




Large-Scale Contact Tracing, Hotspot Detection, and Safe Route Recommendation

Chandresh Kumar Maurya^(✉), Seemandhar Jain, and Vishal Thakre

IIT Indore, Indore, India
cse180001062@iiti.ac.in

Abstract. Recently, the COVID-19 pandemic created a worldwide emergency as it is estimated that such a large number of infections are due to human-to-human transmission of the COVID-19. As a necessity, there is a need to track users who came in contact with users having travel history, asymptomatic and not yet symptomatic, but they can be in the future. To solve this problem, the present work proposes a solution for contact tracing based on assisted GPS and cloud computing technologies. An application is developed to collect each user's assisted GPS coordinates once all the users install this application. This application periodically sends assisted GPS data (coordinates) to the cloud. To determine which devices are within the permissible limit of 5 m (tunable parameter), we perform clustering over assisted GPS coordinates and track the clusters for about t mins (tunable parameter) to allow the measure of spread. We assume that it takes around 3–5 mins to get the virus from an infected object. For clustering, the proposed M-way like tree data structure stores the assisted GPS coordinates in degree, minute, and second (DMS) format. Thus, every user is mapped to a leaf node of the tree. The crux of the solution lies at the leaf node. We split the “seconds” part of the assisted GPS location into m equal parts (a tunable parameter), which amount to d meter in latitude/longitude. Hence, two users who are within d meter range will map to the same leaf node. Thus, by mapping assisted GPS locations every t mins (usually $t = 2:5$ mins), we can find out how many users came in contact with a particular user for at least t mins. Our work's salient feature is that it runs in linear time $O(n)$ for n users in the static case, i.e., when users are not moving. We also propose a variant of our solution to handle the dynamic case, that is, when users are moving. Besides, the proposed solution offers potential hotspot detection and safe-route recommendation as an additional feature, and proof-of-concept is presented through experiments on simulated data of 2/4/6/8/10M users.

Keywords: COVID-19 · Contact-tracing · M-way like tree · Clustering · Assisted GPS

1 Introduction

Contact tracing is the problem of identifying users who have come in contact with a user within a certain distance for a specific amount of time. This problem came

into the limelight during the COVID-19 pandemic, which is thought to spread from humans to humans. The main problem is that users are asymptomatic and can pass the virus to other users unknowingly for weeks or months. In such a case, it is challenging to identify users who came in contact with a particular user, given the users dynamic and complex movement patterns. Another prevalent problem during the pandemic situation is that users are interested in knowing the hotspot areas to avoid them while visiting. Therefore, there is a growing need for a recommendation of a safe route for travel.

Current solutions for contact tracing problems can be categorized into (a) human-based approach, (b) Bluetooth (BT)-based approach, and (c) GPS-based approach. In the first approach, police personnel are deployed to follow the user's movement trail and find out users who were present in close proximity of the particular user for a specific duration. This approach is *costly, cumbersome, and time-consuming* [11]. In the second approach, each user is asked to install an app such as India's Aarogya Setu app or Google-Apple privacy-preserving app. Such apps use BT signals for exchanging data between devices. There are certain limitations of these apps. For example, the limitations of Aarogya Setu are: (a) it does not tell if we were in contact with a user a week or month ago who recently got diagnosed with COVID-19, (b) it also does not tell which areas are hotspots and should be avoided, (c) it is based on self-assessment only, and hence reliability is a concern.

The solutions for hotspot detection are currently based on the government's rules and regulations. For example, sites like COVID tracker¹ mark geographical areas as hotspot-based on the directions received from the government. They do not provide real-time hotspot information, such as alerting the user that they are moving through hotspot/containment zones specially in villages where monitoring is difficult. Further, the existing solutions lack the feature of recommending safe routes for travel, such as Google map.

To overcome the limitations of the existing solutions, we propose an assisted GPS location-based solution that receives the assisted GPS coordinates of the user in the backend and without forcing the user to keep on Bluetooth time and again. This solution does not require any user interventions and sends the assisted GPS coordinates to the cloud periodically. Received coordinates are mapped to the proposed **M-way like tree data structure** in degree, minute, and second (DMS) format. Once the mapping is over, the proposed clustering algorithm executes to track the history of users who spent at least t minutes with other users. **Our solution is scalable and can find users who happen to make a contact event in 17 mins among a population of 10M users.**

Major differences between our solution based on assisted GPS and BT-based solution for contact tracing are: (i) our approach can provide location information where the user got **potentially infected** whereas BT-based approaches fail to do so since they do not provide location information. Why this information may be crucial is that let us suppose that the user visits a dairy milk shop every day besides tons of other places. On one day (s)he got COVID +ve.

¹ <https://www.COVIDhotspots.in/?city=Delhi>.

