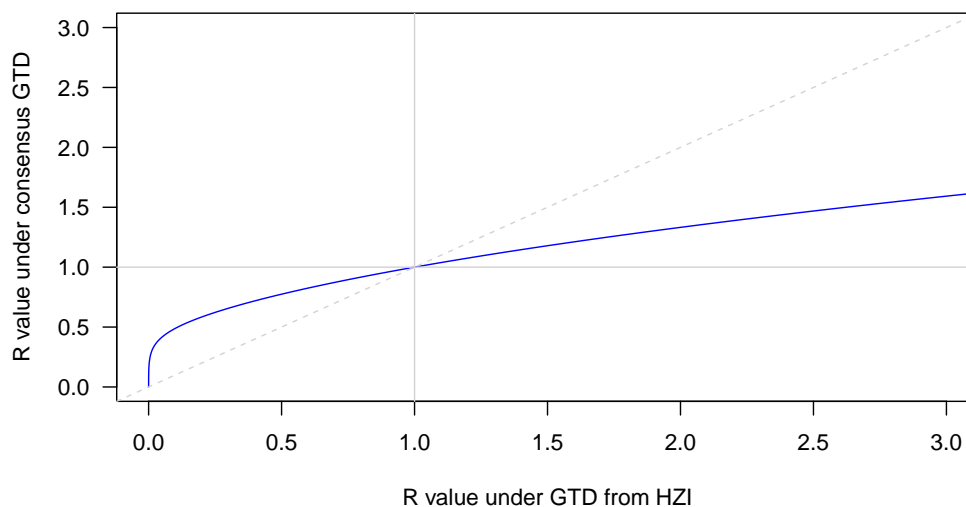


Figure S12: Mapping of R -values under the generation time distribution used by HZI to R -values under the consensus generation time distribution $\text{Exp}(1/4)$.

805 S3.3 Comparison of estimates with and without additional smoothing

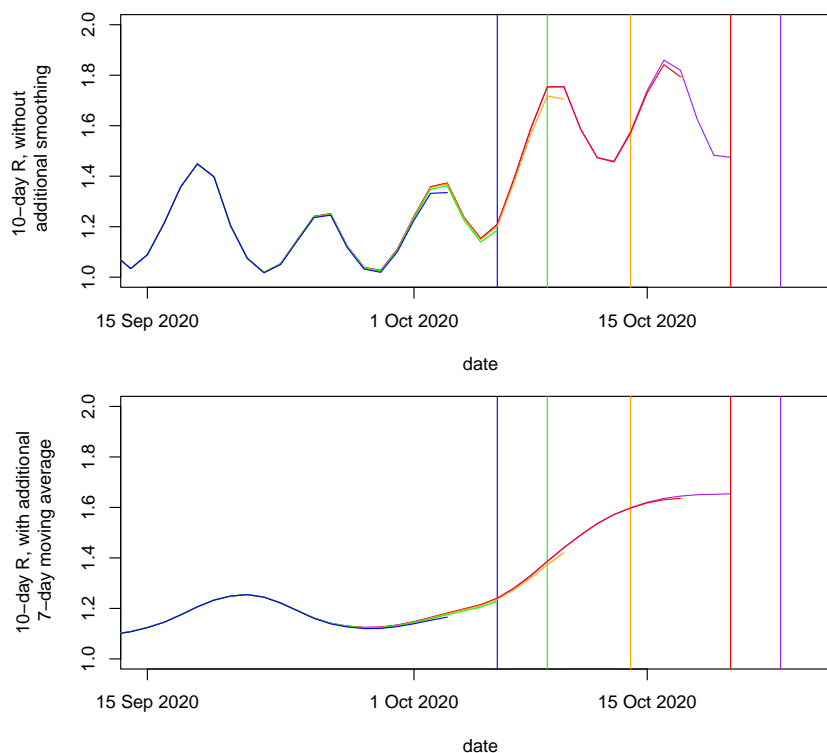


Figure S13: Comparison of real-time estimates from the HZI model with and without an additional smoothing step (7-day trailing moving average). Similarly to Figure 4 estimates issued at different dates are overlaid so that retrospective revisions become visible. It can be seen that the additional smoothing step considerably reduces the (typically upwards) retrospective revisions.

