

CorCast: A Distributed Architecture for Bayesian Epidemic Nowcasting and its Application to District-Level SARS-CoV-2 Infection Numbers in Germany

Anna-Katharina Hildebrandt¹, Konstantin Bob², David Teschner², Thomas Kemmer², Jennifer Leclaire², Bertil Schmidt², and Andreas Hildebrandt^{2,3}

¹MONDATA GmbH, Science Park 1, Saarland University, Saarbrücken

²Institute of Computer Science, Johannes Gutenberg University, Mainz

³Corresponding author: andreas.hildebrandt@uni-mainz.de

May 2021

Timely information on current infection numbers during an epidemic is of crucial importance for decision makers in politics, medicine, and businesses. As information about local infection risk can guide public policy as well as individual behavior, such as the wearing of personal protective equipment or voluntary social distancing, statistical models providing such insights should be transparent and reproducible as well as accurate. Fulfilling these requirements is drastically complicated by the large amounts of data generated during exponential growth of infection numbers, and by the complexity of common inference pipelines. Here, we present CorCast – a stable and scalable distributed architecture for the reproducible estimation of nowcasts suitable for pandemic scenarios – and its application to the inference of district-level SARS-CoV-2 infection numbers in Germany.

1 Introduction

The SARS-CoV-2 pandemic that emerged in late 2019 has clearly shown the need for accurate, timely, and fine-grained infection statistics. To assess infection risks for different parts of the population, e.g., different age groups, the current number of daily infections for

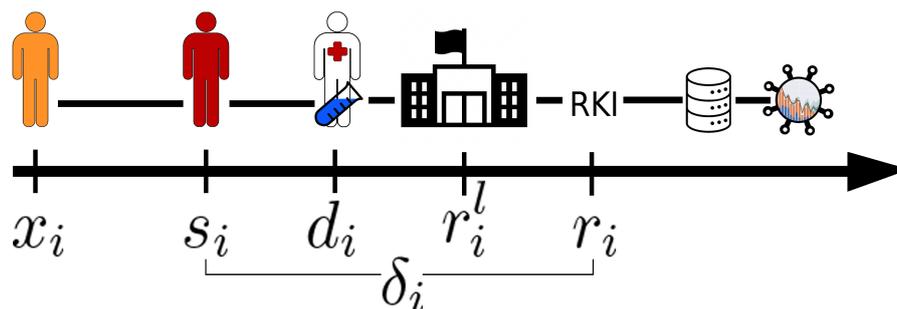


Figure 1: Timeline of dates during infection and reporting. A person gets infected at date x_i , the disease begins at date s_i , is diagnosed at date d_i and gets registered at dates r_i^l at the local health authority and r_i at the RKI, Germany's public health institute. $\delta_i := r_i - s_i$ is called the reporting delay. After registration at RKI, the case is reported in a database and can be fetched by CorCast.

28 that sub-population in the geographic region of interest is obviously crucial. Corresponding
 29 datasets also allow for the estimation of the current replication number R , and thus enable a
 30 judgement about the efficiency of current anti-epidemic measures, such as social distancing
 31 or mandatory wearing of private protective equipment such as face masks.

32 Unfortunately, the data required for these applications is typically not directly available
 33 due to a number of problems of both fundamental and practical nature, some of which are
 34 exacerbated by the particular properties of SARS-CoV-2. To illustrate these problems, let
 35 us assume that a person p_i is infected at date x_i . After an incubation period, the disease
 36 begins at date s_i , is diagnosed at date d_i , and finally registered in a database at date r_i ,
 37 together with relevant information about p_i (such as gender, age, and geographical region).
 38 The whole process is summarized in Figure 1.

39 Obviously, not all infections are identified. In fact, it is generally unknown, at least until
 40 detailed antibody studies can be performed, what the true infection rate is. This is particularly
 41 true for diseases such as SARS-CoV-2, which lead to a large percentage of asymptomatic or
 42 only mildly symptomatic cases. A further complication arises from the fact that, even for
 43 identified cases, the date of disease onset s_i is often unknown, for instance in asymptomatic
 44 cases that were identified through contact tracing or routine mass testing. In particular, it
 45 is not generally possible to assume that $s_i = d_i$, and in particular that $s_i = r_i$. Instead,
 46 each case has an individual *reporting delay* that is composed of the delay between disease
 47 onset and diagnosis as well as the delay between diagnosis and reporting, and this delay is
 48 unknown for a large percentage of the reported cases.

49 Owing to the reporting delay, information about infection numbers at the most recent days
 50 is necessarily inaccurate, or rather, incomplete. The goal of *nowcasting* of infection numbers
 51 is thus to infer an estimate of the true infection numbers based on the current, incomplete
 52 counts and the reported infection history [16]. This task is often decomposed into two
 53 phases [13]: the imputation of disease onsets for cases that have already been reported, but
 54 for which the true onset is unknown, and the nowcasting based on this data. These phases
 55 consist of a number of steps, including data ingestion, data cleanup and preprocessing,

56 imputation of disease onset, nowcasting, postprocessing, validation, and reporting.

57 **1.1 Our contribution**

58 When developing a Bayesian model for imputation and nowcasting of the SARS-CoV-2 in-
59 fection numbers in Germany at a district level, we encountered a number of challenges. For
60 instance, each step or module should be individually exchangeable, which greatly simplifies
61 development and debugging and allows for testing and validation of combinations of module
62 variants. To ensure reproducibility, it is imperative that data, modules, and workflows can all
63 be versioned in a manner that allows exact recreation of the conditions at a specific point in
64 time. Due to the large amounts of data generated in the exponential phase of a pandemic,
65 pipelines further need to scale out sufficiently. Finally, deployment and maintenance of the
66 system needs to be addressed.

67 Despite their crucial importance, these practical steps of implementation, deployment,
68 operation, and maintenance are often ignored in the literature.

69 Here, we report on a stable and scalable distributed system, called CorCast, for the repro-
70 ducible estimation of nowcasts suitable for pandemic scenarios. CorCast's implementation is
71 highly modular – indeed, retargeting CorCast to other geographical regions or other diseases
72 can be achieved by adapting only the data ingest and dashboarding modules. Similarly, the
73 influence of different preprocessing methods, different imputation models, or different now-
74 cast models on the final nowcast can be systematically studied. By basing CorCast on the
75 Pachyderm framework [22], we guarantee reproducibility. Scaling out is achieved through
76 the use of big data concepts and a scalable Kubernetes deployment.

