
COVID-19 India Dataset: Parsing Detailed COVID-19 Data in Daily Health Bulletins from States in India

Mayank Agarwal¹ Tathagata Chakraborti¹ Sachin Grover²

Abstract

While India remains one of the hotspots of the COVID-19 pandemic, data about the pandemic from the country has proved to be largely inaccessible for use at scale. Much of the data exists in an unstructured form on the web, and limited aspects of such data are available through public APIs maintained manually through volunteer efforts. This has proved to be difficult both in terms of ease of access to detailed data as well as with regards to the maintenance of manual data-keeping over time. This paper reports on a recently launched project aimed at automating the extraction of such data from public health bulletins with the help of a combination of classical PDF parsers as well as the state of the art ML-based documents extraction APIs. In this paper, we will describe the automated data-extraction technique, the nature of the generated data, and exciting avenues of ongoing work.

URL: ibm.biz/covid-data-india

1. Introduction

Availability of COVID-19 data is crucial for researchers and policymakers to understand the progression of the pandemic and react to it in real-time. However, unlike countries with uniform and well-defined data reporting mechanisms, pandemic data from India is available either through volunteer-driven initiatives, through special access granted by the government, or manually collected from daily health bulletins published by individual states and cities on their own websites or platforms.

While daily health bulletins from Indian states contain a wealth of data, they are only available in the unstructured form in PDF documents and images. On the other hand, volunteer-driven manual data-curation cannot scale

to the volume of data over time. For example, one of the most well-known sources of COVID data from India: covid19india.org, has manually maintained public APIs for limited data throughout the pandemic. Such approaches, while simultaneously limited in the detail of data made available, are also unlikely to continue in the long term due to the amount of volunteer manual labor required indefinitely. Although this project originally began anticipating that outcome, that eventuality has already come to pass for the aforementioned project, for similar reasons outlined in ([covid19india, 2021b](#)). As such, detailed COVID-19 data from India, in a structured form, remains inaccessible at scale. ([Priyanka Pulla, 2021](#)) notes pleas from researchers in India, earlier this year, for the urgent access to detailed COVID data collected by government agencies.

The aim of this project is to use document and image extraction techniques to automate the extraction of such data in structured (SQL) form from the state-level daily health bulletins; and make this data freely available. Our target is to automate the data extraction process, so that once the extraction for each state is complete, it requires little to no attention after that (other than responding to changes in the schema). The role of machine learning here is to make that extraction automated and robust in coverage and accuracy. This data goes beyond just daily case and vaccinations numbers but also includes comprehensive state-wise metrics such as the hospitalization data, age-wise distribution of cases, asymptomatic and symptomatic cases, and even case information for individuals in certain states.

India, one of the most populous countries in the world, has reported over 33 million confirmed cases of COVID-19 – second only to the United States. The massive scale of this data not only provides intriguing research opportunities in data science, document understanding, and NLP for AI researchers but will also help epidemiologists and public policy experts to analyze and derive key insights about the pandemic in real-time. At the time of this writing, covid19india.org has also released possible alternatives going forward once the current APIs are sunset next month. These suggestions, detailed here: ([covid19india, 2021a](#)), also align perfectly with this current project and give us hope that we can continue providing this data, at scale and with much more detail than ever before.

¹IBM Research ²Arizona State University. Correspondence to: Mayank Agarwal <mayank.agarwal@ibm.com>, Tathagata Chakraborti <tchakra2@ibm.com>.

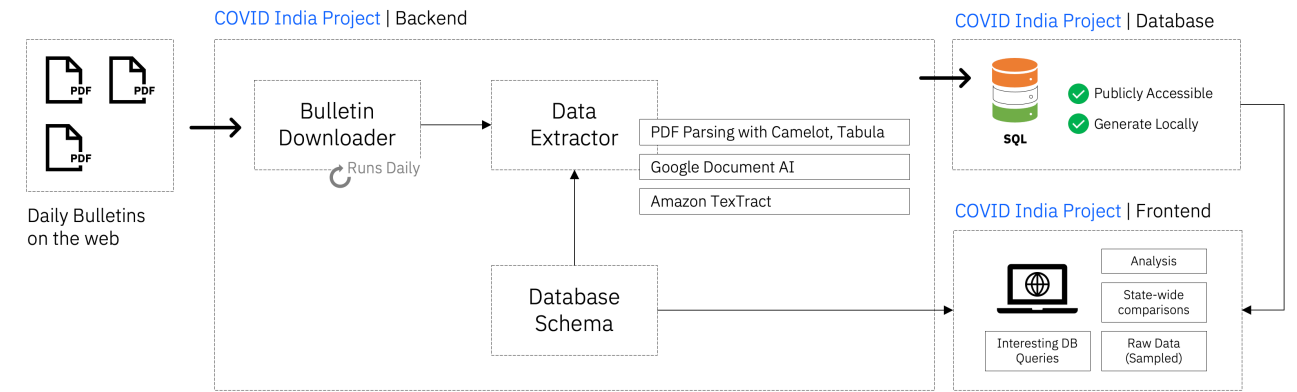


Figure 1. System Architecture for the COVID-19 India Project illustrating the data extraction pipeline from bulletins on the web to a publicly accessible SQL database updated daily.

