

1 A Appendix

2 Datasheet

3 MOTIVATION

4 For what purpose was the dataset created? Was there a specific task in mind? Was there a
5 specific gap that needed to be filled? Please provide a description.

6 Indian-CT was constructed to provide chest CT scans of COVID-19 patients from India for
7 building machine learning models that can aid in analyzing lungs for the diagnosis of the
8 disease. This dataset can also be used as a category while detecting other pulmonary diseases.
9 To our knowledge, this is the only chest CT-scan dataset from the Indian population. Since
10 all the images included in this dataset are from a single hospital and acquired using a single
11 CT scanner, any confounding factors in the data are minimized.

12 Who created this dataset (e.g., which team, research group) and on behalf of which entity
13 (e.g., company, institution, organization)? The research team, including the authors, Suba S,
14 Nita Parekh, Ramesh Loganathan, Vikram Pudi, from International Institute of Information
15 Technology, Hyderabad, collected the data from Gandhi Hospital, Hyderabad.

16 Who funded the creation of the dataset? If there is an associated grant, please provide the
17 name of the grantor and the grant name and number

18 This work was funded by the RAKSHAK (Remedial Action, Knowledge Skimming and
19 Holistic Analysis of COVID-19) project of Department of Science and Technology (DST),
20 India.

21 Any other comments?

22 COMPOSITION

23 What do the instances that comprise the dataset represent (e.g., documents, photos, people,
24 countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people
25 and interactions between them; nodes and edges)? Please provide a description.

26 Each instance in the dataset is a slice of CT scan of the lungs taken from COVID-19 patients.
27 Multiple slices from each patient are available and each slice is stored in PNG format.
28 Except the last two digits, the name of the file is unique to a patient.

29 How many instances are there in total (of each type, if appropriate)?

30 Currently, there are 6174 images taken from 142 COVID-19 patients. Each patient has 40
31 slices (images) in the dataset.

32 Does the dataset contain all possible instances or is it a sample (not necessarily random) of
33 instances from a larger set? If the dataset is a sample, then what is the larger set? Is the
34 sample representative of the larger set (e.g., geographic coverage)? If so, please describe
35 how this representativeness was validated/verified. If it is not representative of the larger set,
36 please describe why not (e.g., to cover a more diverse range of instances, because instances
37 were withheld or unavailable).

38 This dataset does not contain all possible instances. Based on the availability of the data from
39 the hospital, currently, only 142 patient data is available. We will add more eventually. Each
40 CT volume was converted to png format after selecting only slices in the range 40 - 300 as
41 this range was found to consist of the broadest lung window devoid of other internal organs.
42 This heuristic could be applied on all the images as these are obtained from a single CT
43 scanner machine. Every 3rd slice from the chosen range was considered for analysis to
44 reduce the size of the dataset. For a few samples (< 10), however, since sufficient number of
45 slices were not available, every slice in the corresponding range was taken.

46 What data does each instance consist of? "Raw" data (e.g., unprocessed text or images) or
47 features? In either case, please provide a description.

48 Each instance consists of a slice from the CT scan of the lungs of COVID-19 patients in
49 PNG format. The size of each image is 512 by 512 pixels. The raw data from the hospital
50 was in .dcm format, which was converted to PNG.

51 Is there a label or target associated with each instance? If so, please provide a description.
52 The dataset consists of only images of COVID-19 class, and hence no separate label is
53 associated.

54 Is any information missing from individual instances? If so, please provide a description,
55 explaining why this information is missing (e.g., because it was unavailable). This does not
56 include intentionally removed information, but might include, e.g., redacted text

57 The slices in their raw form consisted of patients' personal identification information such
58 as name, age, date of birth, etc., along with the scanner settings which were removed by
59 default when converting to PNG format.

60 Are relationships between individual instances made explicit (e.g., users' movie ratings,
61 social network links)? If so, please describe how these relationships are made explicit.

62 There are no known relationships between individuals except for the fact that they all are
63 COVID-19 patients.

64 Are there recommended data splits (e.g., training, development/validation, testing)? If so,
65 please provide a description of these splits, explaining the rationale behind them.

66 The dataset does not contain any splits. We recommend using this dataset as a test set for
67 understanding the generalizability of machine learning models. This could also be used for
68 a detailed analysis of the lungs of COVID-19 patients from the Indian population.

69 Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a
70 description.

71 None to the best of our knowledge

72 Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,
73 websites, tweets, other datasets)? If it links to or relies on external resources, a) are there
74 guarantees that they will exist, and remain constant, over time; b) are there official archival
75 versions of the complete dataset (i.e., including the external resources as they existed at the
76 time the dataset was created); c) are there any restrictions (e.g., licenses, fees) associated with
77 any of the external resources that might apply to a future user? Please provide descriptions
78 of all external resources and any restrictions associated with them, as well as links or other
79 access points, as appropriate.