In this case, (s)he might be interested in knowing where (s)he got the virus so that they can inform their family members to be vigilant. In this scenario, the BT-based approach has no clue of the location, (ii) assisted GPS-based solution is centralized where BT-based one is decentralized, (iii) our solution is more robust in the sense that BT-based solution communicates with **many heterogeneous devices**. In contrast, ours does not, (iv) BT always needs to be turned on, which usually user forget or do not turn on for power saving purposes. As for assisted GPS is concerned, once the app is installed, it will keep using the location sharing in the background and does not ask the user to turn it on time and again. This option is not available in BT-based solutions currently, (v) the scalability of BT-based apps is also a concern since these apps can communicate to only eight other BT devices at a time. Further, note that we are using *assisted* GPS which is more reliable than GPS, where location information is calculated using satellite signals, mobile towers, and wifi beacons. Therefore, our solution will work in indoor settings as well. From now on, whenever GPS is mentioned will mean assisted GPS. To emphasize to our readers that our solution offers an alternative to BT-based solutions and does not replace them.

In short, we make the following contributions:

- Propose a solution for contact tracing that runs in the time linear to the numbers of users i.e. $O(n)$. the space complexity of the solution is also linear in the number of user identifiers, i.e., $O(n)$.
- Our solution can handle static as well as dynamic case. That means if users stay at the same location (static case) or the users are moving together (dynamic case).
- The proposed solution relies only on location sharing information, which is the property of most android applications such as Google news, maps, etc.
- Our solution can additionally provide hotspot information (areas of large gatherings) and recommend a safe route for travel.

2 Related Works

Contact tracing and regular monitoring of infectious diseases (COVID-19) are essential for public health. Utilizing emerging technologies such as remote sensing, internet-based surveillance, telecommunications, infectious disease modeling, Global Positioning System (GPS), IoT devices, and mobile phones, such contagious diseases can be monitored, prevented (by generating alarms after real-time predictions), and controlled [2].

Few research works [1,12] have introduced the stochastic model used to reduce a deterministic approach for attaining the fundamental dynamics of the epidemics and other associated measures. Mostly, previous models deal only with generic networks. To measure the precision of the trace contact models, there is a need to account for the network of contacts. For example, Huerta and Tsimring [8] present a stochastic model to estimate the effect of random screening and contact tracing as part of the epidemic control strategy in complex networks. This work indicates that a significant outbreak can significantly be reduced via

tracing contacts at a low additional cost. A similar approach is also used by Farrahi et al. [5].

Ferretti et al. [6] state that isolation and contact tracing are being practiced but not preventing the COVID-19, which is why the same high number of asymptomatic infected individuals remain undetected, and cases continue to spread. Therefore, the authors propose mobile apps to trace previous mobile contacts and show mathematically that epidemics can remain even when not all population use the application. Hellewell et al. [7] also propose similar inference through the simulated model. In most scenarios, highly effective case isolation and contact tracing are sufficient to control a new outbreak of COVID-19 within three months; however, only 79% of the contacts are traced. Nevertheless, such conditions create smartphone-based tracing, which is far from a realistic solution. One of the first efforts through mobile phones to trace the contacts was the FluPhone application developed by Cambridge University [18] which uses Bluetooth-based wireless signals to estimate the physical contact. It also asks the users to report on flu-like symptoms to appraise the risk of infections. Next, the Singapore Government also developed a mobile app called Trace-Together for COVID-19². This also uses Bluetooth technology, and it had already been utilized to control the disease spread. Few similar works also focused on privacy as the Pan-European Privacy-Preserving Proximity Tracing (PEPP-PT) [15]. Finally, Google and Apple have teamed up to design, develop and integrate a similar approach into the Android operating Systems and iOS. The proposed solution is more efficient and ubiquitous to users. Aarogya Setu app³ from the Indian government works on location sharing and Bluetooth. Secondly, it tells us that we are in contact with a COVID +ve person based on Bluetooth proximity. There are certain limitations of Aarogya Setu. For example, it does not tell if we were in contact with a person a week or month ago who recently got diagnosed with COVID-19. It also does not reveal which areas are hotspots and should be avoided. It is based on self-assessment only, and hence reliability is a concern. Recently, Mahapatra et al. [10] present a digital contact tracing solution based on a dynamic graph streaming algorithm. Concretely, they use an index-based adjacency list to store graphs whose nodes are users and edges are close contacts. However, they do not mention if they use Bluetooth or GPS to find direct and indirect close contacts. The proposed solution runs in $O(q^L|I|)$ where q , L , $|I|$ denote the number of contacts, level of tracing, and the number of infected users, respectively. All Bluetooth (even GPS) based apps for contact tracing suffer from a severe drawback: they may be subjected to high “false alarm”. For example, when two users are standing opposite side of a wall, contact tracing apps will signal that they are within the infection-proximity [16]. There have been reports of delayed notification (1+ days) to individuals regarding contacts with COVID +ve patients. Such delays can create depression in users later on. Further, many of these apps do not show hotspot areas nor recommend a safe route. To ameliorate/minimize *some of the above problems* and provide additional features, our

² <https://www.tracetogether.gov.sg/>.

³ <https://www.mygov.in/aarogya-Setu-app/>.

solution relies on location sharing via GPS/assisted GPS such as magnetic fingerprints and the BLE beacon RSSI (Received Signal Strength Identifier) nearby, for localization [14] in an indoor setting. Unlike existing solutions, our approach provides *potential hotspot* information as well as *safe route* recommendation which is presented in Sects. 4 and 5 respectively.