2. System Overview

We segment the system into 3 major components: (a) the backend which is responsible for extracting data from health bulletins, (b) the database which stores the parsed structured data, and (c) the frontend which displays key analyses extracted from the parsed data. We describe each of these components in greater detail in the following sections.

2.1. The Backend

Since we aim to extract data from health bulletins published by individual states on their respective websites, there is no standard template that is followed across these data sources in terms of where and how the bulletin is published, and what and how information is included in these bulletins. To account for these variations, we modularize the system into the following 3 main components: a) bulletin download, b) datatable definition, and c) data extraction. We provide an overview of the system in Figure 1 and look at the three components in greater detail. The open-sourced code for the backend can be accessed at: <https://github.com/IBM/covid19-india-data>.

2.1.1. BULLETIN DOWNLOAD

The bulletin download procedure downloads the bulletins from the respective state websites to the local storage and maintaining the state of dates already processed. We use the BeautifulSoup¹ library to parse the states' websites and identify bulletin links and dates for download.

2.1.2. DATATABLE DEFINITIONS

Since each state provides different information, we define table schemas for each state by manually investigating the

¹<https://www.crummy.com/software/BeautifulSoup/>

bulletin (done once per state). We then use the free open-source SQLite² database to interface with the data extractor and store the data.

2.1.3. DATA EXTRACTOR

States typically provide the bulletins in the form of PDF documents. To extract information from them, we use a combination of classical PDF parsers and state of the art Machine Learning based extraction techniques:

1. **Classical PDF parsing:** Since a substantial amount of information in the bulletins are in the form of data tables, we use the Tabula³ and the Camelot⁴ Python libraries to extract these tables in the form of python data structures. While these libraries cover a lot of use cases, they do fail in certain edge case scenarios.
2. **Google Document AI:** For scenarios where Classical PDF parsing fails to identify data tables, we use Google Document AI⁵ API to extract this information. While the exact technologies behind this API are unknown, it uses computer vision and natural language processing techniques for document processing.
3. **Amazon TexTract:** While most of the information provided in the bulletins is text-based, some tables are provided as images that the prior two techniques cannot process. To extract them, we use Amazon TexTract⁶ APIs for table extraction from images.

To process information for a state, a separate data extractor

²<https://www.sqlite.org/index.html>

³<https://tabula.technology/>

⁴<https://camelot-py.readthedocs.io/en/master/>

⁵<https://cloud.google.com/document-ai>

⁶<https://aws.amazon.com/texttract/>

routine is used, which has access to all the three aforementioned APIs. Depending on the format of the particular bulletin, we utilize a combination of the three techniques to extract information.

2.2. The Frontend

The frontend or landing page for the project is generated automatically from the database schema and provides easy access to 1) the raw data (sampled at an appropriate rate to be loaded on the browser); and 2) pages for highlights and analysis based on SQL queries (such as those described in Section 3). The frontend for the project can be accessed at: <https://ibm.biz/covid-data-india>

2.3. The Database

The system described above runs daily and produces a SQL database that is publicly available for download. However, one can also use the source code to generate data customized with their own parameters, and deploy into their local systems. The database can be downloaded at: <https://ibm.biz/covid19-india-db>, and instructions on how to read data from the database can be found at: <https://github.com/IBM/covid19-india-data/wiki/How-to-access-the-data>

2.3.1. CURRENT STATUS

At the time of this writing, we have completely indexed information from two major Indian states.

1. **Delhi (DL):** The National Capital Territory of Delhi provides information on: (A) Daily Cases, (B) Testing, (C) Vaccination, (D) Hospital Infrastructure and Availability, and (E) Containment Zone. We index this information in six tables, which are described in Appendix A.1.
2. **West Bengal (WB):** The eastern state of WB provides information about the (A) number of daily cases at the state and the district level, (B) testing, (C) hospital infrastructure and occupancy, (D) counseling information, among others. We index this information in eight tables, which are described in Appendix A.2.

2.3.2. WORK IN PROGRESS

In addition to the states mentioned above, we have also identified seven other major states – Telangana (TG), Tamil Nadu (TN), Karnataka (KA), Kerala (KL), Madhya Pradesh (MP), Punjab (PB), and Uttarakhand (UK) – which provide a daily health bulletin with information that can be parsed and extracted. Since the release of the extraction package, OSS contributors are currently working on extending the data to these additional states.