80 The data is self-contained.

81 Does the dataset contain data that might be considered confidential (e.g., data that is
82 protected by legal privilege or by doctor-patient confidentiality, data that includes the content
83 of individuals non-public communications)? If so, please provide a description.

84 No. All personal identification details have been removed from the data.

85 Does the dataset contain data that, if viewed directly, might be offensive, insulting, threaten-
86 ing, or might otherwise cause anxiety? If so, please describe why

87 No. The dataset consists of only CT scan images of lungs.

88 Does the dataset relate to people? If not, you may skip the remaining questions in this
89 section.

90 Yes.

91 Does the dataset identify any subpopulations (e.g., by age, gender)? If so, please de-
92 scribe how these subpopulations are identified and provide a description of their respective
93 distributions within the dataset.

94 The dataset consists of CT scan images of the lungs of COVID-19 patients, collected from a
95 hospital in Hyderabad, India. The age of the patients is in the range of 17 to 79 years, with a
96 mean age of 48 years. All the images were acquired using a single CT scanner, and hence
97 the scanner settings are identical for all images.

98 Is it possible to identify individuals (i.e., one or more natural persons), either directly or
99 indirectly (i.e., in combination with other data) from the dataset? If so, please describe how

100 No. It is not possible to identify individuals directly or indirectly.

101 Does the dataset contain data that might be considered sensitive in any way (e.g., data that
102 reveals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or
103 union memberships, or locations; financial or health data; biometric or genetic data; forms of
104 government identification, such as social security numbers; criminal history)? If so, please
105 provide a description.

106 No personal data of the patients whose CT scan images is included in available to the users.

107 Any other comments?

108 COLLECTION PROCESS

109 How was the data associated with each instance acquired? Was the data directly observable
110 (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly
111 inferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or
112 language)? If data was reported by subjects or indirectly inferred/derived from other data,
113 was the data validated/verified? If so, please describe how.

114 The dataset was acquired from the radiology department of Gandhi hospital, Hyderabad.
115 The third-party who assisted in collecting the data named the data COVID-19 CT data
116 collected from the hospital.

117 What mechanisms or procedures were used to collect the data (e.g., hardware apparatus
118 or sensor, manual human curation, software program, software API)? How were these
119 mechanisms or procedures validated?

120 The raw data was acquired using a single CT scanner machine, the details of which are given
121 in Table 1. The images are plain CT scans captured with no contrast and slice thickness of
122 the images are 0.6, 1.5 and 5 mm.

123 If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,
124 probabilistic with specific sampling probabilities)?

125 A total of 533 patient data was obtained from Gandhi hospital. Of this, 255 patient data were
126 considered for this study. On initial screening of these 255, 108 patient data were removed as
127 these did not exclusively belong to chest CT, had missing information like SliceLocation, or
128 came from different CT scanners. The remaining 278 patient data is under the preprocessing
129 stage and will eventually be added to the Indian-CT dataset. Each CT volume was converted
130 to png format after selecting only slices in the range 40 - 300 as this range was found to
131 consist of the broadest lung window devoid of other internal organs. This heuristic was
132 considered because the CT images in this dataset were obtained using only a single CT
133 scanner machine. Every 3rd slice from the chosen range was considered to reduce the size
134 of the dataset. However, for some samples sufficient slices were not available and in such
135 cases every slice in the corresponding range was taken.

136 Who was involved in the data collection process (e.g., students, crowdworkers, contractors)
137 and how were they compensated (e.g., how much were crowdworkers paid)?

138 A third party was involved in the collection process and were funded through RAKSHAK
139 project of Department of Science and Technology (DST), India.

140 Over what timeframe was the data collected? Does this timeframe match the creation
141 timeframe of the data associated with the instances (e.g., recent crawl of old news articles)?
142 If not, please describe the timeframe in which the data associated with the instances was
143 created.

144 Indian-CT data is collected from Gandhi Hospital, Hyderabad, India from the COVID-19
145 isolated patients during the period April - September, 2020.

146 Were any ethical review processes conducted (e.g., by an institutional review board)? If so,
147 please provide a description of these review processes, including the outcomes, as well as a
148 link or other access point to any supporting documentation.

149 The study was approved by the institutional review board (IRB) of Gandhi Hospital, Hyder-
150 abad. The institutional review boards waived the requirement to obtain written informed
151 consent for this retrospective study. For confidentiality, no link between the patients and the
152 researchers was made available.