77 To validate our architecture, we build a full pipeline for the nowcasting of Covid-19 infec-
78 tions in Germany on a district, state, and country-wide level (cf. Figure 2).

79 **1.2 Related work**

80 Related work, such as [16] followed, e.g., by [28] and [13], mainly focused on statistical
81 modeling for nowcasting, while [7] and [1] dealt also with forecasting. Software packages for
82 disease monitoring are covered, e.g., in [25], [26], and [15], but to the best of our knowledge,
83 no general framework for the development of such prediction techniques has been proposed
84 prior to this work.

85 **2 Results**

86 The main result of this work is the CorCast system as a computational pipeline. In the
87 following, we will describe its essential components and the application to Germany as a
88 proof of concept.

89 Our proposed architecture is sketched schematically in Figure 2. Inference is based on
90 probabilistic programming using Stan [5] and Turing [9]. All system services are deployed
91 on a Kubernetes cluster [2] for standardized deployment and operations, including simple
92 horizontal scaling on demand. On the cluster, a Pachyderm installation [22] orchestrates

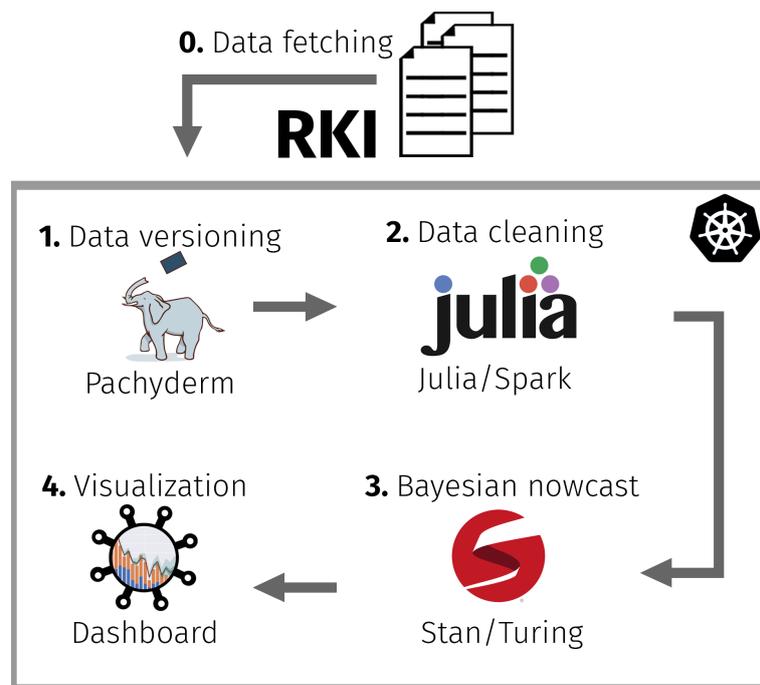


Figure 2: Graphical representation of the nowcast workflow. Data is fetched from the Robert Koch Institute, Germany's public health institute, once per day (0). All versions are saved and checked in for tracking with Pachyderm (1). Data needs to be cleaned and standardized by a set of custom scripts written in the Julia programming language (2). Afterwards, imputation and nowcast can be performed (3). Lastly, processed data and nowcast results are visualized via a custom designed dashboard (4). All CorCast components are orchestrated through Kubernetes.

93 compute pipelines and datasets. Pachyderm also keeps track of revisions of individual com-
94 pute modules and pipelines as well as of input, intermediate, and result datasets. Finally, a
95 Flask- [23] and Dash-based [24] web service, deployed on the Pachyderm cluster as a service,
96 provides a graphical interface to the ingested and processed data as well as to the imputation
97 and nowcast results.

98 Owing to the high case numbers of the pandemic, each individual module needs to be highly
99 efficient. At the same time, to facilitate extensions and improvements of the pipelines, the
100 module implementations should be intuitively readable and high-level to hide many imple-
101 mentation details required for scaling to modern compute hardware. We achieve this goal
102 through the use of the Julia language [3].

103 In the following, we will discuss the individual components of this system in more detail.

104 2.1 Data ingest and preprocessing

105 The first step in our pipeline is the collection of relevant infection data. In Germany, for
106 instance, infection data is collected locally at the health departments at the district level
107 (a political organizational unit below the state). Those departments collect the relevant
108 information and further transmit it to the Robert Koch Institute (RKI), Germany's public
109 health institute. At the RKI, collected information is processed to ensure compliance with
110 privacy regulations, and then published daily. The format of this publication is rather unusual
111 and cumbersome. In addition, the reported dates of interest (infection date, reporting date,
112 and dates of recovery or death) cannot, in general, be uniquely assigned to individual cases¹.

113 Expanding on our notation from Section 1, let us assume that an individual p_i contracts
114 the disease at date x_i and is diagnosed at date d_i . This diagnosis is reported locally, i.e.,
115 at district level, at date r_i^l and transmitted to the RKI, which receives the case information
116 at date² r_i . Each day, the RKI then publishes a spreadsheet containing information on all
117 previously received cases.

118 Extracting the required information from these data publications is a challenging task in
119 itself. For more details, c.f. Sec. 4.1.

120 2.2 Imputation of disease onset

121 The goal of this stage is to impute missing disease onset information. Let C denote the set
122 of all cases, $C \supseteq K := \{i \in C | s_i \text{ is known}\}$ the set of cases with known disease onset, and
123 $C \supset U := C \setminus K$ the set of cases with unknown onset. Our task is now to infer $s_i \forall i \in U$
124 based on the information contained in K . To this end, we build Bayesian (hierarchical)
125 models for the *reporting delay* $\delta_i := r_i - s_i$. Inference is then achieved by sampling δ_i from
126 the posterior for each case in U and using δ_i to finally predict³ $s_i = r_i - \delta_i$.

¹This is probably due to conform with privacy regulations.

²Owing to Germany's slightly archaic setup of monitoring of infectious diseases, r_i often differs from r_i^l , with a gap that can vary quite significantly from district to district, weekday to weekday, or week to week.

