

# 1 A citizen science initiative for open data and 2 visualization of COVID-19 outbreak in Kerala, India

## 3 Collective for Open Data Distribution-Keralam (CODD-K)

### 4 CODD-K authors list

5 Jijo Pulickiyil Ulahannan <sup>1#\*</sup>, Nikhil Narayanan <sup>2#</sup>, Nishad Thalhath <sup>3#</sup>, Prem  
6 Prabhakaran <sup>4</sup>, Sreekanth Chaliyeduth <sup>5</sup>, Sooraj P Suresh <sup>6</sup>, Musfir Mohammed <sup>7</sup>,  
7 Rajeevan E <sup>8</sup>, Sindhu Joseph <sup>9</sup>, Akhil Balakrishnan <sup>10</sup>, Jeevan Uthaman <sup>11</sup>, Manoj  
8 Karingamadathil <sup>12</sup>, Sunil Thonikkuzhiyil Thomas <sup>13</sup>, Unnikrishnan Sureshkumar <sup>14</sup>,  
9 Shabeesh Balan <sup>15</sup>, Neetha Nanoth Vellichirammal <sup>16</sup>

10 <sup>1</sup> Department of Physics, Government College Kasaragod, Kannur University, Kerala, India, <sup>2</sup>  
11 Open Data Researcher, Bengaluru, India, <sup>3</sup> School of Library, Information and Media Studies.,  
12 University of Tsukuba, Japan <sup>4</sup>Department of Advanced Materials and Chemical Engineering,  
13 Hannam University, Daejeon, South Korea <sup>5</sup> Centre for Cognitive and Brain Sciences, Indian  
14 Institute of Technology Gandhinagar, Gujarat, India, <sup>6</sup> Department of Humanities and Social  
15 Sciences, National Institute of Technology Tiruchirappalli, Tamil Nadu, India, <sup>7</sup> Embedded  
16 Analytics, ML and Data Sciences, Experion Technologies, TechnoPark, Thiruvananthapuram,  
17 India, <sup>8</sup> Department of Philosophy, Government Brennen College, Kannur University, Kerala,  
18 India, <sup>9</sup> Department of Travel and Tourism Management, Govinda Pai Memorial Government  
19 College, Kannur University, Kerala, India, <sup>10</sup> Crowcon - A Halma Company, ITPB, Whitefield,  
20 Bangalore, <sup>11</sup> Department of Marine Geophysics, Cochin University of Science and  
21 Technology, Kochi, Kerala, India, <sup>12</sup> Swathantha Malayalam Computing, Thrissur, Kerala,  
22 India, <sup>13</sup> Department of Electronics, College of Engineering Attingal, APJ Abdul Kalam  
23 Technical University, Thiruvananthapuram, Kerala, India, <sup>14</sup> Astronomical Observatory of the  
24 Jagiellonian University, Kraków, Małopolska, Poland, <sup>15</sup> Laboratory for Molecular Psychiatry,  
25 RIKEN Center for Brain Science, Wakoshi, Saitama, Japan, <sup>16</sup> Department of Genetics, Cell  
26 Biology and Anatomy, University of Nebraska Medical Center, Omaha, NE 68198, USA.

27 \*Address for correspondence:

28 Jijo Pulickiyil Ulahannan Ph.D.,  
29 Government College Kasaragod,  
30 Kannur University, Kerala, India 671123

31 Email: [jijo@gck.ac.in](mailto:jijo@gck.ac.in)

32 # Contributed equally

33 **Word count:** Abstract: 249, Main text: 3753

34 **Keywords:** Kerala, India, COVID-19, Open data, Visualization

35 **ABSTRACT**

36 **Objective:** India reported its first COVID-19 case in the state of Kerala and an  
37 outbreak initiated subsequently. The Department of Health Services, Government of  
38 Kerala, initially released daily updates through daily textual bulletins for public  
39 awareness to control the spread of the disease. However, this unstructured data limits  
40 upstream applications, such as visualization, and analysis, thus demanding  
41 refinement to generate open and reusable datasets.

42 **Materials and Methods:** Through a citizen science initiative, we leveraged publicly  
43 available and crowd-verified data on COVID-19 outbreak in Kerala from the  
44 government bulletins and media outlets to generate reusable datasets. This was  
45 further visualized as a dashboard through a frontend web application and a JSON  
46 repository, which serves as an API for the frontend.

47 **Results:** From the sourced data, we provided real-time analysis, and daily updates of  
48 COVID-19 cases in Kerala, through a user-friendly bilingual dashboard  
49 (<https://covid19kerala.info/>) for non-specialists. To ensure longevity and reusability,  
50 the dataset was deposited in an open-access public repository for future analysis.  
51 Finally, we provide outbreak trends and demographic characteristics of the individuals  
52 affected with COVID-19 in Kerala during the first 138 days of the outbreak.

53 **Discussion:** We anticipate that our dataset can form the basis for future studies,  
54 supplemented with clinical and epidemiological data from the individuals affected with  
55 COVID-19 in Kerala.

56 **Conclusion:** We reported a citizen science initiative on the COVID-19 outbreak in  
57 Kerala to collect and deposit data in a structured format, which was utilized for

58 visualizing the outbreak trend and describing demographic characteristics of affected  
59 individuals.

## 60 **BACKGROUND AND SIGNIFICANCE**

61 In December 2019, an outbreak of cases presenting with pneumonia of unknown  
62 etiology was reported in Wuhan, China. The outbreak, caused by a novel severe acute  
63 respiratory syndrome Coronavirus-2 (SARS-CoV-2), later evolved as a pandemic  
64 (coronavirus disease 2019; COVID-19), claiming thousands of lives globally. [1-4]  
65 Initial studies revealed the clinical and prognostic features of COVID-19 along with its  
66 transmission dynamics and stressed the need for implementing public health  
67 measures for containment of infection and transmission among the population at high-  
68 risk. [2 5-9] In response to this, several countries have implemented measures  
69 including travel restrictions and physical distancing by community-wide quarantine. [2  
70 6 10] These extensive measures were imposed, taking into consideration the lack of  
71 adequate testing kits for detection, a vaccine, or proven antivirals for preventing or  
72 treating this disease along with reports of considerable strain on the health system  
73 leading to unprecedented loss of human life.