153 Does the dataset relate to people? If not, you may skip the remaining questions in this
154 section.

155 Yes.

156 Did you collect the data from the individuals in question directly, or obtain it via third parties
157 or other sources (e.g., websites)?

158 This is a retrospective study, and the data was collected via a third party, from the hospital.

159 Were the individuals in question notified about the data collection? If so, please describe (or
160 show with screenshots or other information) how notice was provided, and provide a link or
161 other access point to, or otherwise reproduce, the exact language of the notification itself.

162 The study was approved by the institutional review board (IRB) of Gandhi Hospital, Hyder-
163 abad. The institutional review boards waived the requirement to obtain written informed
164 consent for this retrospective study. For confidentiality, no link between the patients and the
165 researchers was made available

166 Did the individuals in question consent to the collection and use of their data? If so, please
167 describe (or show with screenshots or other information) how consent was requested and
168 provided, and provide a link or other access point to, or otherwise reproduce, the exact
169 language to which the individuals consented.

170 This is a retrospective study, approved by the institutional review board (IRB) of Gandhi
171 Hospital, Hyderabad. The institutional review boards waived the requirement to obtain
172 written informed consent.

173 Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a
174 data protection impact analysis) been conducted? If so, please provide a description of this
175 analysis, including the outcomes, as well as a link or other access point to any supporting
176 documentation

177 There is no link connecting the dataset to the data subjects and hence no potential impact.

178 Any other comments?

179 PREPROCESSING/CLEANING/LABELLING

180 Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,
181 tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, pro-
182 cessing of missing values)? If so, please provide a description. If not, you may skip the
183 remainder of the questions in this section.

184 Yes. A total of 533 patient data was obtained from Gandhi hospital. Of this, 255 patient
185 data were considered for this study. On initial screening of the 255 samples, 108 patient data
186 were removed as these did not exclusively belong to chest CT, or had missing information
187 like SliceLocation, or came from different CT scanner. Each CT volume was converted to
188 png format after selecting only slices in the range 40 - 300 as this range was found to consist
189 of the broadest lung window devoid of other internal organs. This heuristic was considered
190 only because the CT images in this dataset was obtained using only a singly CT scanner
191 machine. Every 3rd slice from the chosen range was considered for analysis to reduce the
192 size of the dataset. However, for few samples sufficient slices were not available and in such
193 cases every slice in the corresponding range was taken.

194 Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to
195 support unanticipated future uses)? If so, please provide a link or other access point to the
196 “raw” data.

197 Yes. The raw data is saved but will only be shared at a later stage.

198 Is the software used to preprocess/clean/label the instances available? If so, please provide a
199 link or other access point.

200 Preprocessing was done using a python code.

201 Any other comments?

202 USES

203 Has the dataset been used for any tasks already? If so, please provide a description.
204 The dataset has been used as a test set for analyzing the performance of deep learning
205 models on this dataset. The details are given in the paper.

206 Is there a repository that links to any or all papers or systems that use the dataset? If so,
207 please provide a link or other access point.

208 No.

209 What (other) tasks could the dataset be used for?

210 The dataset is the only one of its kind dataset from the Indian population. Apart from using
211 the images for COVID-19 diagnosis problem, this dataset can be used for the analysis of
212 the lungs of patients belonging to the Indian population and can be used as a class in lung
213 disease classification problems. This can be used to test the generalizability of models
214 trained on a different population, as done in this study.

215 Is there anything about the composition of the dataset or the way it was collected and
216 preprocessed/cleaned/labeled that might impact future uses? For example, is there anything
217 that a future user might need to know to avoid uses that could result in unfair treatment
218 of individuals or groups (e.g., stereotyping, quality of service issues) or other undesirable
219 harms (e.g., financial harms, legal risks) If so, please provide a description. Is there anything
220 a future user could do to mitigate these undesirable harms?

221 NO

222 Are there tasks for which the dataset should not be used? If so, please provide a description.
223 None to the best of our knowledge.

224 Any other comments?

225 DISTRIBUTION

226 Will the dataset be distributed to third parties outside of the entity (e.g., company, institution,
227 organization) on behalf of which the dataset was created? If so, please provide a description.
228

229 The access will be only provided on request for research purposes, after completing an
230 enrollment form in which requesters agree to the collaborative effort for confidential use of
231 the data, and to cite the dataset in their publications. Ethics agreements have been procured
232 for this purpose.