3 Contact Tracing

As discussed in the introduction section, our solution for contact tracing relies on the location sharing data, i.e., GPS location of the users collected through our app. The naive solution will be to compute the Euclidean distance (or *Haversine distance* as shown in (1), where ϕ_i and λ_i for $i \in \{1, 2\}$ are lat/long of two users, if one wants to be more exact) every t mins for finding the contact event.

$$d = 2r \arcsin\left(\sqrt{\sin^2\left(\frac{\phi_2 - \phi_1}{2}\right) + \cos(\phi_1) \cos(\phi_2) \sin^2\left(\frac{\lambda_2 - \lambda_1}{2}\right)}\right) \quad (1)$$

A contact event is defined as when two users are within a distance of d meter for at least t mins. The naive solution for contact tracing requires computing pairwise distances and has the complexity of $O(n^2)$ for n users and hence not feasible for large n . Therefore, we discuss a data structure which is M-way *like* tree data structure for storing GPS locations (latitude/longitude) as shown in Fig. 1. In other words, we map latitude/longitudes to the proposed tree as discussed later in details. Further, latitude and longitudes are mapped to two separate trees for parallel computation of the contact distances. How do we arrive at the direct distance between two users from latitude and longitude distances is presented in detail in Sect. 3.2. Note that the M-way like tree is not the same as k-d tree which is a binary tree and requires comparison for insertion. Whereas M-way like tree involves look-up operation and no comparison. In this spirit, it is a kind of multi-level hashing.

A GPS coordinate that consists of latitude and longitude is usually expressed in the degree, minute, and second format (DMS), e.g., $(28^\circ 50' 30.12462'')$. For latitude and longitude, degree ranges in $[0, 1, \dots, 89]$ and $[0, 1, \dots, 179]$ respectively while minutes and seconds vary in $[0, 1, \dots, 59]$ for both which are shown in the top 3 gray levels in Fig. 1. For notational brevity, we do not show degree ($^\circ$), minute ($'$) and seconds ($''$) symbols in the fig. The last gray level is the partition of the seconds into m equal parts. As the most interesting property of M-way like tree data structure is the partition at the leaf, i.e., partition every *second* of GPS (which is equal to approx. 30 m) to m equal parts so that $m = 30/d$. Here, d is our contact distance which is a tunable parameter ($d = 5$ in our case). The benefit of the partition is that all users within a distance of d meters on latitude will fall into the same leaf (bucket) and form natural clusters. For example, consider $d = 5$ and hence $m = 6$, then the intervals for the fractional part of the seconds (which varies in $[0, 1]$) at the leaf node will be $[0, 1/6], [1/6, 2/6], \dots, [5/6, 1]$. Users u_1 and u_2 whose latitudes are $(28^\circ 50' 30.12462'')$ and $(28^\circ 50' 30.05462'')$ respectively will map to the same leaf node because their fractional part of

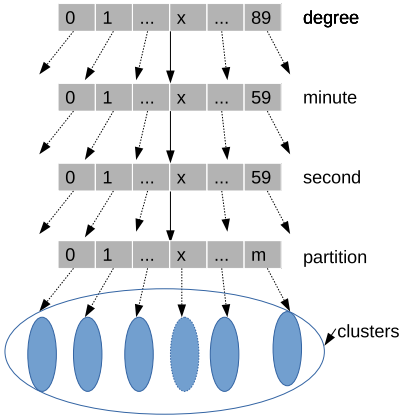


Fig. 1. M-way tree like data structure for storing coordinates



Fig. 2. Example to handle cases when users in the neighboring clusters may be within contact event distance d and not captured during mapping of lat/long to the M-way like tree.

the second (.12462 and .05462) will map to the same leaf and thus are in the same cluster. The approach presented above misses users who are within contact event distance d but fall in the neighboring clusters. For example, users u_1 and u_3 whose latitudes are ($28^\circ 50' 30.12462''$) and ($28^\circ 50' 30.17462''$) will map to different clusters but are within the contact event distance $d = 5m$. An example of such a case is shown in Fig. 2 with four users. Out of these four users, two fall in one cluster, and the other two fall in a different neighboring cluster. To capture missed cases during the mapping of lat/long to the tree, we perform pairwise distance computation. For example, we need 4 distances computed ($\{u_1, u_3\}, \{u_1, u_4\}, \{u_2, u_3\}, \{u_2, u_4\}$) in the worst case (worst case occurs when half of the users fall in one cluster and the other half fall in the neighboring cluster). Had we not mapped users to the tree, $\binom{n}{2}$ computations are required with quadratic complexity $O(n^2)$ and not preferred for large n . To see this in Fig. 2, instead of 6 pairwise distance computations in our running example, we needed only 4 pairwise distance computations in the worst case and hence saving two distance computations. Such a saving becomes paramount for large n .

We collect the GPS locations using our in-house app and send it to the cloud every $t/2$ mins for tracking users at least for t mins ($t = 5$ mins in our experiments). That means we map GPS locations to M-way like tree structure every $t/2$ mins. We have two trees for every interval of $t/2$ mins: one at time step $t/2$ and the other at t . At this step, we perform the intersection of leaves, and the ID of the users found in the intersection is stored. This method guarantees that two users in the same cluster have been in contact for at least $t/2$ mins. It also ensures that it will capture all users who have spent at least t mins together within d meter distance. One limitation is that it might miss users who have been in contact for 3 or 4 mins. Since our goal is to find all users who came in contact for at least 5 mins or more, we can ignore such corner cases.