74 India—the second most populated country in the world—reported its first case  
75 in the state of Kerala on January 30, 2020, among individuals with travel history from  
76 Wuhan, the epicenter of the COVID-19 outbreak. [11] With the subsequent reports of  
77 an outbreak in the Middle East and Europe, Kerala has been on high-alert for a  
78 potential outbreak, as an estimated 10% of the population work abroad and being an  
79 international tourist destination. [12 13] The state has a high population density, with  
80 a large proportion of the population falling in the adult and older age group. [14] This  
81 population also shows a high incidence of COVID-19-associated comorbidities such  
82 as hypertension, diabetes, and cardiovascular disease. [9 15-17] As evidenced by  
83 reports of other countries, these factors pose a significant threat for an outbreak and  
84 would exert a tremendous burden on the public healthcare system. [18-20]

85 Severe public health measures were implemented in the state of Kerala and  
86 across India to prevent an outbreak. International flights were banned by March 22,  
87 2020, and a nation-wide lockdown was initiated on March 25, 2020. [21] However,  
88 before these measures were implemented, several cases (including travelers from  
89 Europe and the Middle East), along with a few reports of secondary transmission, were  
90 reported in Kerala. Since the first case was reported, the Department of Health  
91 Services (DHS), Government of Kerala, initiated diagnostic testing, isolation, contact  
92 tracing, and social distancing through quarantine, and the details of cases were  
93 released for the public through daily textual bulletins.

94 For pandemics such as COVID-19, public awareness via dissemination of  
95 reliable information in real-time plays a significant role in controlling the spread of the  
96 disease. Besides, real-time monitoring for identifying the magnitude of spread helps in  
97 hotspot identification, potential intervention measures, resource allocation, and crisis  
98 management. [22] The lack of such a real-time data visualization dashboard for the  
99 public with granular information specific to Kerala in the local language (Malayalam),  
100 during the initial days of the outbreak, was the motivation for this work.

101 To achieve this, the collection of relevant information on infection and refining  
102 the dataset in a structured manner for upstream purposes such as visualization and/or  
103 epidemiological analysis is essential. Open or crowd-sourced data has immense  
104 potential during the early stage of an outbreak, considering the limitation of obtaining  
105 detailed clinical and epidemiological data in real-time during an outbreak. [23-25]  
106 Furthermore, the structured datasets, when deposited in open repositories and  
107 archived, can ensure longevity for future analytical efforts and policymaking. The  
108 unavailability of such structured, reusable, and crowd-verified datasets on natural

109 disasters in Kerala, documented in the public domain, also motivated us to generate  
110 a resource for the COVID-19 outbreak. This initiative was volunteered by the Collective  
111 for Open Data Distribution- Keralam (CODD-K), a group of technologists,  
112 academicians, students, and the public advocating for open data. This collective, in a  
113 primitive form, was initiated during the devastating 2018 Kerala floods, which brought  
114 together the experts and general public through social media platforms to coordinate  
115 rescue missions through citizen-led open/crowd-sourcing strategies.

116 Here, we report a citizen science initiative to leverage publicly available data on  
117 COVID-19 cases in Kerala from the daily bulletins released by the DHS, Government  
118 of Kerala, and various news outlets. The multi-sourced data was refined to make a  
119 structured live dataset to provide real-time analysis and daily updates of COVID-19  
120 cases in Kerala through a bilingual (English and Malayalam) user-friendly dashboard  
121 (<https://covid19kerala.info/>). We aimed to disseminate the data of the outbreak trend,  
122 hotspots maps, and daily statistics in a comprehensible manner for non-specialists  
123 with bilingual (Malayalam and English) interpretation. Next, we aimed for the longevity  
124 and reusability of the datasets by depositing it in a public repository, aligning with open  
125 data principles for future analytical efforts. [26] Finally, to show the scope of the  
126 sourced data, we provide a snapshot of outbreak trends and demographic  
127 characteristics of the individuals affected with COVID-19 in Kerala during the first 138  
128 days of the outbreak.

129 **METHODS**

130 **The citizen-led collective for data sourcing and curation**

131 The CODD-K constituting, members from different domains, who shared the interest  
132 for sourcing data, building the dataset, visualizing, distributing, and interpreting the  
133 data on infection outbreak volunteered this effort (<https://team.covid19kerala.info/>).  
134 This initiative was in agreement with definitions proposed by different citizen-science  
135 initiatives.[27 28] The CODD-K invited participation in this initiative from the public  
136 through social media. The domain experts in the collective defined the data of interest  
137 to be collected, established the informatics workflow, and the web application for data  
138 visualization. The volunteers contributed by sourcing data from various media outlets  
139 for enriching the data. Dedicated social media (dedicated Telegram channels and  
140 WhatsApp groups) channels were used for data collection, which was verified  
141 independently and curated by data validation team members.

142 **Definition and Scope of Datasets**

143 The collective defined the data of interest as minimal structured metadata of the  
144 COVID-19 infections in Kerala, covering the possible facets of its spatial and temporal  
145 nature, excluding the clinical records ([Supplementary Methods](#)). The resulting  
146 datasets should maintain homogeneity and consistency, assuring the privacy and  
147 anonymity of the individuals. The notion of this data definition is to make the resulting  
148 datasets reusable and interoperable with similar or related datasets. A set of controlled  
149 vocabularies were formed as a core of this knowledge organization system to reduce  
150 anomalies, prevent typographical errors, and duplicate entries. Together with the  
151 controlled vocabularies, identifiers of individual entries in each dataset make the  
152 datasets interlinked. An essential set of authority control is used in populating spatial

153 data to make it accurate in the naming and hierarchy. A substantial set of secondary  
154 datasets were also produced and maintained along with the primary datasets,  
155 including derived and combined information from the primary datasets and external  
156 resources.

157 **Data Collection**

158 We primarily sourced publicly available de-identified data, released daily as textual  
159 bulletins (from January 31, 2020) by the DHS, Government of Kerala, India  
160 (<https://dhs.kerala.gov.in>), of the individuals diagnostically confirmed positive for  
161 SARS-CoV-2 by reverse transcription-polymerase chain reaction (RT-PCR) at the  
162 government-approved test centers. We also collected and curated reports from print  
163 and visual media for supplementing the data ([Supplementary methods](#)). The quality of  
164 the data in terms of veracity and selection bias has been ensured as described  
165 ([Supplementary Methods](#)). Utmost care was taken to remove any identifiable  
166 information to ensure the privacy of the subjects. Entries were verified independently  
167 by CODD-K data validation team members and rectified for inconsistencies ([Figure 1](#)).  
168 Since the data collected were publicly available, no individual consent and ethical  
169 approval were required for the study. To demonstrate the utility of the collected  
170 dataset, we provided the status of the first 138 days (between January 30, 2020, and  
171 June 15, 2020) of the COVID-19 outbreak in Kerala, and also described demographic  
172 characteristics of the individuals affected with COVID-19. We ensured that the sourced  
173 dataset complied with the Open Definition 2.1 laid down by Open Knowledge  
174 Foundation. [26]

175 **Implementation of Web Application**

176 A publicly accessible dashboard for the project is developed from a similar open-  
177 source project covid19japan.com. [29] The dashboard and related source codes are  
178 adapted and released as open-source software under MIT license, a permissive open-  
179 source software license. The dashboard has two distinctive components, a single page  
180 frontend web application accessible at <https://covid19kerala.info/>, and a JavaScript  
181 Object Notation (JSON) repository, which serves as an application programming  
182 interface (API) for the frontend. The API fetches data from the Google sheet and  
183 generates JSON files periodically with GitHub Actions. Both the web application and  
184 the API were created with JavaScript as the programming language and maintained  
185 using NodeJS. These portals use static-file assets without any server-side  
186 technologies. The website and the API are served through GitHub Pages, a free static  
187 web hosting service provided by GitHub ([Figure 2](#)).