3. Preliminary Analysis

In this section, we perform some preliminary analysis on the data collected from the health bulletins of Delhi and West Bengal. We would like to emphasize that some of these analyses (to the best of our knowledge) are the first such analyses available for the two states. However, these are still preliminary but provide an insight into the power of such data available to researchers interested in the subject. We provide the SQL queries to reproduce the results in this section in Appendix C.

3.1. Weekly Case Fatality Rate (CFR)

India has seen two major waves of COVID-19, with the second wave fuelled primarily by the Delta variant (Yang & Shaman, 2021) being more deadly than the first (Budhiraja et al., 2021; Gupta et al., 2021).

We aim to understand the difference between the two waves by computing the Weekly Case Fatality Rate as the ratio of total fatalities to total newly confirmed cases in a particular week. The charts for Delhi and West Bengal are presented in Figure 2. While the weekly CFR for the first wave seems to be comparable for the two states, there appears to be a stark difference in the numbers for the second wave. While the weekly CFR for Delhi increased from 4% to about 12% during the second wave, the numbers for West Bengal do not show any increase.

3.2. Percentage of RT-PCR tests

Currently, India uses the reverse-transcriptase polymerase-chain-reaction (RT-PCR) tests and the Rapid Antigen Tests (RATs) to detect COVID-19 cases. While RT-PCR tests are highly accurate and are considered gold-standard tests for detecting COVID-19 (Brihn et al., 2021), they are more expensive and time-consuming than the less accurate RATs. While the official advisory is to prefer RATs for surveillance in areas with known higher rate of transmission, and prefer RT-PCRs in other settings (Indian Council of Medical Research, 2020), there exists a discrepancy in how different states have been using the two testing methods (Cherian et al., 2021), and how this ratio affects the reported case results (Chatterjee, 2020).

The state government of Delhi has in the past been called out for over-reliance on RATs as opposed to the preferred RT-PCR tests (Sirur, 2020). Following this criticism, the government increased the share of RT-PCR tests. We compute this ratio of RT-PCR tests to total tests conducted in the state, and show it in Figure 2. As is evident, in 2020, less than 50% of the total tests conducted in the state were RT-PCR tests. However, starting 2021, and especially during the second wave of COVID-19 in India, this ratio increased to over 70%.

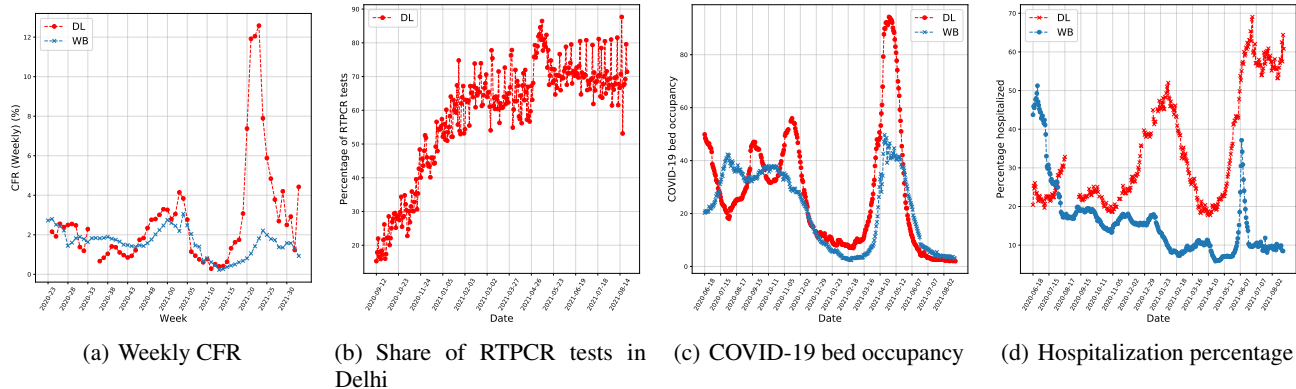


Figure 2. Preliminary analysis illustrating the depth of data available from the daily health bulletins.

3.3. COVID-19 bed occupancy

Both states (DL and WB) report the dedicated COVID-19 hospital infrastructure and occupancy information in their bulletins. Using these numbers, we compute the COVID-19 bed occupancy as the ratio of occupied beds to total, and plot these numbers in Figure 2. Similar to the results in Section 3.1, bed occupancy for Delhi shows a steep increase – reaching about 90% occupancy – during the second wave, while the occupancy for West Bengal does not show any significant difference during the two waves.