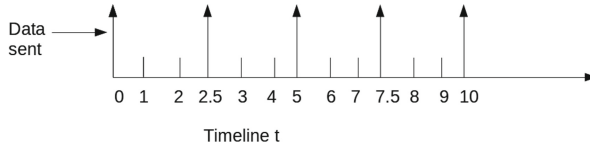


Fig. 3. Timeline to show when to send GPS coordinates to the cloud

3.1 Intuition Behind $t/2$ Mins

It is not necessary to send GPS coordinates every $t/2$ mins. However, to guarantee that two users make contact for at least t mins for the infection to occur, it is necessary to send GPS data every $t/2$ mins as shown in Fig. 3 for $t = 5$ mins. We can see in the Fig. 3 that if two users met at time $t_1 = 1$ and remain in contact until $t_2 = 6$ ($\Delta t = t_2 - t_1$) then we can see that we have sent their data two times (at $t = 2.5$ and $t = 5$ mins). As a result, we can capture the contact event. On the other hand, if we are sending data every 5 mins, it is not possible to capture the aforementioned contact event. This approach guarantees that all users whose contact event duration is at least 5 mins will be definitely captured. However, it may miss contact events of duration 3 or 4 mins.

3.2 How Lat/long Distances Map to Circular d m?

We want to compute the distance $d = AB$ as shown in the Fig. 4. That is, all users at a distance d from user A . For small distances ($d < 10$ or 15 m in our case), we can safely assume that the points A and B lie on Euclidean coordinate system. Therefore, we can draw a simple circle at A whose radius is d (otherwise, it is **great circle**).

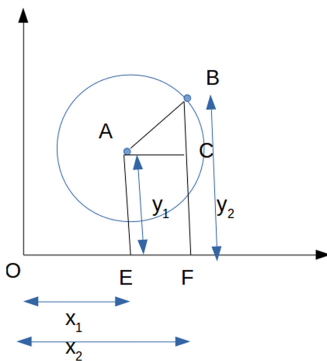


Fig. 4. Mapping lat/long to circular distances.

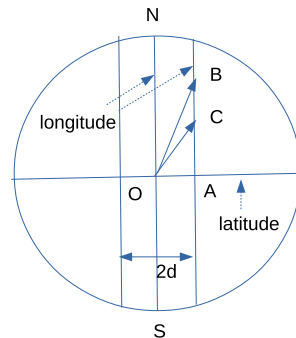


Fig. 5. Figure showing the case when many users are within contact distance d along latitude but can be far away along longitude.

Main Challenge: computing d is easy (we can take Euclidean distance or Haversine to be more precise). But, Euclidean computing distance for a large number of users (n) is computationally challenging. It takes $O(n^2)$ for pairwise distances. When $n = 1\text{M}$, it takes 10^{12} distances to compute and storing them requires $O(n^2)$ space. Thus, pairwise Euclidean distance computation is impractical since we want a scalable solution where n can be as large as 1.3B. Therefore, our main idea is to avoid computing d directly and instead compute AC and BC distances, which are lat/long distances between A and B .

Pythagorean theorem says that $AB^2 = AC^2 + BC^2$. Thus, if we can compute lat/long distances, we can find d . For $d = 5$ m. We can find either $(\text{lat}, \text{long}) = (3 \text{ m}, 4 \text{ m})$ or $(4 \text{ m}, 3 \text{ m})$. That means we find all users who are 3 m away from user in question on latitude and 4 m on longitude or vice-versa. Choosing $(\text{lat}, \text{long}) = (3 \text{ m}, 4 \text{ m})$ gives a more accurate solution since computing distance along longitude depends on the accuracy of the distance along latitude [17]. Further, notice that we map latitude/longitude of a user to two separate M-way like trees. That means we create a tree for mapping latitude and a tree for longitude. Such trees are created every $t/2$ mins, and users making contact events are found by the intersection of leaf nodes. More details of the procedure is described in the next section. Approximating diagonal distance via distances along latitude/longitude does not always hold. In the previous example, when $(\text{lat}, \text{long}) = (3\text{m}, 4\text{m})$, we are finding users who are within 3 m (4 m) along latitude (longitude) and clustering them through mapping their latitude (longitude) on the M-way like tree data structure. However, the catch here is that there is a possibility that users within the contact distance along latitude but far away along longitude (or vice versa) are mapped to the same leaf node. As shown in Fig. 5, users at locations B and C will fall in the same leaf node when we map longitudes to the tree; however, locations B and C can be miles away along latitude (parallel lines passing through B and C but not shown in the figure). This can degenerate into a case when many users are mapped to the same leaf node of either tree. Such a situation can be avoided by running our proposed methodology over each city instead of the whole country for contact tracing. What will happen is that the probability that many users align to the same latitude (longitude) simultaneously for a short duration will be very low due to population dynamics. Therefore, we assume that a leaf node will not have many users mapped to it in the worst case. In the next section, we present solutions when the users stay at the same location for at least t mins (called the static case) and moving (called dynamics case). Though dynamic case subsumes the static case, we bifurcate the mapping into two cases for speed-up purposes.