³In practice, we replace δ_i by $\delta_i + 1$ during inference and correct accordingly in the imputation stage. This facilitates inference, as our target distributions typically predict a probability of 0 for delays of 0. These do, however, occur in the data.

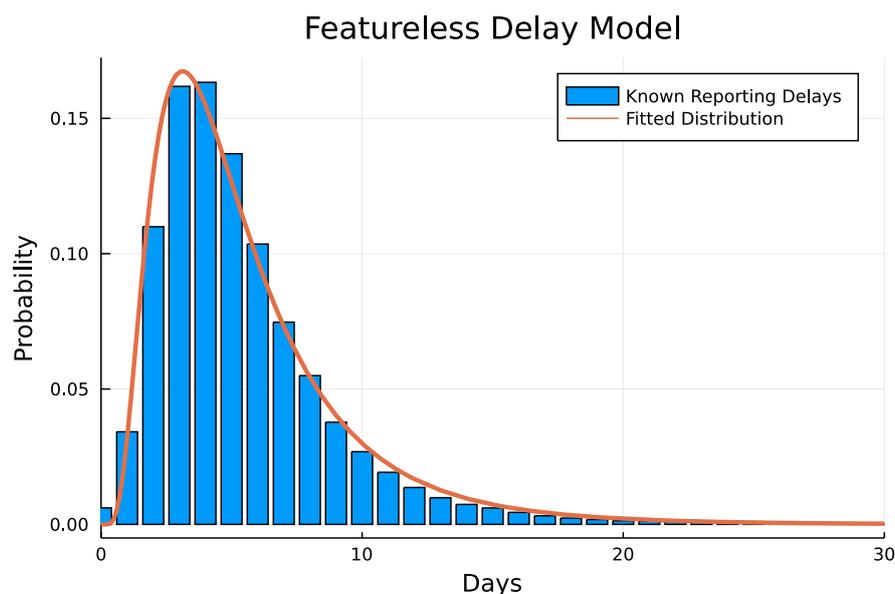


Figure 3: The delay histogram (delays shifted by 1 day to prevent values of zero) against a LogNormal distribution with parameters generated from the mean of 4 MCMC chains with 2,000 samples each, generated by the NUTS [17] sampler fed with a subset of 1,918,308 known delays.

127 In our framework, we implement several different statistical models for the delay distribu-
128 tion in two probabilistic programming frameworks (Stan [5] and Turing [9]). The baseline is
129 a trivial model that does not take any covariates of interest into account, see Section 4.3 for
130 details. However, even the baseline model fits the data remarkably well. Figure 3 shows an
131 example.

132 In a more complex model we included additional covariates, such as age group, gender,
133 or district of the cases. Since it is reasonable to assume that reporting delays might change
134 over time (e.g., owing to varying test rates or overloads of institutions during high-incidence
135 intervals), the more sophisticated model also includes a slowly varying time-dependent term,
136 e.g., modelled by a spline.

137 Additionally, the data has an obvious hierarchical component – in Germany, districts are
138 located within states. A multi-stage hierarchical model thus seems appropriate to represent
139 dependencies by partial pooling: the means and variances of the delay distribution in each
140 district should ideally be able to learn from the means and variances of each state, and those
141 of the states from those of the country as a whole.

142 Figure 4 shows a graphical representation of the model. See Section 4.3 for a detailed
143 description of the model.

144 We implemented our models in both Stan [5] (using CmdStan.jl [6]) and Turing.jl [9]
145 and draw samples from the posterior using the NUTS [17] sampler to determine probability
146 distributions for the model parameters. From these distributions, we can finally impute

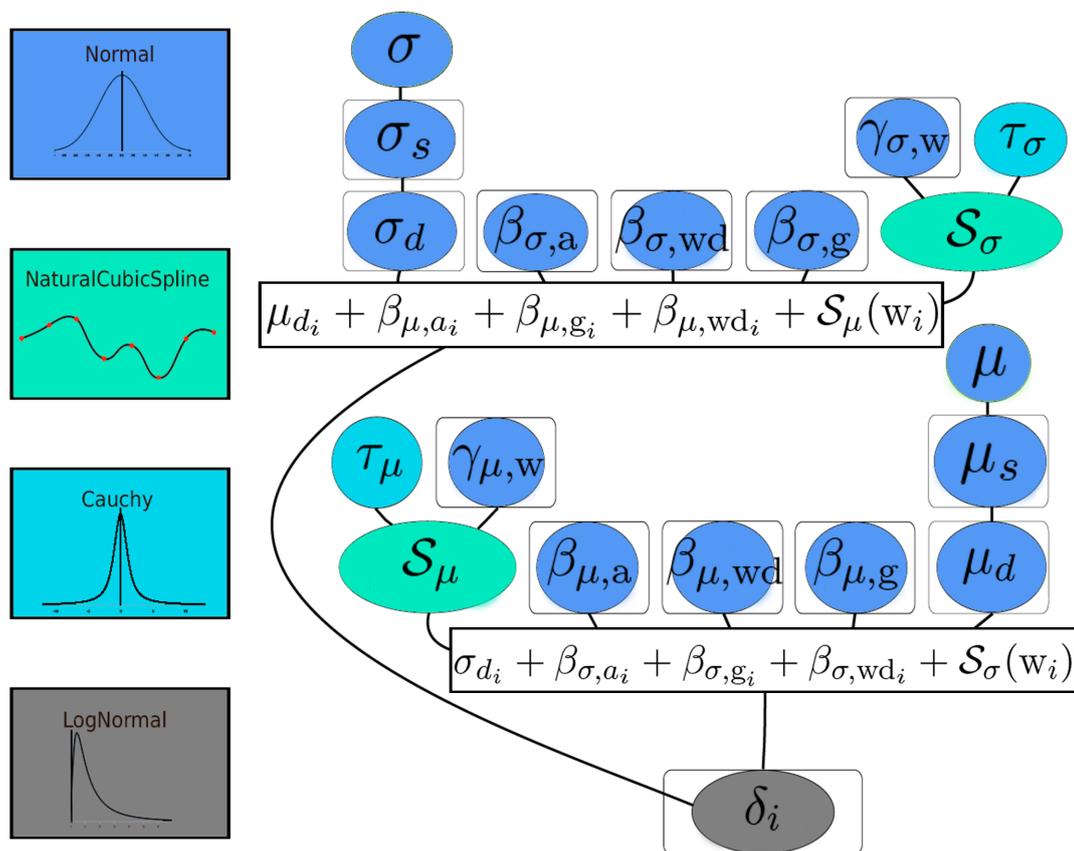


Figure 4: Imputation model with covariates age group, gender, district, weekday, and number of weeks since the start of the pandemic. Ellipses denote random variables, ellipses inside boxes indicate multidimensional variables, the terms in boxes indicate dependent random variables and lines indicate dependencies. The prior distributions used are encoded by color. Note the hierarchical structure of the spatial covariate. Due to computational constraints the model was evaluated in several steps. Details are given in Section 4.3.