188 **Hotspot Mapping**

189 COVID-19 hotspots for the Local Self Government (LSG) administration area—  
190 Panchayats, Municipalities, and Corporations were notified by the Government of  
191 Kerala, based on the recommendations ([Supplementary Methods](#)) of the Kerala State  
192 Disaster Management Authority and were updated daily through DHS bulletin as text  
193 data. A set of metadata for the LSGs, manually derived from multiple official sources  
194 with labels in both English and Malayalam, was made as an authority control for  
195 hotspots. Hotspots declared in daily bulletins are mapped to the identifiers in the LSG  
196 authority control, and containment zones were added as additional information. The  
197 LSG controlled vocabulary ensures location accuracy as well as eliminates duplicates  
198 and spelling irregularities. An independent generator periodically fetches the created

199 hotspot list, adds spatial geometry along with the LSG metadata and generates the  
200 hotspot dataset for the dashboard. The spatial geometry of the LSGs are also  
201 manually sourced from different public resources and optimized for minimal visual  
202 indication of the boundaries of the LSGs. On the dashboard, Mapbox service renders  
203 this GeoJSON as an interactive map. [30]

204 **RESULTS**

205 **Open-Data Release**

206 The resulting open-data sets are published under Open Data Commons Attribution  
207 License v1.0 (ODC-BY 1.0). A manually curated data archive is maintained as a  
208 GitHub repository for the provenance. [31] The datasets are provided with the schema  
209 definition and an actionable data-package declaration. [32] Periodic versioned  
210 snapshots were released as 'Covid19Kerala.info-Data' through Zenodo  
211 (<https://zenodo.org/>). [33] CODD-K manages the longevity and stewardship of the  
212 data. Sufficient documentation is provided to increase the adaptability of the datasets.  
213 We ensured that the datasets complied with the Open Definition 2.1, which would  
214 enable findability, easy access, sharing, reuse, and interoperability. [26] Additionally,  
215 as per the 5-Star Linked Open Data concepts, an incremental framework for deploying  
216 data, the dataset which we sourced, enriched, and disseminated, when complied with  
217 Open Definition 2.1, evolved to 3-star open data from the 1-Star open data released  
218 by the DHS. [34] Thus, our effort by aligning to Open Definition 2.1 significantly  
219 increased openness of the data.

220 **Visualization of the COVID-19 data through a dashboard**

221 Here we have collected, cleaned, and visualized publicly available data in a user-  
222 friendly bilingual progressive web application (PWA) designed to be both device and

223 browser agnostic. For the convenience of the public, the dashboard mainly highlighted  
224 the number of individuals who are hospitalized, tested, confirmed, currently active,  
225 deceased, recovered, and people under observation (State-wise and District Data),  
226 updated daily. We also visualized maps for hotspots, and active patients, along with  
227 outbreak spread trend (new, active, and recovered cases), new cases by day,  
228 diagnostic testing trend, patients—age breakup, confirmed case trajectories at the  
229 district administration level ([Figure 3A, B, and Supplementary Figure 1](#)). To the best  
230 of our knowledge, our dashboard was the first one to be online (March 22, 2020) with  
231 a bilingual dashboard with English and the local language Malayalam, featuring  
232 outbreak map, hotspot map, and trend line map with reports of new, active, and  
233 recovered cases, along with COVID-19 related deaths in Kerala. The official  
234 dashboard version by DHS followed later. We regularly received feedback from the  
235 users and added new plots and visual tools based on user recommendations. Till June  
236 15, 2020, the web application has seen 37,205 unique users, with an average of 2,000  
237 visitors per day. The source code and data were open for the public to fork and  
238 analyze, thus providing a framework for a data collection, analysis, and visualization  
239 platform for future disease outbreaks.

240 **Mapping of Hotspots for early outbreak identification**

241 Since the SARS-CoV-2 infection outbreak occurs in clusters, early identification and  
242 isolation of these clusters are essential to contain the outbreak. Accurate tracking of  
243 the new cases and real-time surveillance is essential for the effective mitigation of  
244 COVID-19. However, the daily public bulletins by DHS did not have any unique  
245 identification code for the COVID-19 infected individuals and also for secondary  
246 contacts who have contracted the infection through contact transmission. This limited  
247 us from tracking the transmission dynamics. As an alternative, we resorted to mapping

248 hotspots for infection as a proxy measure to indicate possible outbreak areas. Initially,  
249 red, orange, and green zones based on the number of cases were designated to each  
250 district by the Government of India. Later, the Government of Kerala started releasing  
251 COVID-19 hotspot regions of the LSG administration area. We manually curated the  
252 hotspot information from the DHS bulletins, and the dataset was published as a static  
253 JSON file in the GeoJSON format, which improves the browser caching and drops the  
254 requirement of server-sided API services. The hotspot locations were highlighted as  
255 red dots with descriptions, and when zoomed, the LSG administration area will be  
256 displayed on the map. In order to improve the visual clarity of hotspots with varying  
257 sizes of the LSGs and different zoom levels in browsers, an identifiable spot is placed  
258 on the visual center of the LSG area polygon. This inner center of the polygon was  
259 calculated with an iterative grid algorithm. To the best of our knowledge, this feature  
260 is unique to our dashboard. We also provided a toggle bar to visualize district  
261 boundaries and areas declared as hotspots at LSG resolution ([Figure 3C](#)). Owing to  
262 the lack of data, additional information such as the number of active cases in these  
263 hotspots could not be plotted.