3.4. Hospitalization percentage

To treat COVID-19 patients, India adopted a two-pronged strategy of hospitalization along with home isolation, where patients with a mild case of COVID-19 were advised home isolation whereas hospitals were reserved for patients with more severe cases of COVID-19 (Varghese & John, 2020; Bhardwaj et al., 2021). We compute the hospitalization percentage as the ratio of the number of occupied hospital beds to the number of active cases. This is an estimate of how many of the currently active COVID-19 patients are in hospitals versus home isolation. The results are shown in Figure 2. The peaks we see for the two states relate to time periods after the respective wave has subsided, the minima and the subsequent rise in hospitalization relates to the onset of the particular wave.

4. Future work

The primary aim of this project is to extract as much information about the pandemic as possible from all available sources with the hope that providing this data in an easy and structured form to researchers will allow them to utilize such data (from one of the most populous and heavily COVID-affected countries in the world) in their research. We foresee two main areas of future work for this project:

1. In the immediate future, we aim to integrate information for states mentioned in Section 2.3 into the dataset. Additionally, the project currently relies on health bulletins alone to extract all the data. There are however other platforms where the authorities release data, such as Twitter and Government APIs (covid19india, 2020). We hope to integrate these additional sources of information in the dataset.
2. The Tamil Nadu (TN) bulletins has been partially parsed. The database schema and sample queries are provided in Appendix C. The classical APIs failed to parse some of the tables of TN for which we are currently looking at alternatives like Google Document AI. The TN bulletins provide information about RTPCR and district-wise information for bed occupancy and hospitalization or at-home quarantining. It also includes information about gender-based numbers for testing and the number of positive cases, the number of people who tested positive that traveled to the state using different modes (car, train, ships, or flights).
3. We anticipate that this data can be helpful in (A) validating or extending models developed for other countries to India (Friedman et al., 2021; Borchering et al., 2021), (B) developing pandemic models which integrate additional variables available in our dataset (Hethcote, 2000; Agrawal et al., 2021; Adiga et al., 2020; Bharduri et al., 2020), (C) Understanding various aspects of the pandemic (Ray et al., 2020; Ghosh et al., 2020), among others.

References

Adiga, A., Dubhashi, D., Lewis, B., Marathe, M., Venkatraman, S., and Vullikanti, A. Mathematical models for COVID-19 pandemic: a comparative analysis. *Journal of the Indian Institute of Science*, pp. 1–15, 2020.