147 infection dates by drawing – for each case with unknown disease onset – delay values from
148 the resulting posterior predictive distribution, conditioned on the appropriate covariates (e.g.,
149 the age group, gender, or district associated with the case).

150 Here it is worth noting that the delay distributions (cf. Figure 3) tend to have a rather
151 pronounced peak, typically close to 7 days. If we would sample delay values multiple times for
152 each case with unknown disease onset to average over the posterior predictive distribution,
153 we would thus almost always assign delays of approx. 7 days. Instead, delay values are
154 drawn once per case and used to compute disease onsets. We then group by disease onset
155 and aggregate cases. Repeating this process multiple times allows to compute confidence
156 intervals for the number of new infections per day.

157 **2.3 Nowcasting infection numbers**

158 During an epidemic, infection numbers cannot be observed in real time. Due to the reporting
159 delay between the date x_i of disease onset and r_i of receiving the case information at the
160 central case registry (the RKI in Germany), information about current case numbers will
161 necessarily be incomplete. If we assume that the reporting delay δ is bounded by a reasonable
162 maximum D , we can expect, at each day during the epidemic, to be notified about additional
163 cases with disease onset in the last D days (we assume that reporting delays are generally
164 > 1).

165 We adopted a model from [13], which in turn extends earlier work by [16]. The details are
166 given in section 4.4.

167 To provide nowcasts at the different levels of the hierarchy (district, state, and country),
168 we group the data by district and run NUTS-sampling on each of the resulting 412 datasets
169 (one for each district), followed by one run for each of the 16 states and one for the whole of
170 the country. This process is quite time consuming and takes roughly 5 hours on three AMD
171 EPYC™ 2nd Gen. worker nodes with 8 cores and 32 GiB RAM each.

172 **2.4 Evaluation of nowcast**

173 The evaluation of nowcast results is not a trivial task. In Section 4.5, we summarize the
174 relevant literature. CorCast contains an evaluation module that implements the WIS-score
175 described in that section.

176 **2.5 Dashboard functionality**

177 CorCast provides a dashboard feature through Pachyderm's service pipeline mechanism. The
178 dashboard is implemented in Flask [23] and Plotly's dash framework [24], and can be easily
179 extended and adapted. The dashboard for the CorCast model of district-level imputation and
180 nowcasting in Germany can be found at <https://dashboard.covid19.hildebrandtlab.org/>
181 org/ (cf. Figure 5).

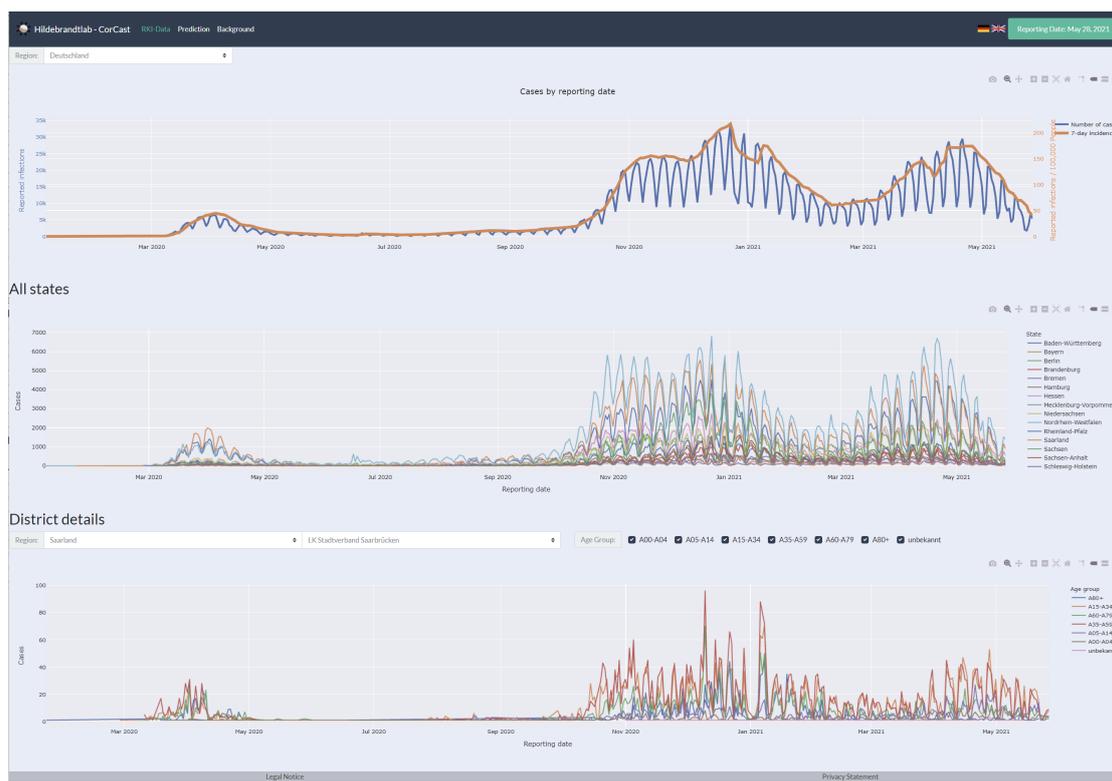


Figure 5: The CorCast dashboard.

182 **3 Discussion**

183 Several advantages distinguish the CorCast system described in this work. Here, we want
184 to summarize the most important of these, namely (i) reproducibility (ii) modularity and
185 extensibility (iii) scalability and (iv) standardized deployment and operations.