264 **Outbreak trend and demographic characteristics of individuals affected with  
265 COVID-19 in Kerala from the dataset**

266 To understand the outbreak trend and demographic characteristics of the COVID-19  
267 infections in Kerala, we analyzed the dataset for the first 138 days of the outbreak,  
268 from January 30, 2020, to June 15, 2020. During this period, Kerala reported 2,543  
269 cases, of which 1,174 individuals recovered during the reported period, along with 20  
270 fatalities. Among the total number of COVID-19 infected individuals reported in Kerala,  
271 72.36% were males, and 26.03% were females, with a large proportion of individuals  
272 falling in the age group of 20-40 ([Table 1](#)). The median age of affected individuals was

273 36 (0-93) (male; 38 (0-93), female; 33 (0-88)). Around 84.66% of cases had a travel  
274 history to places with reported infection, and 15.30% were infected through secondary  
275 transmission. However, even as the number of reports during this time frame  
276 increased, there was no official report of community spread. During the reported  
277 period, the state declared 163 hotspots for infection, and currently (June 15, 2020),  
278 this number has reduced to 125. Kerala has established 34 testing centers (22  
279 government and 12 private) across the state and performed 151,686 tests during the  
280 period January 31, 2020, to June 15, 2020, which accounted for 4,359 tests per million  
281 of the population (TPR = 1.68%). In addition to routine testing, the DHS implemented  
282 additional targeted testing and testing based on random sampling in the hotspot areas.  
283 The median duration of illness was 13 days, with a trend that showed increasing  
284 recovery time for the older age group ([Table 1](#)). Oldest individuals to recover were 93  
285 and 88 years old. The fatality rate of Kerala was 0.79%.

286 **Table 1:** Demographic characteristics of the individuals affected with COVID-19 in Kerala, India between January 30, 2020 to  
 287 June 15, 2020

	All cases; n (% of males)	Cases with travel history; n (% of males)	Secondary transmission; n (% of males)	Recovery; n (% of males)	Fatality; n (% of males)	Duration of illness; median (range)
Total	2543 (72.41)	2153 (75.6)	389 (53.5)	1174 (64.5)	20 (65.0)	13 (2-45)
Age break down						
<10	89 (43.82)	62	27	33	1	14 (5-32)
10 - 19	92 (54.34)	64	28	47	--	12 (7-27)
20 - 29	555 (69.19)	492	63	244	1	12 (4-42)
30 - 39	589 (78.78)	501	88	277	1	13 (4-37)
40 - 49	467 (79.44)	403	64	200	2	12 (3-44)
50 - 59	328 (76.83)	279	49	132	2	13 (3-45)
60 - 69	162 (76.54)	133	29	61	7	12 (5-35)
70 - 79	32 (50.0)	24	8	12	5	18 (4-23)
> 80	18 (44.44)	5	13	14	1	11 (4-41)
Unspecified	211 (63.51)	190	21	154	--	15 (2-38)
All cases		Cases with travel history	Secondary transmission	Recovery	Fatality	
Age; median (range)	36 (0-93)	36 (0-83)	36 (0-93)	35 (2-93)	63.5 (0-87)	

288 **DISCUSSION**

289 In this report, we describe a citizen science initiative that leveraged publicly available  
290 unstructured COVID-19 data released daily by the Government of Kerala supplemented  
291 with news from media outlets and structured this into a knowledge bank for quick and  
292 easy interpretation through a user-friendly bilingual dashboard. The motivation for such  
293 an initiative arose due to the paucity of a real-time data visualization dashboard specific  
294 to Kerala during the initial stages of the outbreak. To the best of our knowledge, we were  
295 the first to host a visualization dashboard for COVID-19 outbreak in Kerala, with a user-  
296 friendly bilingual interface and unique features such as hotspot map. We reason that  
297 accurate information about the pandemic has made the public vigilant to adopt  
298 appropriate precautionary measures in controlling the outbreak. Our dashboard also has  
299 contributed to achieving this feat, as evidenced by the usage statistics within days of the  
300 launch. Furthermore, this open/multi-sourced dataset with a set of correlated temporal  
301 and spatial metadata was also made available for the public through an open repository,  
302 enabling retrospective analyses.

303 The framework developed for dataset generation and visualization can potentially  
304 be a model for advancing biomedical informatics, from a citizen-science/open data  
305 perspective. Specifically, our initiative rapidly established an easily adaptable platform  
306 and workflow for potential disease outbreaks and similar calamities, especially in  
307 resource-limited settings. With a reasonably minimal definition of data/metadata,  
308 adhering to the Open Definition 2.1, our dataset permits data-driven research on the  
309 epidemiology of the COVID-19 outbreak in Kerala and also increased openness as per  
310 5-Star Linked Open Data concepts. Furthermore, the temporal and spatial metadata

311 might aid in future studies involving genetic lineage diversity of SARS-CoV-2 in Kerala, in  
312 relation to the demographic characteristics and clinical phenotypes. [35 36] Thus, our  
313 model also sets an example for efficient data management in such citizen-science  
314 initiatives.

315 While the real-time information serves the public for assessing potential risk based  
316 on the outbreak trend/containment in a specific location; the inferences made from the  
317 emerging demographic data such as gender, age, recovery, and mortality statistics can  
318 help in refining our responses and understanding the epidemiology of COVID-19  
319 outbreak. Also, it provides helpful insights into a rapidly developing novel pandemic for  
320 policymaking, social awareness, and enhancing compliance with the Government  
321 policies. Additionally, retrospective analysis can give insights on how policy changes or  
322 other events altered the dynamics of the COVID-19 outbreak.

323 Kerala has effectively utilized open/crowd-sourcing platform using citizen-led  
324 initiatives to coordinate rescue missions through social media platforms during the floods  
325 that devastated the state during 2018 and 2019. [37-39] Our collective, CODD-K evolved  
326 as a result of crowd-sourced volunteering and coordination during the floods in Kerala  
327 from 2018. Our experience during flood volunteering and the lack of appropriate data  
328 archiving during this disaster prompted us to design a real-time dashboard for COVID-19  
329 pandemic proactively. This experience enabled us to assemble a team and launch the  
330 dashboard as rapid response during this pandemic. Experts from various domains and  
331 the general public assembled and volunteered to source data, build the dataset, visualize,  
332 distribute, and interpret the data on the outbreak through this collective. A series of recent  
333 studies involving crowd/open-source visualization of COVID-19 outbreak statistics have

334 indicated wide popularity and impact of these community-led initiatives, including in India.  
335 [23-25] However, our approach differed from those as we sourced unstructured official  
336 data released by the government, supplemented by the information from media outlets.  
337 This strategy not only ensures authenticity but also enriches the data available in the  
338 public domain into a structured dataset, though it depends on the data release policies  
339 adopted by the different state governments. Kerala is one of the many states in India with  
340 a transparent data release policy, which ensured the authenticity of data collected through  
341 our initiative. Furthermore, the granularity of the data at the LSG levels, which are  
342 manually verified (as released in local language) gives an added advantage, in terms of  
343 data depth, over other Pan-Indian dashboards that rely on APIs to fetch cumulative data.