- Agrawal, M., Kanitkar, M., and Vidyasagar, M. SUTRA: An approach to modelling pandemics with asymptomatic patients, and applications to COVID-19. *arXiv preprint arXiv:2101.09158*, 2021.
- Bhaduri, R., Kundu, R., Purkayastha, S., Kleinsasser, M., Beesley, L., and Mukherjee, B. Extending the Susceptible-Exposed-Infected-Removed (SEIR) model to handle the high false negative rate and symptom-based administration of COVID-19 diagnostic tests: SEIR-fansy. *medRxiv*, 2020.
- Bhardwaj, P., Joshi, N. K., Gupta, M. K., Goel, A. D., Saurabh, S., Charan, J., Rajpurohit, P., Ola, S., Singh, P., Bisht, S., et al. Analysis of Facility and Home Isolation Strategies in COVID 19 Pandemic: Evidences from Jodhpur, India. *Infection and Drug Resistance*, 14:2233, 2021.
- Borcherding, R. K., Viboud, C., Howerton, E., Smith, C. P., Truelove, S., Runge, M. C., Reich, N. G., Contamin, L., Levander, J., Salerno, J., et al. Modeling of future COVID-19 cases, hospitalizations, and deaths, by vaccination rates and nonpharmaceutical intervention scenarios – United States, April–September 2021. *Morbidity and Mortality Weekly Report*, 70(19):719, 2021.
- Brihn, A., Chang, J., OYong, K., Balter, S., Terashita, D., Rubin, Z., and Yeganeh, N. Diagnostic Performance of an Antigen Test with RT-PCR for the Detection of SARS-CoV-2 in a Hospital Setting—Los Angeles County, California, June–August 2020. *Morbidity and Mortality Weekly Report*, 70(19):702, 2021.
- Budhiraja, S., Indrayan, A., Aggarwal, M., Jha, V., Jain, D., Tarai, B., Das, P., Aggarwal, B., Mishra, R., Bali, S., et al. Differentials in the characteristics of COVID-19 cases in Wave-1 and Wave-2 admitted to a network of hospitals in North India. *medRxiv*, 2021.
- Chatterjee, P. Is India missing COVID-19 deaths? *The Lancet*, 396(10252):657, 2020.
- Cherian, P., Krishna, S., and Menon, G. I. Optimizing Testing for COVID-19 in India. *medRxiv*, pp. 2020–12, 2021.
- covid19india. Hornbill. <https://blog.covid19india.org/2020/06/15/hornbill/>, 2020.
- covid19india. Operations. <https://blog.covid19india.org/2021/08/24/operations/>, August 2021a.
- covid19india. When the curtains come down. <https://blog.covid19india.org/2021/08/07/end/>, August 2021b.
- Friedman, J., Liu, P., Troeger, C. E., Carter, A., Reiner, R. C., Barber, R. M., Collins, J., Lim, S. S., Pigott, D. M., Vos, T., et al. Predictive performance of international COVID-19 mortality forecasting models. *Nature Communications*, 12(1):1–13, 2021.
- Ghosh, K., Sengupta, N., Manna, D., and De, S. Inter-state transmission potential and vulnerability of COVID-19 in India. *Progress in Disaster Science*, 7:100114 – 100114, 2020.
- Gupta, N., Kaur, H., Yadav, P., Mukhopadhyay, L., Sahay, R. R., Kumar, A., Nyayanit, D. A., Shete, A. M., Patil, S., Majumdar, T. D., et al. Clinical characterization and Genomic analysis of COVID-19 breakthrough infections during second wave in different states of India. *medRxiv*, 2021.
- Hethcote, H. W. The mathematics of infectious diseases. *SIAM review*, 42(4):599–653, 2000.
- Indian Council of Medical Research. Advisory on Strategy for COVID-19 Testing in India. https://www.icmr.gov.in/pdf/covid/strategy/Testing_Strategy_v6_04092020.pdf, 2020.
- Priyanka Pulla. “There are so many hurdles.” Indian scientists plead with government to unlock COVID-19 data. <https://www.science.org/news/2021/05/there-are-so-many-hurdles-indian-scientists-plead-government-unlock-covid-19-data>, May 2021. *Science*.
- Ray, D., Salvatore, M., Bhattacharyya, R., Wang, L., Du, J., Mohammed, S., Purkayastha, S., Halder, A., Rix, A., Barker, D., Kleinsasser, M., Zhou, Y., Bose, D., Song, P. X. K., Banerjee, M., Baladandayuthapani, V., Ghosh, P., and Mukherjee, B. Predictions, role of interventions and effects of a historic national lockdown in India’s response to the COVID-19 pandemic: data science call to arms. *Harvard Data Science Review*, 2020 Suppl 1, 2020.
- Sirur, S. It isn’t just Delhi. Kerala, Bihar & UP also conduct more than 50% rapid antigen tests. <https://theprint.in/health/it-isnt-just-delhi-kerala-bihar-up-also-conduct-more-than-50-rapid-antigen-tests/550255/>, Nov 2020. *ThePrint*.
- Varghese, G. M. and John, R. COVID-19 in India: Moving from containment to mitigation. *The Indian journal of medical research*, 151(2-3):136, 2020.
- Yang, W. and Shaman, J. COVID-19 pandemic dynamics in India and impact of the SARS-CoV-2 Delta (B. 1.617. 2) variant. *medRxiv*, 2021.

A. Database schema

A.1. Delhi

We define the following 6 tables for the state of Delhi (see Figure 3):

1. `DL_case_info`: Daily cases, recovered, and fatality information
2. `DL_containment`: Information about the containment zones, Ambulance calls, and helpline calls
3. `DL_cumulative`: Cumulative positive, recovered, and fatalities data
4. `DL_patient_mgmt`: Hospital occupancy information
5. `DL_vaccination`: Vaccination information
6. `DL_testing_status`: RTPCR / Rapid Antigen test information table

A.2. West Bengal

We define the following 8 tables for the state of West Bengal (see Figure 4):

1. `WB_case_info`: Daily and cumulative new cases, discharge, and fatality information
2. `WB_district_cases`: Daily and cumulative case, discharge, and fatality information for all districts in West Bengal state
3. `WB_testing`: Testing information for the state
4. `WB_hospital`: Hospital infrastructure, ICU, hospital beds, and home isolation patient numbers
5. `WB_testing_labs`: Information regarding testing lab in the state
6. `WB_vaccination`: Vaccination information for the state
7. `WB_equipment`: PPE, N95 masks, and gloves availability in the state
8. `WB_counselling`: Tele consultations, ambulance calls, and general query information

A.3. Tamil Nadu

We define the following 10 tables for the state of Tamil Nadu (see Figure 5):

1. `TN_positive_cases_detail`: Detailed daily information of new cases, discharges, and fatalities for the state
2. `TN_cumulative_info`: Summary of daily new positive cases, discharges, and fatalities
3. `TN_district_details`: District-wise details of new cases, discharges, and fatalities
4. `TN_district_hospital_bed_details`: District-wise hospital infrastructure details
5. `TN_comorbidities_deaths`: Aggregate statistics of fatalities with and without comorbidities
6. `TN_airport_surveillance`: Arrival, Testing, and positivity details of passengers arriving through railways
7. `TN_airport_surveillance_details`: Airport surveillance record for each flight arriving in the state
8. `TN_incoming_people_till_yesterday`: Aggregated record of total people entering the state through various transportation means
9. `TN_railway_surveillance`: Arrival, Testing, and positivity details of passengers arriving through railways
10. `TN_seaport_surveillance`: Arrival, Testing, and positivity details of passengers arriving through the sea route