186 **3.1 Reproducibility**

187 Reproducibility is a cornerstone of science that can be very challenging to achieve. In com-
188 plex computational systems, it is typically insufficient to merely guarantee accessibility to
189 input data sets and informal descriptions of the employed programs and workflows. Modern
190 scientific computer programs often rely on a large number of third-party dependencies, such
191 as support libraries. Even small variations in the computational infrastructure can lead to
192 differing outcomes, either through the introduction of bugs [4] or through numerical effects
193 that lead, e.g., to convergence to different minima. True reproducibility hence requires reli-
194 able versioning of data, programs, and workflows. CorCast achieves this goal by relying on
195 three key technologies: Git, for versioning source code, Docker containers [20] for versioning
196 applications, and Pachyderm [22] for versioning data and workflows.

197 A side effect of incorporating reproducibility tightly into the system design is a greatly
198 simplified evaluation and comparison of different steps in the nowcasting pipeline. Imagine,
199 for instance, that we want to compare a novel Bayesian model for imputation of disease
200 onset with respect to its influence on the final nowcast. Through the versioning of data
201 and workflows, we can recreate the exact information used to perform the nowcasts at any
202 former date since the initial data commit of the CorCast system, exchange the imputation
203 module through our new candidate, run the full nowcasting pipeline, and finally evaluate the
204 performance and compare it to the historic performance of our former model, which we can
205 again query from the corresponding data repository.

206 **3.2 Modularity and Extensibility**

207 CorCast has been designed as a highly modular pipeline, with components responsible for
208 ingesting data, processing and cleaning of the raw inputs, imputation of disease onsets,
209 nowcasting, evaluation, and interactive visualization. Each of these components can be easily
210 replaced, but it is also a simple task to add steps to the pipeline, or to provide alternative
211 implementations for each step in the pipeline (such as competing nowcasting models) and
212 to run them in parallel. To give an example of the implications of this extensibility, we want
213 to stress that the current implementation of CorCast is focused on nowcasting SARS-CoV-2
214 infection numbers in Germany, but that it would be a simple matter of extending or replacing
215 the ingest nodes of the pipeline to extend it to other countries, or to adapt it for different
216 infectious diseases.

217 **3.3 Scalability**

218 While it might not be obvious at first glance, the data generated during a pandemic grows
219 very rapidly. Merely reporting summary statistics, i.e., the number of new infections per

220 day, is typically insufficient, as the progression of the epidemic can vary significantly in
221 different sub-populations (differentiated, e.g., by age, gender, or location). To account
222 for these differences, all differentiating features (including the date of disease onset) will
223 have to be reported. In practice, much of the reporting hence happens at the level of
224 individual cases, which grows rapidly during the pandemic phases of exponential growth.
225 Further complications arise through inconsistencies and retrospective modifications that have
226 to be treated carefully in the data preparation stage. In the case of German infection
227 numbers as provided by the Robert Koch Institute, this requires iterating over all previous
228 case reports, where each case report contains information about every single case since the
229 start of the pandemic, leading to quadratic effort. At the time of writing, the uncompressed
230 data size for this correction step is approximately 43.4 GiB. The large case numbers also
231 render Bayesian inference rather expensive. At the time of writing, imputation of disease
232 onset takes roughly 5 hours, nowcasting roughly 4 hours, on a Kubernetes cluster with
233 one dedicated master and three AMD EPYC™ 2nd Gen. worker nodes with 8 cores and
234 32GiB RAM each. To handle the large amounts of data and the computational effort,
235 we have taken care to expose opportunities for parallelism, e.g., by separating inference
236 models by geographic region, by running multiple shorter MCMC chains in parallel, or by
237 vectorization- and parallelization primitives inside the computational nodes. Scheduling of
238 parallel computation is achieved through Pachyderm's parallelization primitives. Without
239 these measures, it would be essentially infeasible to produce daily nowcast results at our
240 level of granularity.

241 **3.4 Standardized Deployment and Operations**

242 CorCast uses DevOps principles to facilitate deployment and operations. Following an infra-
243 structure-as-code (IaC) approach, all components are provisioned and managed using Ter-
244 raform [14] and Helm [8] files. All compute nodes are realized as Docker containers and all
245 workflows encoded as Pachyderm pipelines in JSON format. Orchestration and managing
246 of system resources (monitoring, up- and down-scaling, upgrade handling, etc.) are realized
247 through Kubernetes and Prometheus. As a result of the IaC approach, all information re-
248 quired to set up a CorCast cluster is under version control in the CorCast Git repositories.
249 Scaling the existing instance or setting up a new instance on a variety of cloud providers
250 or on premise is trivial. Hence, CorCast can be easily adapted and deployed by researchers
251 interested in epidemics now- or forecasting.

252 **3.5 Summary**

253 CorCast is a highly flexible and scalable framework for the implementation of computational
254 statistics and machine learning approaches to epidemiology that provides automatic repro-
255 ducibility. In addition, the CorCast instance described in this work is, to the best of our
256 knowledge, the only currently available implementation for nowcasting of German Sars-CoV-
257 2 infection numbers that can provide daily updates for disease onset imputation and nowcasts
258 on a district level. In future work, we intend to add additional models to the CorCast system,
259 such as the stable inference of replication numbers, and want to develop pipelines for other

260 geographic regions.

261 4 Methods

262 4.1 Preprocessing

263 The RKI provides so-called “publications” that contain sufficient information to query the
264 number of infections, recoveries, and deaths as a function of r_i^l , i.e., the date at which
265 they were recorded locally at the health department of the district. Unfortunately, the
266 files do not contain the date at which they were recorded centrally by the RKI, r_i . Since
267 much of the information provided by the RKI, such as imputed disease onsets and nowcasts,
268 seems to be predicted based on r_i rather than on r_i^l , the lack of this information is rather
269 unfortunate. To recreate r_i , it is necessary to keep an archive of all previous daily data
270 publications of the RKI (a task for which Pachyderm is ideally suited), and iterate backwards
271 to identify at what date a case first appeared in the RKI publications. Since the data
272 publications also contain corrections, which are modelled as deletions of cases in previous
273 publications and new insertions of these cases, recreating r_i is non-trivial and error-prone.
274 In addition, the publications contain numerous inconsistencies and errors that have to be
275 identified and corrected for as accurately as possible. Finally, the process is made considerably
276 more cumbersome by the fact that the data format varies over time. Dates, in particular, are
277 stored in at least four different formats, depending on the date of the publication and the
278 data column. Hence, parsing and preprocessing the input data is relatively complex. At the
279 same time, the preprocessing needs to be highly efficient: individual cases have to be tracked
280 through hundreds of files, each of which needs several steps of data conversion, processing,
281 and quality control. Coupled with the large case numbers of the pandemic, preprocessing
282 becomes quite resource intensive.