344 Although this approach seems to be efficient, an unexpected surge in cases can  
345 jeopardize the data collection, thus limiting the feasibility. During such a scenario, a trade-  
346 off between depth and breadth of data collected has to be decided. Moreover, this  
347 approach also has inherent limitations, including issues with the veracity of data, owing  
348 to the anonymity, and depth of the data released, including clinical symptoms. Since each  
349 infected case identified in Kerala was not provided with a unique ID, it was impossible to  
350 track these cases for the assessment of vital epidemiological parameters like the  
351 reproduction number ( $R_0$ ). Based on our experience of collating and analyzing COVID-19  
352 data from the public domain in Kerala, we propose to frame specific guidelines for the  
353 public data release for COVID-19 or other epidemics. We recommend the release of  
354 official COVID-19 data in a consistent, structured and machine-readable format, in  
355 addition to the bulletins, which could be provided with a permanent URL and also archived  
356 in a public repository for future retrospective analyses. We also suggest releasing the

357 assigned unique ID for the individuals affected with COVID-19, to avoid inconsistencies  
358 in reporting and to enable tracking the secondary transmission. Furthermore, providing  
359 COVID-19 associated symptomatic information, without compromising the privacy of the  
360 infected individuals will also aid in the basic understanding of the disease through  
361 analytical approaches.

362 Our dataset, compiled between January 30, 2020, to June 15, 2020, indicates that  
363 the infections reported in Kerala were mainly among working-age men, with a travel  
364 history of places with COVID-19 outbreak. The absence of reported community spread in  
365 this period emphasizes the effectiveness of government implemented rapid testing and  
366 quarantine measures. Active tracking and isolation of cases with travel history lead to  
367 better management with minimal COVID-19-associated death. Since the majority of  
368 cases reported in Kerala were within the age group of 20-40 years, and the patients being  
369 in constant inpatient care possibly contributed to a better outcome and lesser mortality  
370 rate, respectively. Kerala implemented vigorous COVID-19 testing, and even though the  
371 test rate was relatively low (4,359 tests per million of the population), early testing  
372 combined with strict quarantine policies for individuals with travel history prevented  
373 community spread. However, the average number of positives detected for 1,000 tests  
374 (individuals) was lesser compared to other states in India, thus negating community  
375 spread. Data from Kerala also provides insights about the mean duration of illness and  
376 the effect of increasing age on this parameter.

377 Collectively, we report a citizen science initiative on the COVID-19 outbreak in  
378 Kerala to collect data in a structured format utilized for visualizing the outbreak trend and  
379 describing demographic characteristics of affected individuals. While the core aim of this

380 initiative is to document COVID-19 related information for the public, researchers, and  
381 policymakers, the implemented data visualization tool also alleviates the citizen's anxiety  
382 around the pandemic in Kerala. We anticipate that the dataset collected will form the basis  
383 for future studies, supplemented with detailed information on clinical and epidemiological  
384 parameters from individuals with COVID-19 infection in Kerala.

385 **Acknowledgments:**

386 We acknowledge Shane Reustle for his help and support for forking the Japan COVID-

387 19 Coronavirus Tracker repository and implementation of the dashboard. We thank Jiahui

388 Zhou for the original concept and design of the tracker. We also thank Sajjad Anwar for

389 generously providing the administrative boundary shapefiles and geoJSONS for Kerala.

390 Maps were generously provided by the Mapbox community team.

391 **Competing Interests:**

392 The authors declare no competing interests

393 **Funding:**

394 This study was not funded by any agencies and was purely a voluntary effort during the

395 community-wide quarantine period by a team of technologists, academicians, students,

396 and the general public advocating open data and citizen science.