3.3 Static Case

A case is static when users who participate in a contact event are not moving at least for t mins. A static case can be easily handled by mapping lat/long to the M-way like tree every $t/2$ mins. The intuition is that two users who stay within a distance of d m for at least t mins will map to the same leaf node of the M-way like tree. We then find the intersection of the two M-way like trees so

obtained (intersection of two corresponding leaves to be more precise since users are static). Users in the intersection will be our output candidates (see Fig. 6 where leaves from two trees are shown facing each other).

Illustration of the static case: In the Fig. 6, lat/longs are mapped to two trees as discussed before. The user u_1, u_2 and u_3 were at a location at time t (top leaf node). After $t + \Delta t$ ($\Delta t = 2.5$ mins), users u_1, u_2 and u_3 were found at the same location resulting in the leaf node shown in the bottom. When the intersection is performed, u_1, u_2 , and u_3 are found to be present at the same location. That means, user u_1, u_2 and u_3 participated in a contact-event. The entire process of the static case is described step by step in Algorithm 1.

Algorithm 1: Static Case. $[]$ denotes indexing in an array/dictionary

Input: T_t - M-way tree at time t , T_{t+1} - M-way tree at time $t + \Delta t$
Output: CT - Contact Tracing Vector

```

1 for  $d$  in degrees do
2   for  $m$  in minutes do
3     for  $s$  in seconds do
4       for  $p$  in partitions do
5          $ct = T_{t+1}[d][m][s][p] \cap T_t[d][m][s][p]$ 
6         Insert  $ct$  in CT
7         for  $n$  in  $[p - 1, p + 1]$  do
8           if  $T_t[d][m][s][n]$  and  $T_t[d][m][s][p]$  has peoples who are less
           than  $5mt$  apart then
9              $end\_case_t = findpairs(T_t[d][m][s][n], T_t[d][m][s][p])$ 
10          if  $T_{t+1}[d][m][s][n]$  and  $T_{t+1}[d][m][s][p]$  has peoples who are
           less than  $5mt$  apart then
11             $end\_case_{t+1} = findpairs($ 
               $T_{t+1}[d][m][s][n], T_{t+1}[d][m][s][p])$ 
12          Insert  $end\_case_t \cap end\_case_{t+1}$  in CT
13 return CT
```

3.4 Dynamic Case

The dynamic case is when users move together in close proximity (within a radius of d m). Due to the movement, their lat/log are continuously changing, resulting in their mapping to two different (non-corresponding) leaf nodes in the M-way like tree. To find out where users in close contact map to the leaf require performing the intersection of each leaf node against all nodes in the latter tree. This is computationally intractable, requiring $O(n^2)$ intersection. Therefore, we solve this issue via an approximation. First, we discuss when users are walking on foot (movement in a car can be handled similarly). On avg, a pedestrian speed is 1.4 m/s. That means, in 5 mins, they walk around 420 m which maps

to 14" in lat/long. Therefore, we perform the intersection of each leaf against all leaves in the latter tree which are 14" left and right of the corresponding leaf as shown in the Fig. 7 (left/right is considered since users can move left side or right side of the current latitude. Similarly, up/down movement might happen for longitude). In either case (static/dynamic), once we find users who made contact along the latitude and longitude, we need to compute the distance AB (the actual contact direction) as discussed in Sect. 3.2. To build the final contact list, we proceed as follows. Firstly, we create a set of pairwise contacts from the contact list found by the intersection operation as mentioned previously. This step is repeated for the longitude contact list as well. Secondly, the two sets are intersected to generate the final contact list, i.e., the list of users who actually made a contact event along the direction AB. This operation cost $O(n)$ (due to sets being a hashmap) where the hidden constant is the length of the longest contact list, which is bounded above by a constant since a user can not meet more than a certain number of users in t ($=5$) mins. From the contact list, we can easily create a **contact vector** of each user and store in an adjacency matrix as discussed in the next section. The entire process of the dynamic case is described step by step in Algorithm 2.

Algorithm 2: Dynamic Case. $[]$ denotes indexing in an array/dictionary

Input: T_t - M-way tree at time t , T_{t+1} - M-way tree at time $t + \Delta t$

Output: CT - Contact Tracing Vector

```

1  $q = 14$                                 # 420 meter = 14"
2 for  $d$  in degrees do
3   for  $m$  in minutes do
4     for  $s$  in seconds do
5       for  $p$  in partitions do
6          $p1 \leftarrow []$ 
7          $p2 \leftarrow []$ 
8         for  $n$  in  $[p - q, p + q]$  do
9            $ct = T_{t+1}[d][m][s][n] \cap T_t[d][m][s][p]$ 
10          Insert  $ct$  in  $CT$ 
11          for  $n1$  in  $[n - 1, n + 1]$  do
12             $end\_case_t =$  find peoples who are less than 5mt apart in
13               $T_t[d][m][s][n1]$ 
14               $end\_case_{t+1} =$  find peoples who are less than 5mt apart
15                in  $T_{t+1}[d][m][s][n1]$ 
16              Insert  $end\_case_t$  in  $p1$ 
17              Insert  $end\_case_{t+1}$  in  $p2$ 
18              Insert  $p1 \cap p2$  in  $CT$ 
19 return  $CT$ 