283 The ingestion pipelines at the start of our Pachyderm workflows start with a module that
284 downloads the current RKI data publication and archives it in a Pachyderm repository. In
285 the next step, the newest publication is processed and converted into a dataframe, which is
286 then serialized to another Pachyderm repository. Simultaneously, we attempt to reconstruct
287 the missing r_i values through iteration over all previous data publications⁴. The results are
288 again stored in a Pachyderm repository. This approach allows us to base further processing
289 steps on either the RKI-provided r_i^l or on the missing r_i values.

290 4.2 Bayesian hierarchical modeling

291 Bayesian modeling is statistical modeling that is centered around Bayes' theorem⁵: One starts
292 with formulating possible hypotheses and assigns probabilities (prior probabilities $p(H)$) to
293 them before comparison to the data D and then computes the conditional probabilities of
294 each hypothesis given the data (posterior probabilities $p(H|D)$).

⁴The details of the implementation are rather complex to handle a variety of special cases.

⁵and on a more philosophical level, a different interpretation of probability. [19] provides further insights into the foundations and philosophical implications of Bayesian thinking.

295 This approach has the advantage that already available information can be leveraged in a
296 transparent way and the uncertainty of estimated quantities can be quantified as well. Fur-
297 thermore, the robustness of predictions can be greatly improved by using posterior weighted
298 predictions.

299 When dealing with data that arises in systems with a hierarchical nature, i.e., follow some
300 tree-like structure, a common technique is to split the prior probability distributions into a
301 product that resembles this structure, giving rise to hierarchical models [10].

302 Exact inference of Bayesian hierarchical models is hindered by fact that the integral
303 $\int dH p(D|H)p(H)$, where $p(D|H)$ denotes the likelihood, is typically intractable by symbolic
304 expressions. Thus, numerical methods are used in practice, which are mostly sampling-based.
305 These Markov chain Monte Carlo (MCMC) methods generate samples that are distributed
306 according to the posterior and thus allow to compute approximations to quantities of interest
307 on that sample.

308 While the classical Metropolis-Hastings algorithm [21] provides asymptotically correct sam-
309 ples, more advanced algorithms like [17] from the family of Hamiltonian Monte Carlo algo-
310 rithms yield a better performance in terms of required sample size and correlation of samples.

311 As Hamiltonian Monte Carlo requires the computation of gradients, several software pack-
312 ages, e.g., [5], [9], and [27], were developed to leverage automatic differentiation of model
313 expressions.

314 4.3 Imputation models

315 Our baseline model is defined as

$$\mu \sim \text{Normal}_{[0.1,10]}(1, 2) \quad (1)$$

$$\sigma \sim \text{Normal}_{[0.1,10]}(1, 2) \quad (2)$$

$$\delta \sim \text{LogNormal}(\mu, \sigma), \quad (3)$$

316 where $\text{Normal}_{[l,u]}$ denotes a truncated normal distribution and the priors for μ and σ are
317 only weakly informative and positive.

318 The fully hierarchical model shown in Figure 4 turned out to be infeasibly costly to compute
319 on the large amounts of data generated during the pandemic. As a less sophisticated variant,
320 we use a manual approximation to the hierarchical model by first fitting a global model for
321 the delay distribution to all cases in Germany, then a model with means and variances per
322 state of the form

$$\mu_s \sim \text{Normal}_{[0.1,10]}(\mu_g, 0.5), \quad 1 \leq s \leq 16 \quad (4)$$

$$\sigma_s \sim \text{Normal}_{[0.1,10]}(\sigma_g, 0.5), \quad 1 \leq s \leq 16 \quad (5)$$

$$\delta_i \sim \text{LogNormal}(\mu_s, \sigma_s), \quad s = \text{state}(\text{case}(i)), 1 \leq i \leq \#cases \quad (6)$$

323 and, finally, fitting a more complex model with additional covariates (here: age group,
324 gender, district, weekday, and number of weeks since the start of the pandemic):

$$\mu_d \sim \text{Normal}_{[0.1,10]}(\mu_s, 0.5), \quad s = \text{state}(\text{district}(d)), d \in \text{districts} \quad (7)$$

$$\sigma_d \sim \text{Normal}_{[0.1,10]}(\sigma_s, 0.5), \quad s = \text{state}(\text{district}(d)), d \in \text{districts} \quad (8)$$

$$\beta_{\mu,a} \sim \mathcal{N}(0, 0.5), \quad a \in \text{age groups} \quad (9)$$

$$\beta_{\sigma,a} \sim \mathcal{N}(0, 0.5), \quad a \in \text{age groups} \quad (10)$$

$$\beta_{\mu,g} \sim \mathcal{N}(0, 0.5), \quad g \in [\text{male, female, unknown}] \quad (11)$$

$$\beta_{\sigma,g} \sim \mathcal{N}(0, 0.5), \quad g \in [\text{male, female, unknown}] \quad (12)$$

$$\beta_{\mu,wd} \sim \mathcal{N}(0, 0.5), \quad wd \in \text{weekdays} \quad (13)$$

$$\beta_{\sigma,wd} \sim \mathcal{N}(0, 0.5), \quad wd \in \text{weekdays} \quad (14)$$

$$\tau_\mu \sim \text{Cauchy}(0, 1) \quad (15)$$

$$\tau_\sigma \sim \text{Cauchy}(0, 1) \quad (16)$$

$$\gamma_{\mu,w} \sim \mathcal{N}(0, 1), \quad w \in \text{week nodes} \quad (17)$$