397 **Authors contribution:**

398 Conceptualization; JiU,

399 Data collection and curation; JiU, NN, PP, SC, SPS, MM, SJ, JeU, MK, US

400 Formal analysis; JiU, NN, NT,

401 Methodology; JiU, NN, NT, SPS, AB, MK,

402 Resources; NT, MK, AB

403 Software; NT, AB, MK,

404 Supervision; JiU, STT, RE, SB

405 Visualization; NT, AB, PP, JiU, NN, SB

406 Roles/Writing - original draft; SB, NNV

407 Writing - review & editing; SB, NNV, JiU, NN, NT

- 408 **References**
- 409
- 410 1. Bedford J, Enria D, Giesecke J, et al. COVID-19: towards controlling of a pandemic. Lancet  
411 2020;395(10229):1015-18 doi: 10.1016/S0140-6736(20)30673-5[published Online First:  
412 Epub Date]].
- 413 2. Zhou F, Yu T, Du R, et al. Clinical course and risk factors for mortality of adult inpatients with  
414 COVID-19 in Wuhan, China: a retrospective cohort study. The Lancet  
415 2020;395(10229):1054-62 doi: [https://doi.org/10.1016/S0140-6736\(20\)30566-3](https://doi.org/10.1016/S0140-6736(20)30566-3)[published Online First: Epub Date]].
- 416
- 417 3. Zhou P, Yang X-L, Wang X-G, et al. A pneumonia outbreak associated with a new  
418 coronavirus of probable bat origin. Nature 2020;579(7798):270-73 doi: 10.1038/s41586-  
419 020-2012-7[published Online First: Epub Date]].
- 420 4. Wu F, Zhao S, Yu B, et al. A new coronavirus associated with human respiratory disease in  
421 China. Nature 2020;579(7798):265-69 doi: 10.1038/s41586-020-2008-3[published  
422 Online First: Epub Date]].
- 423 5. Li Q, Guan X, Wu P, et al. Early Transmission Dynamics in Wuhan, China, of Novel  
424 Coronavirus-Infected Pneumonia. New England Journal of Medicine  
425 2020;382(13):1199-207 doi: 10.1056/NEJMoa2001316[published Online First: Epub  
426 Date]].
- 427 6. Arons MM, Hatfield KM, Reddy SC, et al. Presymptomatic SARS-CoV-2 Infections and  
428 Transmission in a Skilled Nursing Facility. N Engl J Med 2020 doi:  
429 10.1056/NEJMoa2008457[published Online First: Epub Date]].
- 430 7. Bai Y, Yao L, Wei T, et al. Presumed Asymptomatic Carrier Transmission of COVID-19.  
431 JAMA 2020;323(14):1406-07 doi: 10.1001/jama.2020.2565[published Online First: Epub  
432 Date]].
- 433 8. Wu Z, McGoogan JM. Characteristics of and Important Lessons From the Coronavirus  
434 Disease 2019 (COVID-19) Outbreak in China: Summary of a Report of 72 314 Cases  
435 From the Chinese Center for Disease Control and Prevention. JAMA  
436 2020;323(13):1239-42 doi: 10.1001/jama.2020.2648[published Online First: Epub Date]].
- 437 9. Yang J, Zheng Y, Gou X, et al. Prevalence of comorbidities and its effects in patients infected  
438 with SARS-CoV-2: a systematic review and meta-analysis. International Journal of  
439 Infectious Diseases 2020;94:91-95 doi:  
440 <https://doi.org/10.1016/j.ijid.2020.03.017>[published Online First: Epub Date]].
- 441 10. Devi S. Travel restrictions hampering COVID-19 response. The Lancet  
442 2020;395(10233):1331-32 doi: [https://doi.org/10.1016/S0140-6736\(20\)30967-3](https://doi.org/10.1016/S0140-6736(20)30967-3)[published Online First: Epub Date]].
- 443
- 444 11. Yadav PD, Potdar VA, Choudhary ML, et al. Full-genome sequences of the first two SARS-  
445 CoV-2 viruses from India. The Indian journal of medical research 2020
- 446 12. Rajan SI, Zachariah KC. Emigration and Remittances: New Evidences from the Kerala  
447 Migration Survey 2018, 2019.
- 448 13. Thimm T. The Kerala tourism model—An Indian state on the road to sustainable  
449 development. Sustainable Development 2017;25(1):77-91
- 450 14. Board KSP. Economic Review 2019. Economic Review 2019. Kerala: Government of  
451 Kerala, 2020.
- 452 15. Ghosh S, Kumar M. Prevalence and associated risk factors of hypertension among persons  
453 aged 15–49 in India: a cross-sectional study. BMJ open 2019;9(12)
- 454 16. Prabhakaran D, Jeemon P, Sharma M, et al. The changing patterns of cardiovascular  
455 diseases and their risk factors in the states of India: the Global Burden of Disease Study  
456 1990–2016. The Lancet Global Health 2018;6(12):e1339-e51

- 457 17. Vijayakumar G, Manghat S, Vijayakumar R, et al. Incidence of type 2 diabetes mellitus and  
458 prediabetes in Kerala, India: results from a 10-year prospective cohort. BMC public  
459 health 2019;19(1):140
- 460 18. Ji Y, Ma Z, Peppelenbosch MP, Pan Q. Potential association between COVID-19 mortality  
461 and health-care resource availability. The Lancet Global Health 2020;8(4):e480
- 462 19. Boccia S, Ricciardi W, Ioannidis JP. What other countries can learn from Italy during the  
463 COVID-19 pandemic. JAMA Internal Medicine 2020
- 464 20. Wadhera RK, Wadhera P, Gaba P, et al. Variation in COVID-19 Hospitalizations and Deaths  
465 Across New York City Boroughs. JAMA 2020
- 466 21. Lancet T. India under COVID-19 lockdown. The Lancet 2020;395(10233):1315 doi:  
467 [https://doi.org/10.1016/S0140-6736\(20\)30938-7](https://doi.org/10.1016/S0140-6736(20)30938-7)[published Online First: Epub Date]].
- 468 22. Rivers C, Chretien J-P, Riley S, et al. Using “outbreak science” to strengthen the use of  
469 models during epidemics. Nature Communications 2019;10(1):3102 doi:  
470 10.1038/s41467-019-11067-2[published Online First: Epub Date]].
- 471 23. Dong E, Du H, Gardner L. An interactive web-based dashboard to track COVID-19 in real  
472 time. The Lancet Infectious Diseases 2020 doi: [https://doi.org/10.1016/S1473-  
473 3099\(20\)30120-1](https://doi.org/10.1016/S1473-3099(20)30120-1)[published Online First: Epub Date]].
- 474 24. Sun K, Chen J, Viboud C. Early epidemiological analysis of the coronavirus disease 2019  
475 outbreak based on crowdsourced data: a population-level observational study. The  
476 Lancet Digital Health 2020;2(4):e201-e08 doi: [https://doi.org/10.1016/S2589-  
477 7500\(20\)30026-1](https://doi.org/10.1016/S2589-7500(20)30026-1)[published Online First: Epub Date]].
- 478 25. Xu B, Kraemer MU, Group DC. Open access epidemiological data from the COVID-19  
479 outbreak. The Lancet. Infectious Diseases 2020
- 480 26. Open Knowledge Foundation. Open Definition 2.1. Secondary Open Definition 2.1 2017.  
481 <http://opendefinition.org/od/2.1/en/>.
- 482 27. Robinson LD, Cawthray J, West SE, Bonn A, Ansine J. Ten principles of citizen science:  
483 UCL Press, 2018.
- 484 28. Heigl F, Kieslinger B, Paul KT, Uhlik J, Dörler D. Opinion: Toward an international definition  
485 of citizen science. Proceedings of the National Academy of Sciences  
486 2019;116(17):8089-92 doi: 10.1073/pnas.1903393116[published Online First: Epub  
487 Date]].
- 488 29. covid19japan.com. Secondary 2020. <https://github.com/reustle/covid19japan>.
- 489 30. Butler H, Daly M, Doyle A, Gillies S, Hagen S, Schaub T. The geojson format. Internet  
490 Engineering Task Force (IETF) 2016
- 491 31. covid19kerala.info. Secondary 2020. <https://purl.org/codd-k/c19k/data/v1.0>
- 492 32. Paul Walsh RP. Frictionless Data Specs. Secondary Frictionless Data Specs May 2, 2017  
493 2007. <https://specs.frictionlessdata.io/data-package/>.
- 494 33. Jijo U, Narayanan N, Suresh SP, et al. Covid19Kerala.info-Data: A collective open dataset  
495 of COVID-19 outbreak in the south Indian state of Kerala. Zenodo, 2020.
- 496 34. W3C Working Group. Linked Data Glossary. In: Bernadette Hyland, Ghislain Atemezing,  
497 Michael Pendleton, Srivastava B, eds. Technical Report. W3C Working Group Note:  
498 W3C, 2013.
- 499 35. Lu J, du Plessis L, Liu Z, et al. Genomic Epidemiology of SARS-CoV-2 in Guangdong  
500 Province, China. Cell 2020;181(5):997-1003.e9 doi:  
501 <https://doi.org/10.1016/j.cell.2020.04.023>[published Online First: Epub Date]].
- 502 36. Saez-Rodriguez J, Costello JC, Friend SH, et al. Crowdsourcing biomedical research:  
503 leveraging communities as innovation engines. Nat Rev Genet 2016;17(8):470-86 doi:  
504 10.1038/nrg.2016.69[published Online First: Epub Date]].
- 505 37. Ajay A. Role of technology in responding to disasters: insights from the great deluge in  
506 Kerala. Curr Sci India 2019;116(6):913-18 doi: 10.18520/cs/v116/i6/913-918[published  
507 Online First: Epub Date]].