```

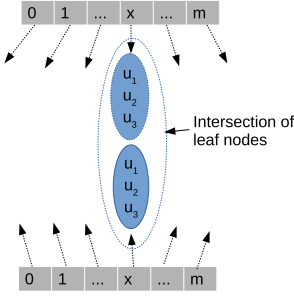


Fig. 6. Finding users in close proximity through intersection of corresponding leaves of two M-way like tree in the static case

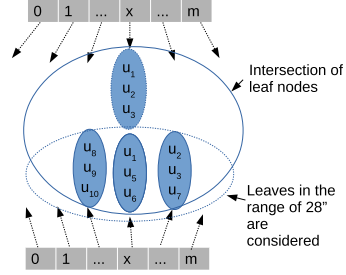


Fig. 7. Finding users in close proximity through intersection of corresponding leaves of two M-way like tree in the dynamic case

4 Potential Hotspot Detection

Our methodology can be used to detect *potential* hotspot areas in a city. Note that there is no unanimous definition of what constitutes a hotspot in epidemiology. As per [9], hotspot can be defined in three ways: (a) as areas of disease elevated occurrence or risk, (b) as areas of frequent disease emergence or reemergence, and (c) an area of elevated transmission efficiency. We take the third definition that is based on the assumption that large gatherings can stimulate the rate of transmission. Therefore, we identify areas in a city where large gatherings are happening along with COVID +ve cases and predict those areas as potential hotspots. The potential hotspot information can help safeguard users and discourage them from freely moving in those areas. Note that our method does declare gatherings as *potential hotspot* if there are no +ve cases in that gathering (cluster). Our assumption is that either COVID +ve user is surrounded by normal users who are not quarantined, and in some cases where COVID +ve user might also be moving, say coming back from a testing center. So, hotspot may contain COVID +ve users who are moving.

In order to locate areas of a potential hotspot, we proceed as follows. Contact list of each user obtained from the contact tracing methodology presented in Sect. 3 from the last 14 d is stored in the form of an adjacency matrix as shown in Eq. 2.

$$\begin{matrix} & u_1 & u_2 & \dots & u_n \\ \begin{matrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{matrix} & \begin{bmatrix} 1 & 1 & \dots & 1 \\ 0 & 0 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 1 \end{bmatrix} \end{matrix} \tag{2}$$

(In actual implementation, it can be implemented using a hashmap where keys are user ids and values are queue/list which stores the ids of users in the order the contact happens). We call each row as a **contact vector** for that user. We

can easily maintain the contact vector by adding /deleting from the queue/list the new/old contacts. Note that the adjacency matrix stores users who came in contact with other users. It does not store the COVID +ve users. In other words, we rely on official data from government authority to declare a user as a COVID +ve based on the diagnostic report. Once we receive such information, we can mark that user as a COVID +ve by setting a bit along that row or by maintaining another bit vector which is indexed by the user id as shown in Fig. 8.

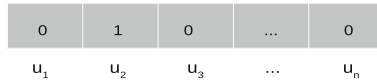


Fig. 8. Bit vector to store COVID +ve users

To detect a *potential hotspot* area, we provide a reference GPS location. This can be a city location as searched by users interested in knowing hotspot areas. The reference GPS is mapped to the lat/long tree created in that period. Assuming a 10 km radius from the reference point though, this can be dynamically adjusted according to the city area. Firstly, all the users are identified who are within the 10 km radius from the reference point. The 10 km of distance maps to around 5.39 min in lat/long. Therefore, all the users in the leaves, which are 5.39' left and right of the reference point, are identified. The user ids of such users are mapped to the bit vector (Fig. 8) to find out users who are COVID +ve within our search area (10 km, for example). Using the adjacency matrix, we also identify users who came in contact with users who are found COVID +ve in the search area. Such users might form *susceptible users*. Further, user ids of the susceptible users are cross-tallied in the search area (since users who came in contact with the COVID +ve user in the search area t days ago may have moved to a different city). Finally, we have built a set of COVID +ve users and suspected users. We can mark COVID +ve and susceptible users on the city map, and if the number of COVID +ve users exceeds a threshold, it can be marked as a *potential hotspot city*.