$$\gamma_{\sigma,w} \sim \mathcal{N}(0, 1), \quad w \in \text{week nodes} \quad (18)$$

$$\mathcal{S}_\mu = \text{NaturalCubicSpline}(\text{week nodes}, \tau_\mu * \gamma_{\mu,w}) \quad (19)$$

$$\mathcal{S}_\sigma = \text{NaturalCubicSpline}(\text{week nodes}, \tau_\sigma * \gamma_{\sigma,w}) \quad (20)$$

$$\mu_i = \mu_{d_i} + \beta_{\mu,a_i} + \beta_{\mu,g_i} + \beta_{\mu,wd_i} + \mathcal{S}_\mu(w_i) \quad (21)$$

$$\sigma_i = \sigma_{d_i} + \beta_{\sigma,a_i} + \beta_{\sigma,g_i} + \beta_{\sigma,wd_i} + \mathcal{S}_\sigma(w_i) \quad (22)$$

$$\delta_i \sim \text{LogNormal}(\mu_i, \sigma_i) \quad (23)$$

325 4.4 Nowcasting model

326 In this work, we use a nowcasting model based on [13], which in turn extends earlier work by
 327 [16]. These approaches use a hierarchical Bayesian model composed of two parts: a model
 328 for the delay distribution δ as a function of time ($P(\delta_i = d | x_i = t) := p_{t,d}$) and a model for
 329 the expected number of infections per day ($E[N(t, \infty)] := \lambda_t$).

330 For the first part, we can either plug in the delay distribution obtained during imputation
 331 (cf. Section 2.2) or fit a new model, such as the one proposed in [13], which starts from a
 332 discrete time hazard model with $h_{t,d} = P(\delta = d | \delta \geq d, W_{t,d})$ with time- and delay-dependent
 333 covariates $W_{t,d}$ and with $\text{logit}(h_{t,d}) = \gamma_d + W'_{t,d}\eta$ for $d = 0, \dots, D-1$, with bias γ_d and with
 334 $h_{t,D} = 1$. In this approach, the covariates $W_{t,d}$ can represent different features of interest.

335 Like [13], we use a first-order spline to include a general time dependence and factor variables
336 to represent weekdays and holidays. The probabilities $p_{t,d}$ can be obtained from the hazard
337 model through $p_{t,0} = h_{t,0}$ and $p_{t,d} = \left(1 - \sum_{d=0}^{D-1} p_{t,d}\right) \times h_{t,d}$.

338 For the second part of the hierarchical nowcasting model, we again follow [13] and use the
339 Gaussian random walk

$$\lambda_0 \sim \text{LogNormal}(0, 1) \quad (24)$$

$$\lambda_t \sim \text{LogNormal}(\log(\lambda_{t-1}), \sigma^2), \quad t = 1, \dots, T \quad (25)$$

$$n_{t,d} | \lambda_t, p_{t,d} \sim \text{NegativeBinomial}(\lambda_t \times p_{t,d}, \phi), \quad t = 1, \dots, T \quad (26)$$

340 where ϕ is the overdispersion parameter of the negative binomial distribution. Appropriate
341 priors for the model are again taken from [13] and the implementation referenced therein.

342 We finally obtain $N(t, \infty) = \sum_{d=0}^D n_{t,d}$.

343 The number of parameters of this model grows with $T \times D$, where T is the day at which
344 the nowcast is to be evaluated. Hence, the number of variables grows linearly with each day
345 of the pandemic, rendering sampling increasingly time- and memory-consuming. To prevent
346 this growth from turning nowcasting infeasible, we limit the amount of history we track in
347 the Gaussian random walk, i.e., instead of setting $t = 0$ to the first day of the Covid-19
348 pandemic, we choose a fixed time interval of size $\Delta > D$ (set to the last $\Delta = 50$ days in
349 our implementation). Since $\lambda_0 = E[N(0, \infty)]$ equals the number of infections happening at
350 day $t = 0$, we need to adapt Eq. (24) accordingly.

351 4.5 Evaluating nowcasts

352 Evaluating probabilistic predictions requires different scoring functions than point estimates
353 do. For a comprehensive discussion of the topic, we refer the interested reader to [18]. Here,
354 we briefly summarize the main evaluation measures used in this work.

Evaluation of now- and forecasts is relatively straightforward if the computational model yields the full predictive probability distribution. Given a test datum y , common measures are the logarithmic score [12]

$$\log S(y) = \log(p(y)), \quad (27)$$

where $p(\cdot)$ denotes the probability density function used for prediction, or the continuous ranked probability score [11]

$$\text{CRPS}(F, y) = \int_{-\infty}^{\infty} dx \{F(x) - 1_{x \geq y}\}^2, \quad (28)$$

355 where $F(\cdot)$ denotes the cumulative distribution function used for prediction. In [13] the
356 coverage frequency of the 95% prediction interval is used as well.

357 In many practical situations, the full predictive distribution is not available. In those cases,
358 we typically have access to the predictive mean or median and, optionally, several quantiles
359 of the cumulative distribution function. Assuming that we are given the predictive median
360 m and a collection of K central prediction intervals (PIs) at levels $(1 - \alpha_1) < (1 - \alpha_2) <$

361 $\dots < (1 - \alpha_K)$, and weights $w_k, k = 1 \dots K$, the weighted interval score (WIS) is defined
362 as

$$\text{WIS}_{\alpha_{\{0:K\}}} := \frac{1}{K + \frac{1}{2}} \times \left(w_0 \times |y - m| + \sum_{k=1}^K \{w_k \times \text{IS}_{\alpha_k}(\text{F}, y)\} \right) \quad (29)$$

363 where the interval score IS_α for the central $(1 - \alpha) \times 100\%$ PI is defined as

$$\text{IS}_\alpha(\text{F}, y) := (u_F - l_F) + \frac{2}{\alpha} \times (l_F - y) \times 1_{y \leq l_F} + \frac{2}{\alpha} \times (y - u_F) \times 1_{y \geq u_F}$$

364 and where l_F and u_F denote the $\frac{\alpha}{2}$ -quantile and $(1 - \frac{\alpha}{2})$ -quantile of F , respectively.
365 Following [18], we use $w_0 = \frac{1}{2}$ and $w_k = \frac{\alpha_k}{2}$ for $k = 1 \dots K$, as this choice approximates
366 the CRPS score for large values of K and equidistant α_k .