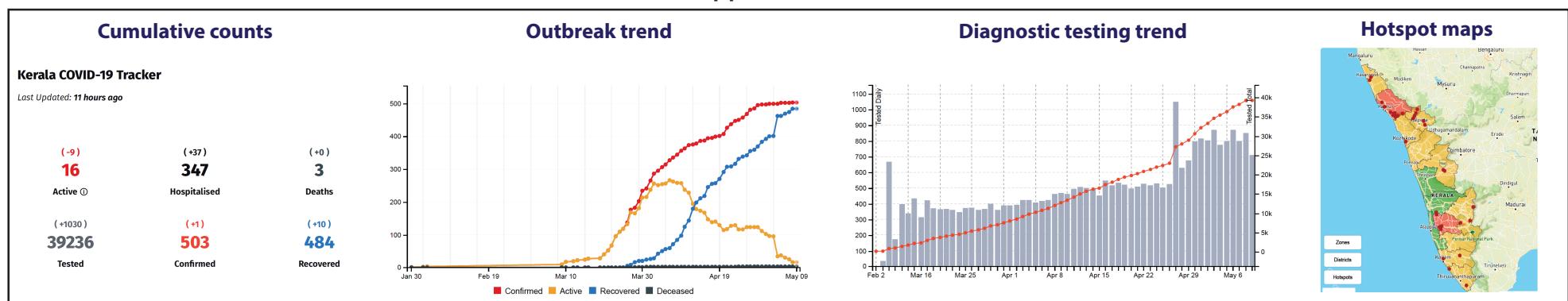
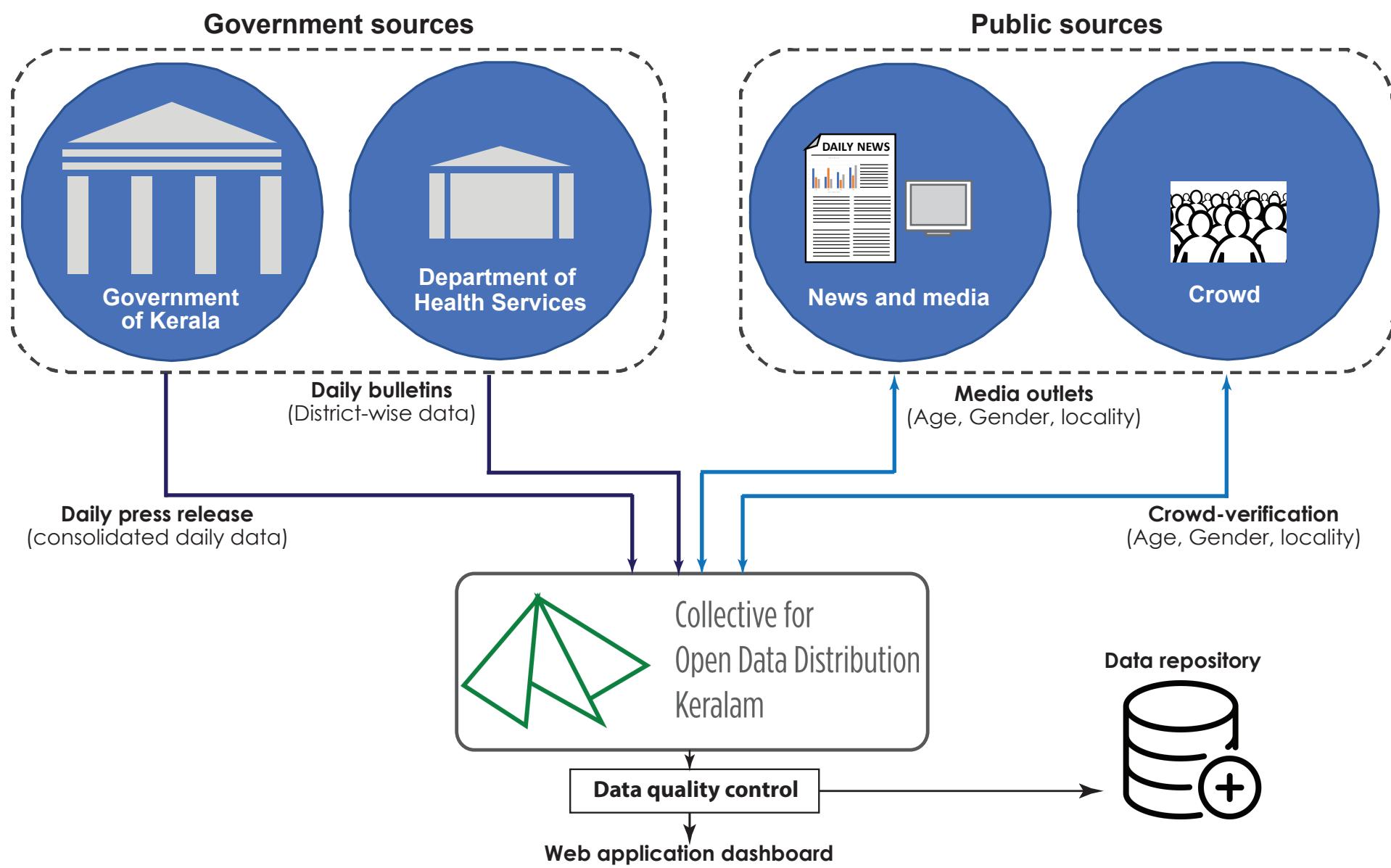
- 508 38. Architectural Considerations for Building a Robust Crowdsourced Disaster Relief  
509 Application. 2020 International Conference on COMmunication Systems & NETworkS  
510 (COMSNETS); 2020. IEEE.
- 511 39. Mishra AK, Nagaraju V. Space-based monitoring of severe flooding of a southern state in  
512 India during south-west monsoon season of 2018. Natural Hazards 2019;97(2):949-53
- 513
- 514