5 Safe Route Recommendation

Based on the potential hotspot areas, we can recommend a safe route for travel. Users currently use google map/Apple Maps as a primary navigation medium for traveling from one place to another. However, these maps have limited information about whether a route is safe for travel. For example, unless someone marks specifically that a particular route is blocked or traffic is slow, current navigation systems have limited information to tell users in advance that route is not safe. This situation was observed during the recent exodus of people from red zones such as Mumbai/New Delhi. People were following google maps to travel to their hometowns. However, the map took them to other hotspot cities

Algorithm 3: Safe Route Recommendation Algorithm

Input: G - Graph, S - Source, D - Target, $HOTSPOT$ - list of hotspot nodes.
Output: Path from S to T

```

1 for each vertex  $V$  in  $G$  do
2    $distance[V] \leftarrow \infty$ 
3    $last[V] \leftarrow None$ 
4   if  $V \neq S$  then
5     add  $V$  to Priority Queue  $Q$  (Using Min-Heap)
6 while  $Q$  is Not Empty do
7    $U \leftarrow$  Extract min From  $Q$ 
8   for each unvisited neighbour  $V$  of  $U$  do
9     if  $V$  not in  $HOTSPOT$  then
10       $temporary\_Distance \leftarrow distance[U] + edge\_Weight(U, V)$ 
11      if  $temporary\_Distance < distance[V]$  then
12         $distance[V] \leftarrow temporary\_Distance$ 
13         $last[V] \leftarrow U$ 
14 if  $distance[D]$  is  $\infty$  then
15   return  $NoPath$ 
16 return  $distance[D]$ 

```

on the way. As a result, they have to take a long detour once they reached the hotspot city on the way because entry into the city was not allowed. In such situations and many others, such as avoiding potential gatherings that may be a potential hotspot, our approach can recommend a safe route, thereby avoiding hotspot areas in advance so that users can plan their travel accordingly.

Suppose a user wants to travel from source location s to target location t as shown in Fig. 9. If the user follows Google map, it might show the *green* route since it is the shortest path. However, it may take the user to the hotspot city since it does not have *potential hotspots* information in *real-time*. Following our previous approach, we have the information of *potential hotspots* based on the previous section's discussion. Therefore, we remove the nodes representing the hotspot (red node in the Fig. 9) as well as edges incident to it and run the Dijkstra's algorithm [4, 16] on the remaining subgraph. The algorithm produces the shortest path in the subgraph (marked in blue color). This path does not fall in red zones/hotspot areas and thus safe for travel. We verify our approach on a toy dataset and show its efficacy in the experiment section. The safe route recommendation algorithm is described in Algorithm 3.

In the above algorithm, G is the input graph, S and D are the sources and destination nodes (cities), and $HOTSPOT$ is the list of nodes marked as a hotspot thus not safe for travel.

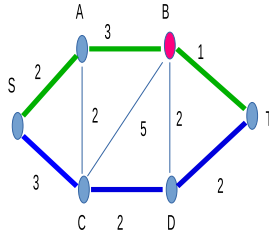


Fig. 9. Safe route recommendation. The shortest route might not be safe. Google Maps recommend the Green route since it is the shortest. Our approach recommends the Blue route though it may not be the shortest, it avoids the potential hotspot areas.

6 Complexity Analysis

In this section, time and space complexity analysis of various operations is presented.

Time Complexity of the Mapping a GPS coordinate to a tree leaf is $O(h)$ where h is the height of the tree. In our case, $h = 4$, so mapping GPS coordinates to tree leaf is a constant time operation $O(1)$. Further, append operation of the users into the list attached to the leaf node will have a time complexity of $O(1)$.

After mapping, we need to find a new cluster of users missed in the mapping process as discussed in Sect. 3.2. Such users can be put into new clusters by sliding a window of size 2, moving over adjacent pairs of clusters, and computing Euclidean/Haversine distance for each pair of users from each cluster. Time taken by this operation is $O(ck\ell)$, where k and ℓ are the sizes of two clusters and $c(= 90 * 60 * 60 * 6)$ is a constant. As discussed previously, the size of the clusters can not grow unbounded in the worst case due to population dynamics. Usually, k and ℓ are in the range of 1 to 100 since within a circle of 5 m radius; we assume no more than 100 users can stay. Further, the operation of finding missed users can be computed efficiently using parallelization (presented in the empirical section).

Time Complexity of the Intersection of leaf node (clusters) in the static case is more straightforward. We need to perform the intersection of $O(c)$ corresponding leaf nodes in two trees obtained at time intervals of $t/2$ mins. Here c is as defined previously. Intersection operation in the dynamic case is involved and incurs extra overhead due to one leaf being intersected with 168 ($=28 * 6$) leaf nodes in the other tree for the pedestrian case. Consequently, time cost goes up by a factor of 168 in the dynamic case compared to the static case. Further, note that the intersection operation in both the cases (static and dynamic) can be efficiently parallelized since one leaf node can be intersected with other nodes in parallel. Such parallelization tricks are used in the empirical section to reduce the runtime cost of the intersection operation.

Space Complexity: Total number of internal nodes for mapping latitudes is $90 + 90 * 60 + 90 * 60 * 60 + 90 * 60 * 60 * m = 329490 + 32400m$. Where m is the number of partitions of a second. Take $m = 6$ for $d = 5$, we have total internal

Table 1. Run-time comparison in the static and dynamic cases. All time in seconds. Values in blue and black denote run-time in static and dynamic case respectively. Cells shaded in yellow color show serial run-time whereas cells without shading indicate parallel run-time. Best viewed in color.