233 How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the
234 dataset have a digital object identifier (DOI)?

235 Details of how to access the dataset are currently available on the AIMed website
236 (<http://aimedhub.iiit.ac.in/>). Eventually, this access will be moved to a Data Foundation
237 platform, which is currently being built at IIIT Hyderabad to host several datasets across
238 verticals including health, mobility, smart-cities, etc., as part of a national initiative (details
239 at <https://ihub-data.iiit.ac.in/>).

240 When will the dataset be distributed?

241 The dataset is available upon request.

242 Will the dataset be distributed under a copyright or other intellectual property (IP) license,
243 and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU,
244 and provide a link or other access point to, or otherwise reproduce, any relevant licensing
245 terms or ToU, as well as any fees associated with these restrictions.

246 Access to the dataset will be made available on request for research purposes, after complet-
247 ing an enrollment form in which requesters agree to the collaborative effort for confidential
248 use of the data, and to cite the dataset in their publications. Ethics agreements have been
249 procured for this purpose.

250 Have any third parties imposed IP-based or other restrictions on the data associated with the
251 instances? If so, please describe these restrictions, and provide a link or other access point
252 to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with
253 these restrictions

254 No
255 Do any export controls or other regulatory restrictions apply to the dataset or to individual
256 instances? If so, please describe these restrictions, and provide a link or other access point
257 to, or otherwise reproduce, any supporting documentation.
258 Current ethics agreements have authorized the dataset to be used by researchers and collabora-
259 tors of IIIT Hyderabad and Gandhi Hospital, Hyderabad. The ethics agreements have not
260 authorized the dataset to be hosted directly on any public platform. Thus, it may be made
261 available on request to researchers who agree to the collaborative effort for confidential use
262 of the data.
263 Any other comments?

264 MAINTENANCE

265 Who will be supporting/hosting/maintaining the dataset?
266 IIIT, Hyderabad will be hosting and maintaining the dataset.
267 How can the owner/curator/manager of the dataset be contacted (e.g., email address)?
268 By filling the download request form on the website <http://aimedhub.iiit.ac.in/>.
269 Is there an erratum? If so, please provide a link or other access point
270 No.
271 Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete in-
272 stances)? If so, please describe how often, by whom, and how updates will be communicated
273 to users (e.g., mailing list, GitHub)?
274 The dataset will be updated and the updates will be communicated to collaborators via email.
275
276 If the dataset relates to people, are there applicable limits on the retention of the data
277 associated with the instances (e.g., were individuals in question told that their data would be
278 retained for a fixed period of time and then deleted)? If so, please describe these limits and
279 explain how they will be enforced.
280 No. There is no applicable limits on retention of the data.
281 Will older versions of the dataset continue to be supported/hosted/maintained? If so, please
282 describe how. If not, please describe how its obsolescence will be communicated to users.
283 Any changes will be communicated to the collaborators via email.
284 If others want to extend/augment/build on/contribute to the dataset, is there a mechanism
285 for them to do so? If so, please provide a description. Will these contributions be vali-
286 dated/verified? If so, please describe how. If not, why not? Is there a process for commu-
287 nicating/distributing these contributions to other users? If so, please provide a description.
288
289 The ethics agreement in place has not authorized the dataset to be hosted openly on any
290 public platform, and hence the dataset is not available for extension and redistribution. It
291 is still possible for other datasets to provide references to this dataset, but not include this
292 dataset's records directly.
293 Any other comments?

294 Other Details

- 295 1. Manufacturer details of the CT machine using which the Indian-CT data was collected is
296 given in Table 1.
- 297 2. The training of the model was done using an existing dataset, COVID_x-CT-2. The COVID_x
298 CT-2 dataset is released under a CC BY-NC-SA 4.0 license in accordance with the licenses
299 of its constituent datasets. The details of each individual licence and the link from which the
300 dataset was downloaded is: <https://www.kaggle.com/hgunraj/covidxct>. The new asset that is
301 used for testing is proprietary and is not released yet.

Table 1: Manufacturer details of the CT machine

Key	Value
Manufacturer	SIEMENS
Modality	CT
Manufacturer's Model Name	Emotion 16
Device Serial Number	39306
Software Version(s)	Syngo CT 2014A

- 302 3. The study was approved by the institutional review board (IRB) of Gandhi Hospital, Hyder-
303 abad. The institutional review boards waived the requirement to obtain written informed
304 consent for this retrospective study., For confidentiality, no link between the patients and the
305 researchers was made available.
- 306 4. This study was retrospective and the IRB has waived the requirement of consent.
- 307 5. The data was checked manually, sample by sample for any identifiable information and only
308 dummy name, age and dummy date of births were found without any link to address or any
309 other information so that the patient can be traced.
- 310 6. This was a retrospective study and hence no potential participant risks are involved and no
311 participant compensation is applicable.