806 S4 Details on the handling of temporal shifts

To handle differing assumptions on the incubation period and reporting delay as well as the exact definition of R_t (case vs. instantaneous), we align estimates via simple shifting. We do this in a data-driven way, where for each method we minimize the mean absolute distance to the consensus model. Denoting the days in the considered period by indices $t = 1, \dots, T$, we thus obtain the shift s_m for each model m as

$$s_m = \operatorname{argmin}_{s \in \{-14, -13, \dots, 14\}} \sum_{t=1}^T |\hat{R}_{t-s}^m - \hat{R}_t^{\text{consensus}}|.$$

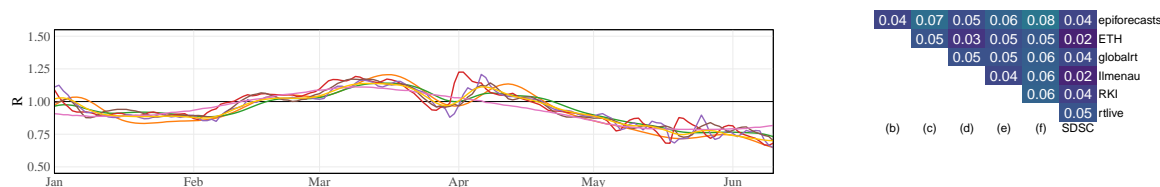
807 Here we denote by \hat{R}_t^m estimates from model m and by $\hat{R}_t^{\text{consensus}}$ estimates from the consensus model. This
 808 approach is based on [Alvarez et al. \(2021\)](#). We note that this is a pragmatic approach and does not minimize
 809 the divergence for each pair of methods nor the sum of all divergences.

810 In [Table S3](#), we compare the optimal shifts determined this way by shifts which we can compute from
 811 the employed mean incubation periods and reporting delays as provided in the respective manuscripts or
 812 code bases. Apart from the HZI model, these agree quite well.

Table S3: Shift which minimizes the mean absolute error to the consensus model and explanatory features.

Method	incubation period	reporting delay	type of R_t	shift for R_t^{case}	resulting expected shift	optimal shift (data-driven)
ETH	5.3 (3.2)	4.4 (3.4)	instant.	0	10	10
RKI	1	3.4 (0.4)	instant.	0	4	4
Ilmenau	5	2	instant.	0	7	7
SDSC		(sum to 7)	instant.	0	7	7
epiforecasts	5.4 (2.2)	5.9 (14.6)	instant.	0	11	10
rtlive	5	7.1 (5.9)	case	4	16	19
globalrt	0	0	case	4	4	3
HZI	5.2	3.7	instant.	0	9	-3
consensus		(sum to 7)	instant.	0	7	7 (fixed)

Step 4a: Shift estimates by the mean of the incubation period and reporting delay distribution.



Step 4b: Shift case reproductive number by mean generation time distribution.

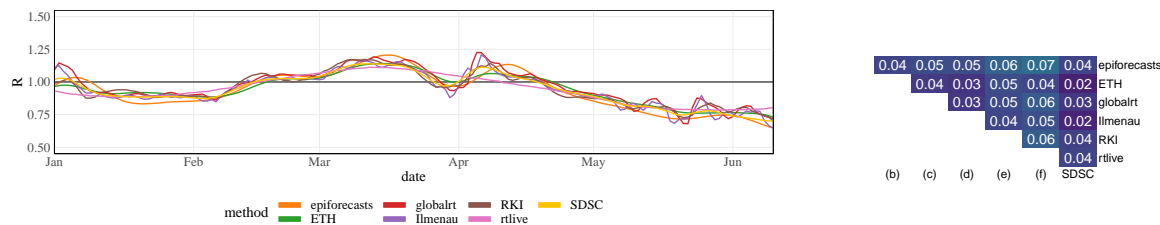


Figure S14: Alternative display on the temporal alignment of estimates as in [Figure 6](#), based on information on generation times, incubation periods, and type of R_t from [Table S3](#). We split this into two steps and omit HZI as the temporal labeling of estimates obviously does not agree with our reasoning.