B. State health bulletin samples

B.1. Delhi

Bulletins for Delhi state can be found at the link: http://health.delhigovt.nic.in/wps/wcm/connect/doit_health/Health/Home/Covid19/Bulletin+September+2021

B.2. West Bengal

Bulletins for West Bengal state can be found at the link: <https://www.wbhealth.gov.in/pages/corona/bulletin>

B.3. Tamil Nadu

Bulletins for Tamil Nadu state can be found at the link: <https://stopcorona.tn.gov.in/daily-bulletin/>

COVID-19 India Dataset

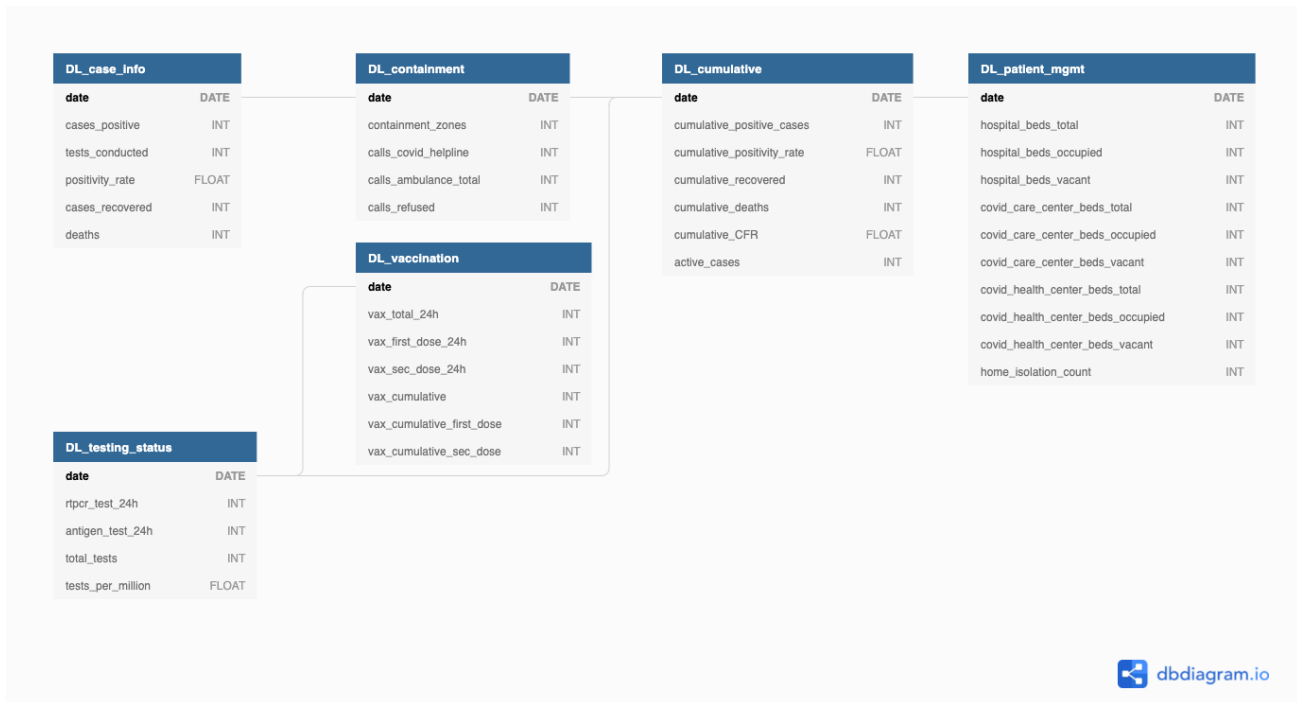


Figure 3. Delhi database schema

C. SQL queries for preliminary analysis

To run these queries would require access to the database file which can be found here: <https://ibm.biz/covid19-india-db>

C.1. Weekly Case Fatality Rate

```
WITH DL_CFR AS (
    select strftime('%Y-%W', date) as
        WeekNo ,
        SUM(cases_positive) AS
            total_positive,
        SUM(deaths) as total_deaths
    from DL_case_info GROUP by WeekNo
    ORDER BY WeekNo
),
WB_CFR AS (
    select strftime('%Y-%W', date) as
        WeekNo ,
        SUM(cases_new) AS
            total_positive,
        SUM(deaths_new) as total_deaths
    from WB_case_info GROUP by WeekNo
    ORDER BY WeekNo
)
SELECT
    DL_CFR.WeekNo,
```

```
DL_CFR.total_deaths * 100.0 /
    DL_CFR.total_positive,
WB_CFR.total_deaths * 100.0 /
    WB_CFR.total_positive
FROM DL_CFR
JOIN WB_CFR ON DL_CFR.WeekNo == WB_CFR.
WeekNo
ORDER BY DL_CFR.WeekNo
```