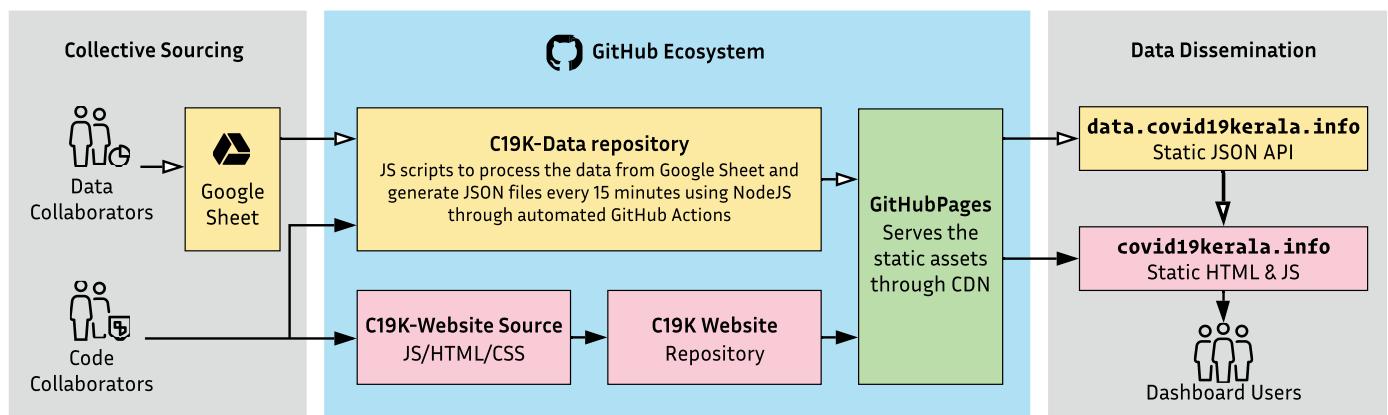
515 **Figures legends**

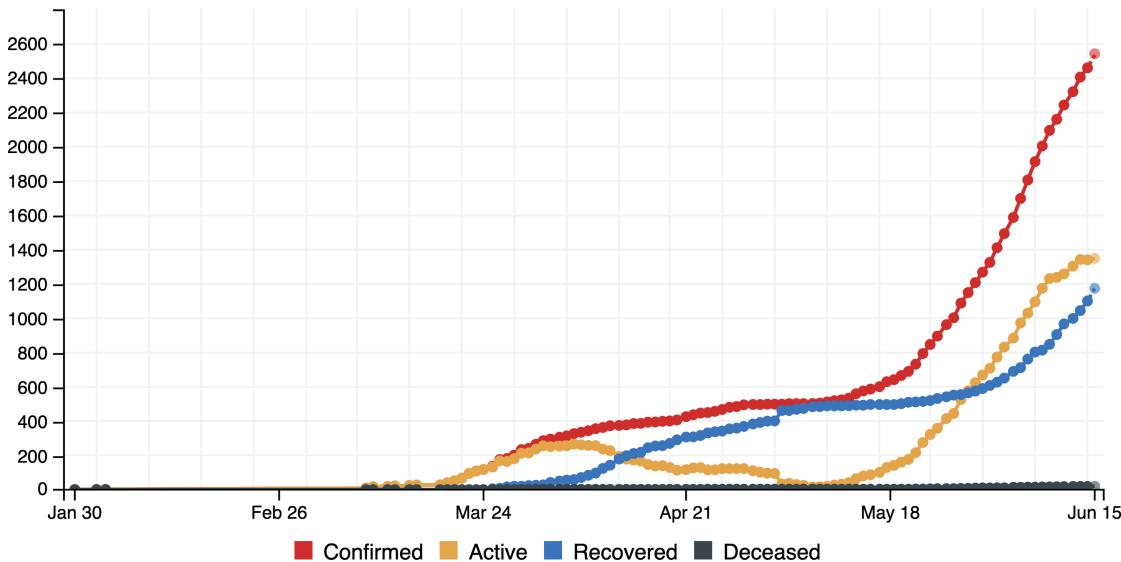
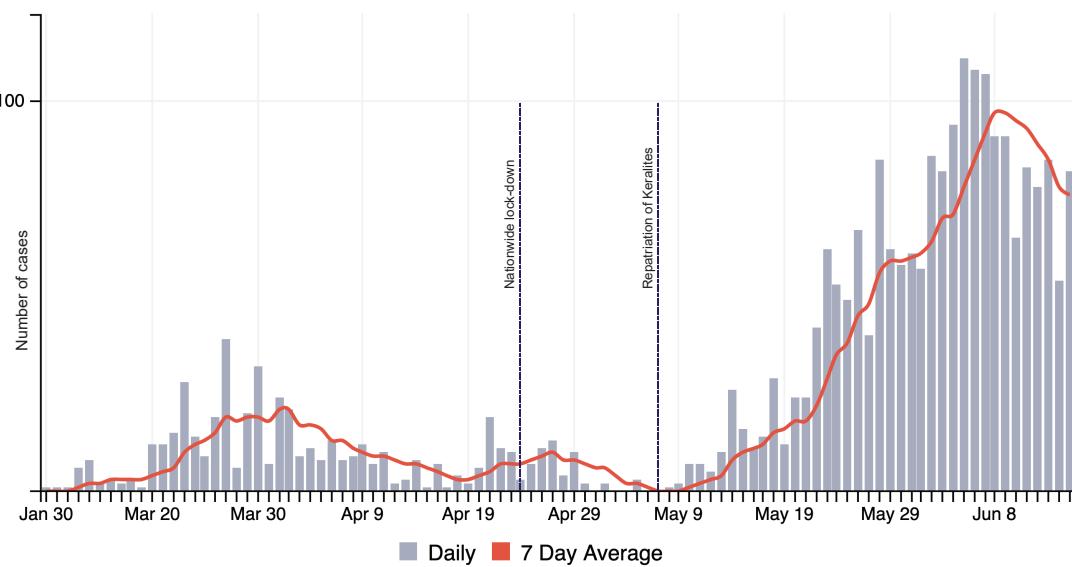
516 **Figure 1:** Outline of data collection, curation, and quality control for generating dataset  
517 and visualization

518 **Figure 2:** Implementation of web application and workflow

519 **Figure 3:** Representative images of COVID-19 outbreak trend for Kerala as visualized  
520 from the sourced data: (A) Plot showing number of confirmed, active, recovered and  
521 deceased cases (B) The trend curve, plotted with daily cases and seven days' average  
522 is shown. The dotted lines shows the initiation of nation-wide lockdown, and repatriation  
523 of Keralites from abroad and other states (C) the hotspot map showing the districts and  
524 hotspot location





**A****B****C**