S5 Supplementary Figures on temporal coherence

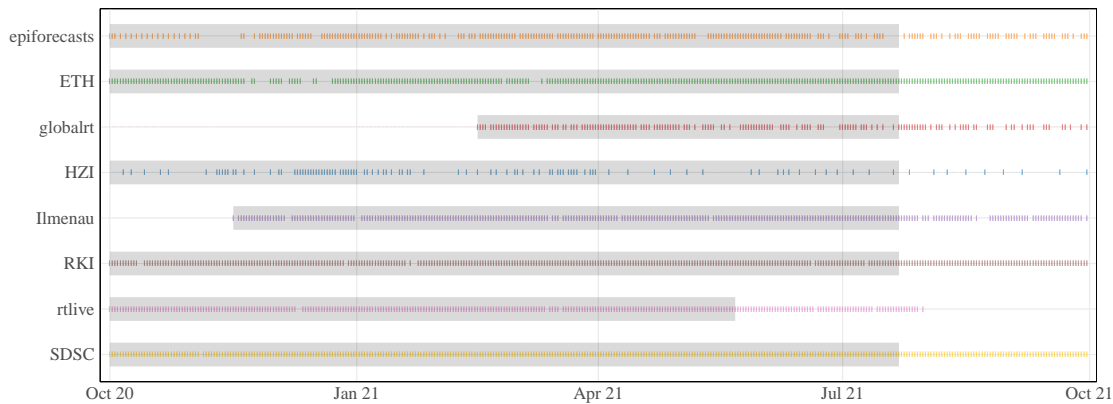


Figure S15: Dates for which estimates were published by the research groups. The shaded areas correspond to estimates which are included in the averages in Figure 5. This does not include estimation dates which are only used as consolidated estimates.

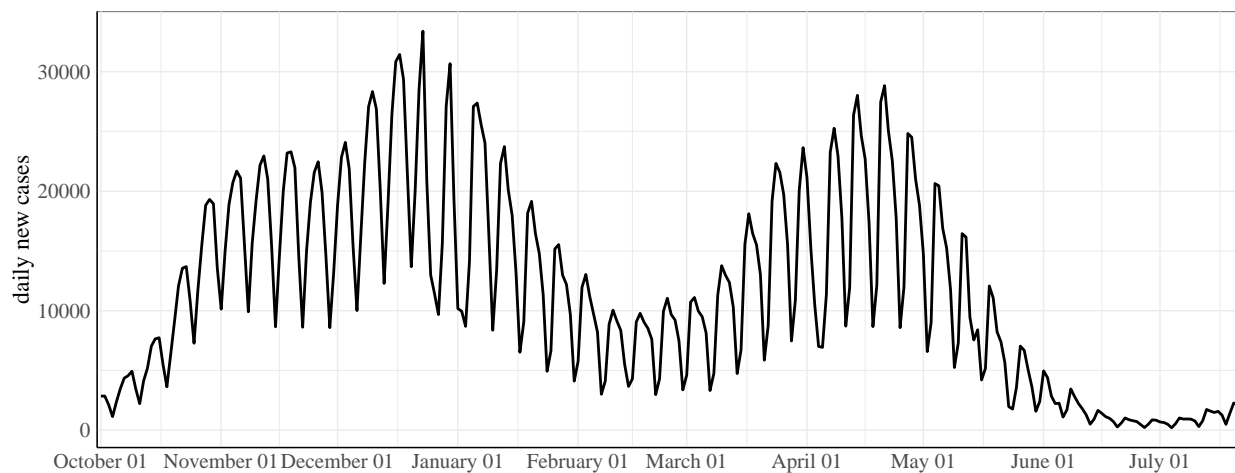


Figure S16: Incidence over the time period considered in Section 3 (RKI, positive test).