C.2. Share of RT-PCR tests in Delhi

```
Select D1.date, D1.rtpcr_test_24h *1.0/
    D2.tests_conducted AS rtpcr_ratio
FROM DL_testing_status D1 JOIN
DL_case_info D2 ON D1.date == D2.
date WHERE rtpcr_ratio is NOT NULL
```

C.3. COVID-19 bed occupancy

```
SELECT D1.date, D2.
    hospital_beds_occupied * 100.0 / D2
    .hospital_beds_total AS
    DL_occupancy,
W1.covid19_bed_occupancy AS
    WB_occupancy
FROM DL_case_info D1
JOIN DL_patient_mgmt D2 ON D1.date ==
    D2.date
```

COVID-19 India Dataset

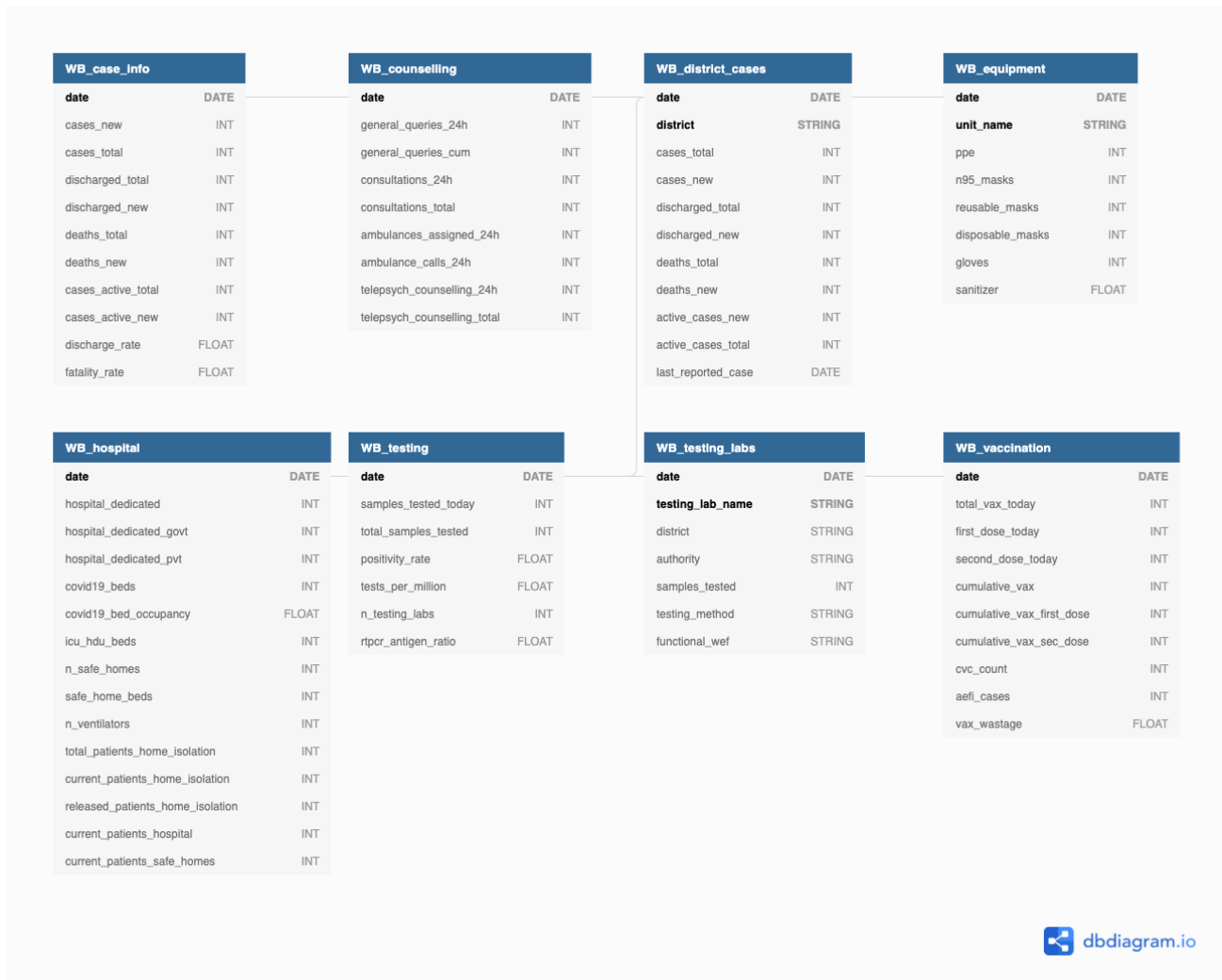


Figure 4. West Bengal database schema

```
JOIN WB_hospital W1 ON D1.date == W1.
    date
ORDER BY D1.date ASC
```

```
JOIN WB_hospital W2 ON D1.date == W2.
    date
ORDER BY D1.date ASC
```