367 Authors Contributions

368 AKH and AH conceptualized the work. AKH, KB, DT, TK, JL and AH discussed and
369 implemented the work, BS contributed work on parallelization and optimization. All authors
370 contributed to the conception, writing and editing of this manuscript.

371 Competing Interests

372 The authors declare that there is no conflict of interest.

373 Data and code availability

374 All source code for the deployment and operation of the CorCast system and for the district-
375 level models for German infection numbers will be made available upon publication at <https://github.com/hildebrandtlab/Covid19Nowcast.jl>.
376

377 References

- 378 [1] Vinicius VL Albani, Roberto M Velho, and Jorge P Zubelli. “Estimating, monitor-
379 ing, and forecasting COVID-19 epidemics: a spatiotemporal approach applied to NYC
380 data”. In: *Scientific Reports* 11.1 (2021), pp. 1–15.
- 381 [2] The Kubernetes Authors. *Kubernetes* (<https://kubernetes.io>).
- 382 [3] Jeff Bezanson et al. “Julia: A fresh approach to numerical computing”. In: *SIAM*
383 *Review* 59.1 (2017), pp. 65–98.
- 384 [4] Jayanti Bhandari Neupane et al. “Characterization of Leptazolines A–D, Polar Oxazo-
385 lines from the Cyanobacterium *Leptolyngbya* sp., Reveals a Glitch with the “Willoughby-
386 Hoye” Scripts for Calculating NMR Chemical Shifts”. In: *Organic Letters* 21.20 (2019),
387 pp. 8449–8453.
- 388 [5] Bob Carpenter et al. “Stan: A probabilistic programming language”. In: *Journal of*
389 *statistical software* 76.1 (2017).

- 390 [6] *CmdStan.jl* (<https://github.com/StanJulia/CmdStan.jl>).
- 391 [7] Jonas Dehning et al. “Inferring change points in the spread of COVID-19 reveals the
392 effectiveness of interventions”. In: *Science* 369.6500 (2020).
- 393 [8] Cloud Native Computing Foundation. *Helm* (<https://www.helm.sh/>).
- 394 [9] Hong Ge, Kai Xu, and Zoubin Ghahramani. “Turing: a language for flexible proba-
395 bilistic inference”. In: *International Conference on Artificial Intelligence and Statistics,
396 AISTATS 2018, 9-11 April 2018, Playa Blanca, Lanzarote, Canary Islands, Spain*. 2018,
397 pp. 1682–1690.
- 398 [10] A. Gelman et al. *Bayesian Data Analysis, Third Edition*. Chapman & Hall/CRC Texts
399 in Statistical Science. Taylor & Francis, 2013.
- 400 [11] Tilmann Gneiting, Fadoua Balabdaoui, and Adrian E. Raftery. “Probabilistic forecasts,
401 calibration and sharpness”. In: *Journal of the Royal Statistical Society: Series B (Sta-
402 tistical Methodology)* 69.2 (2007), pp. 243–268.
- 403 [12] Tilmann Gneiting and Adrian E Raftery. “Strictly Proper Scoring Rules, Prediction,
404 and Estimation”. In: *Journal of the American Statistical Association* 102.477 (2007),
405 pp. 359–378.
- 406 [13] Felix Günther et al. “Nowcasting the COVID-19 pandemic in Bavaria”. In: *Biometrical
407 Journal* 63.3 (2020), pp. 490–502.
- 408 [14] HashiCorps. *Terraform* (<https://www.terraform.io/>).
- 409 [15] Michael Hohle. “Surveillance: An R package for the monitoring of infectious diseases”.
410 In: *Computational Statistics* 22.4 (2007), pp. 571–582.
- 411 [16] Michael Höhle and Matthias an der Heiden. “Bayesian nowcasting during the STEC
412 O104:H4 outbreak in Germany, 2011”. In: *Biometrics* 70.4 (2014), pp. 993–1002.
- 413 [17] Matthew D. Homan and Andrew Gelman. “The No-U-turn Sampler: Adaptively Setting
414 Path Lengths in Hamiltonian Monte Carlo”. In: *J. Mach. Learn. Res.* 15.1 (Jan. 2014),
415 pp. 1593–1623.
- 416 [18] Bracher J et al. “Evaluating epidemic forecasts in an interval format”. In: *PLoS Comput
417 Biol* 17.2 (2021).
- 418 [19] E. T. Jaynes. *Probability Theory - The Logic of Science*. Cambridge: Cambridge Uni-
419 versity Press, 2003.
- 420 [20] Dirk Merkel. “Docker: Lightweight Linux Containers for Consistent Development and
421 Deployment”. In: *Linux J.* 2014.239 (Mar. 2014).
- 422 [21] Nicholas Metropolis et al. “Equation of State Calculations by Fast Computing Ma-
423 chines”. In: *The Journal of Chemical Physics* 21.6 (1953), pp. 1087–1092.
- 424 [22] Jon Ander Novella et al. “Container-based bioinformatics with Pachyderm”. In: *Bioin-
425 formatics* 35.5 (2018), pp. 839–846.
- 426 [23] Pallets. *Flask* (<https://flask.palletsprojects.com/en/2.0.x/>).
- 427 [24] Plotly. *Dash* (<https://dash.plotly.com/>).

- 428 [25] Maëlle Salmon, Dirk Schumacher, and Michael Höhle. “Monitoring Count Time Series
429 in R: Aberration Detection in Public Health Surveillance”. In: *Journal of Statistical*
430 *Software, Articles* 70.10 (2016), pp. 1–35.
- 431 [26] Maëlle Salmon et al. “A system for automated outbreak detection of communicable
432 diseases in Germany”. In: *Eurosurveillance* 21.13, 30180 (2016).
- 433 [27] John Salvatier, Thomas V. Wiecki, and Christopher Fonnesbeck. “Probabilistic pro-
434 gramming in Python using PyMC3”. In: *PeerJ Computer Science* 2 (Apr. 2016), e55.
- 435 [28] Olivera Stojanović et al. “A Bayesian Monte Carlo approach for predicting the spread
436 of infectious diseases”. In: *PloS one* 14.12 (2019), e0225838.