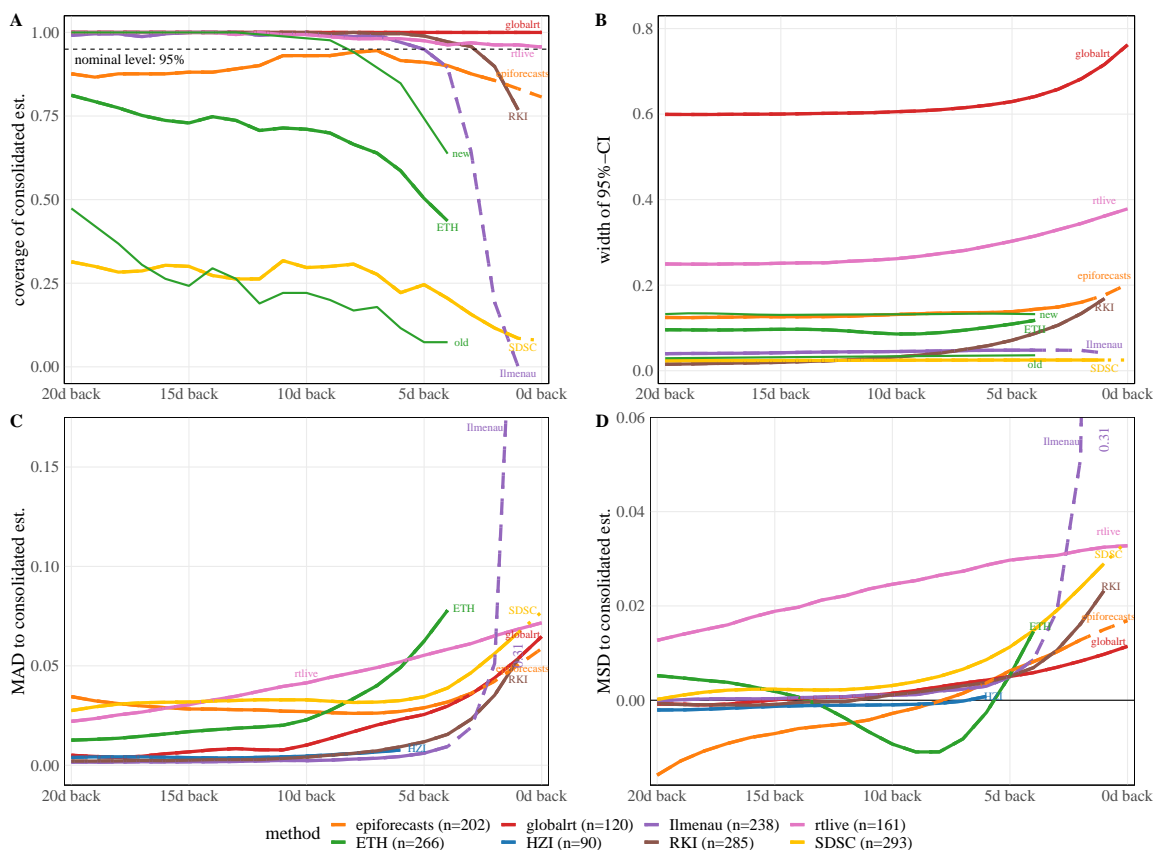


Figure S17: Same display as in Figure 5 but with curves shifted by the “optimal shift” from Section 4. Unlike in Figure 5, the horizons are thus approximately aligned, which facilitates comparison.