C.4. Hospitalization percentage

```
SELECT D1.date, D2.
    hospital_beds_occupied * 100.0 / D1
    .active_cases AS DL_hosp,
(W2.covid19_bed_occupancy * (W2.
    covid19_beds + W2.icu_hdu_beds)) / (
    W1.cases_active_total) AS WB_hosp
FROM DL_cumulative D1
JOIN DL_patient_mgmt D2 ON D1.date ==
    D2.date
JOIN WB_case_info W1 ON D1.date == W1.
    date
```


COVID-19 India Dataset

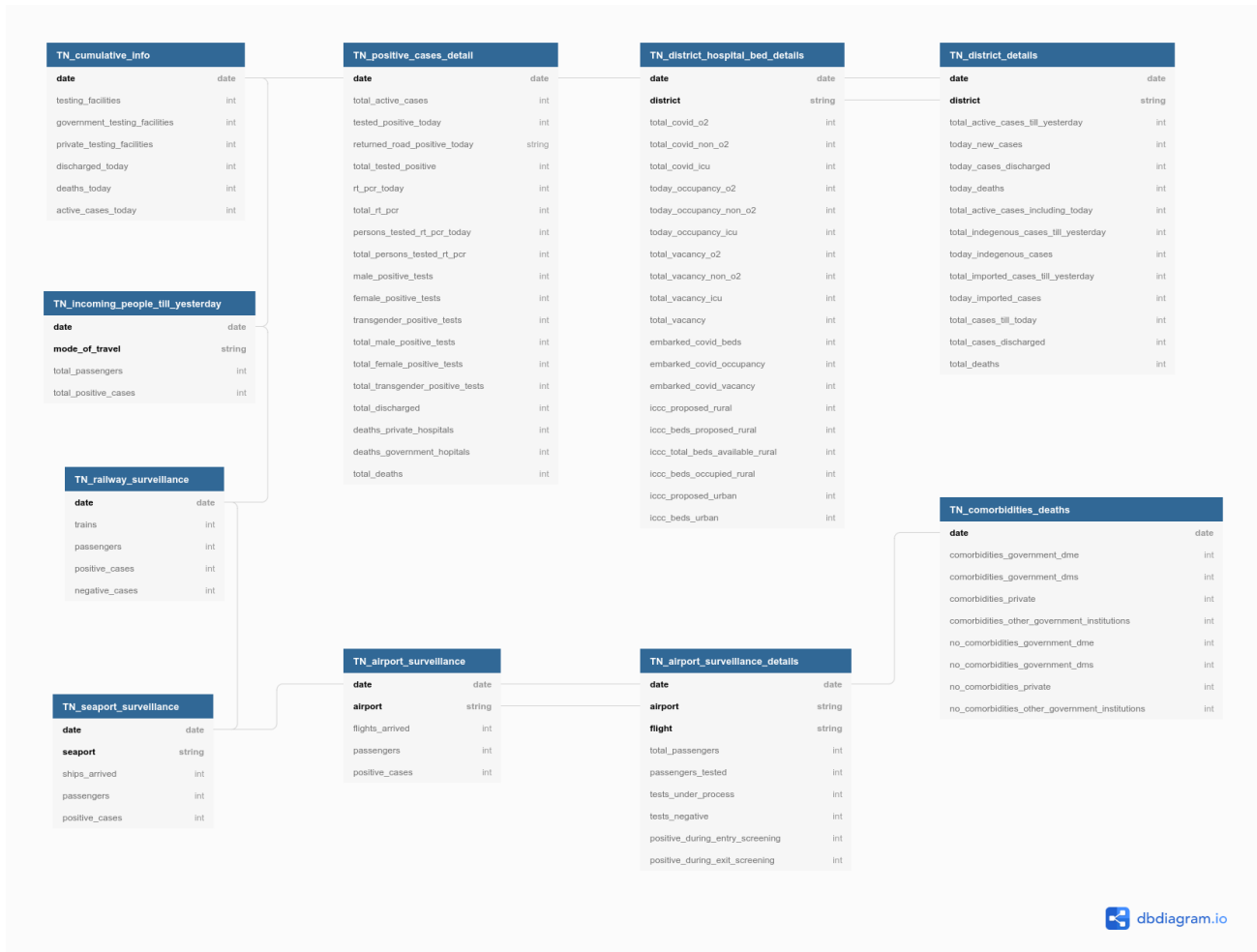


Figure 5. Tamil Nadu database schema