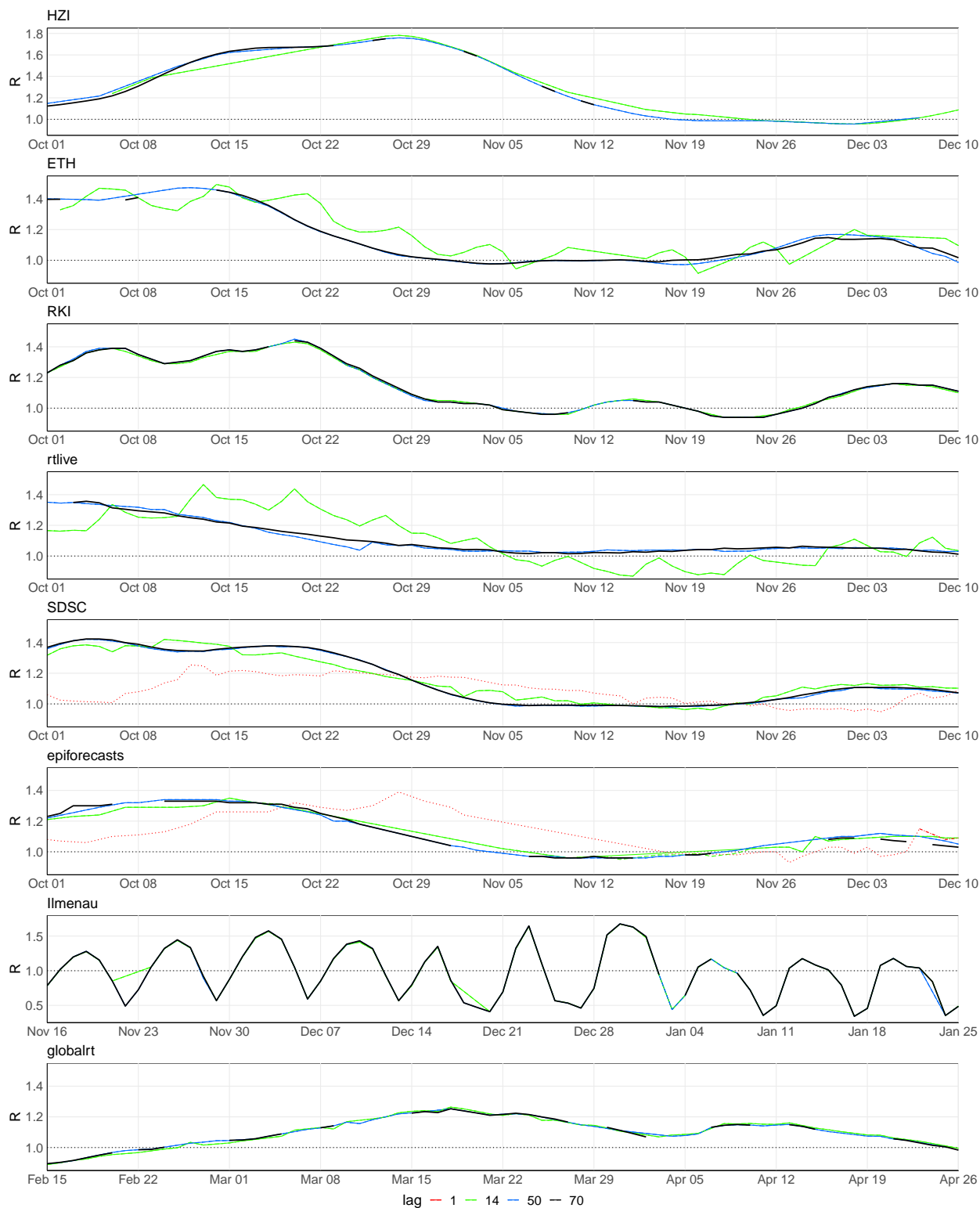


Figure S18: Real-time estimates as published 1, 14, 50, and 70 days after the target date.

814 S6 Extension of the Cori method to a conditional negative bino- 815 mial distribution

816 We provide some details on the extension of the [Cori et al. \(2013\)](#) approach to a conditional negative binomial
817 distribution as used in Section 4.3. As in equation (1) we assume that

$$\mathbb{E}(X_t \mid X_{t-1}, \dots, X_1) = \lambda_t = R_t^{\text{inst}} \times \sum_{i=1}^{t-1} w_i X_{t-i}. \quad (3)$$

In the classical [Cori et al. \(2013\)](#) approach this is combined with a conditional Poisson assumption

$$X_t \mid \lambda_t \sim \text{Pois}(\lambda_t).$$

Instead, we now use a negative binomial distribution

$$X_t \mid \lambda_t \sim \text{NegBin}(\lambda_t, \psi),$$

818 which we parameterize by its mean λ_t and an overdispersion parameter ψ , which we assume to be time-
819 constant. This implies that

$$\text{Var}(X_t \mid X_{t-1}, \dots, X_1) = \lambda_t + \psi \lambda_t^2.$$

820 This parameterization is used, e.g., in the endemic-epidemic modeling framework ([Held et al. 2007](#)) for
821 infectious disease count time series.

To fit this model to data (specifically, data from a time window of length w), we construct a covariate

$$A_t = \sum_{i=1}^{t-1} w_i X_{t-i},$$

such that

$$\mathbb{E}(X_t \mid A_t) = R_t^{\text{inst}} \times A_t.$$

822 Inference for this negative binomial generalized linear model with an identity link and no intercept can be
823 conducted using the function `glm.nb` from the R package `MASS`:

```
824 glm.nb(X ~ -1 + A, link = identity)
```

825 This estimates the parameters R_t^{inst} and ψ simultaneously using maximum likelihood inference, providing
826 confidence intervals for both parameters. The practical implementation underlying Figure 9 can be found in
827 the file [https://github.com/ElisabethBrockhaus/Rt_estimate_reconstruction/blob/main/otherFiles/epiestim_](https://github.com/ElisabethBrockhaus/Rt_estimate_reconstruction/blob/main/otherFiles/epiestim_vs_glm.R)
828 [vs_glm.R](https://github.com/ElisabethBrockhaus/Rt_estimate_reconstruction/blob/main/otherFiles/epiestim_vs_glm.R) in the GitHub repository accompanying this paper.

829 Reference:

830 Held, L., Höhle, M. and Hofmann M. (2005): A statistical framework for the analysis of multivariate
831 infectious disease surveillance counts. *Statistical Modelling*, 5: 187–199.