Time to Map GPS to Tree				Intersection Tree Time	
Before(at T=t)		After(at T=t+ Δ t)			
Latitude	Longitude	Latitude	Longitude	Latitude	Longitude
20.18	20.43	20.67	20.96	49.98	46.66
21.99	22.09	23.45	23.67	150.89	161.13
30.43	32.98	35.87	34.88	113.67	122.32
31.9	33.88	37.98	34.98	587.32	632.49
46.76	48.87	43.42	49.44	496.75	513.23
50.32	47.64	48.98	51.32	1208.98	1267.22
55.43	51.4	59.32	24.55	576.23	612.43
65.12	68.32	66.23	71.42	2198.32	2223.64
68.98	72.7	74.89	79.84	786.67	796.87
74.99	76.65	79.43	80.1	3983.78	4073.27
19.89	19.12	20.22	22.43	8.32	7.39
22.43	18.9	21.32	22.94	22.69	29.32
30.21	33.41	33.78	36.98	32.43	40.98
32.76	36.74	38.01	39.63	78.19	102.23
48.64	46.75	50.43	49.31	60.98	73.92
50.76	53.21	48.98	54.67	365.42	398.71
59.45	62.32	58.43	56.49	86.15	92.43
61.23	65.43	66.32	68.32	598.41	631.42
78.23	71.96	75.63	70.54	122.34	134.43
85.78	79.12	83.45	81.12	945.32	990.87

nodes as 2273490. Assume each integer takes 8 bytes. The space complexity of an empty tree is approx 18.18 MB. Further, we are just storing user ids, which for 10M users takes 80 MB. In total, one tree after mapping 10M GPS coordinates will take 98.18 MB of space. If we also store lat/long for 10M users, it will further take 160 MB. Thus in the worst case, one tree will occupy approximately 260 MB. (During execution, we store two trees every 2.5 mins interval. To generate contact vectors of 10M users takes around 17 mins as given in the experimental section; we maintain a total of 16 trees before we can flush the trees and reuse the space. Hence our approach takes roughly 4 GB of space for generating contact vectors compared to 55 GB in [10].)

In our implementation, we take latitude $\simeq (7^\circ, 37^\circ)$ and longitude $\simeq (68^\circ, 97^\circ)$ for India for which empty tree takes ≈ 39 MB when representing

internal nodes as four-level dictionary keys and leaves as the empty list in the python programming language. We store app IDs or some other useful global identifier such as IMEI number into leaf nodes to track each user. Thus space complexity for storing user id into leaf will be $O(n)$ where n is the number of user ids. Therefore, the dominating part will be used for storing user ids, and it scales linearly with the users.

One crucial point to consider is that the leaf node is implemented as a list or array which has the cost of $O(1)$ for append operation. However, a growing list beyond specified size incurs extra overhead. However, we must mention that *any* list does not hold more than a constant number of users since each leaf maps GPS range of 5 m and in a circle of 5 m radius, we assume that not more than 100 users can stand.

7 Empirical Demonstration

In this section, we show the working of the proposed methodology for contact tracing, hotspot detection, and safe route recommendation on simulated data.

7.1 Contact Tracing Experiment

To show the scalability of the contact tracing approach, uniformly at random GPS coordinates in the range latitude $\simeq (7^\circ, 37^\circ)$ and longitude $\simeq (68^\circ, 97^\circ)$ are generated for 2/4/6/8/10M users. We will map these GPS coordinates to the M-way tree like data structure as discussed and compute the contact vectors. Our goal is to estimate the time it takes to generate the contact lists. The serial and parallel versions for static and dynamic cases are implemented in C++ with openMP [3] (for multi-threading). The experiments are conducted on a personal computer with Intel(R) Core(TM) *i5-7500U* CPU and 8.0 GB memory (RAM) running Ubuntu operating system. The results are shown in the table 1. We have several observation from the results in the table. Firstly, mapping GPS to tree takes at most 80 secs for 10M users using serial implementation, whereas at most 86 secs in the parallel implementation. A little increase in time could be due to thread scheduling overhead in the latter case. In other words, mapping GPS to trees via parallel implementation is not *recommended*, and serial implementation suffices to practical cases. Secondly, compared to mapping GPS coordinates to trees, performing intersection of leaves is time-consuming. Thirdly, intersection time using parallel implementation achieves dramatic speedup. For example, parallel implementation in the dynamic case for 10M users gains $4 \times$ speedup. Finally, generating the final contact vectors for 10M users takes around 16.5 mins (990.87+85.78 secs). In other words, if we want to track contact events every 5 mins, we need to maintain four sets of lat/long trees before we can flush the data in the trees and reuse them again, thereby reducing the memory footprint.

Baselines One of the baselines used to compare the runtime of the contact tracing approach is by performing pairwise comparison (the naive solution).

Table 2. The runtime comparison of our approach with the baseline. The experiment is conducted over 2M users. The number of threads uses in the parallel implementation is set to 10.

	Implementation	CPU time(s)
Naive approach	Serial	1135714.0
	Parallel	127124.59
Ojagh et al. [13]	Serial	428
Our approach	Serial	185.80
	Parallel	52.26

Another baseline we use it from [13] in which the author perform contact tracing in an indoor setting using a graph-based approach. The result from the serial and parallel implementation of the baselines are shown in Table 2. It is obvious that our approach beats the naive solution by a significant order of magnitude as well as $2.3\times$ faster than [13] and favors practical implementation.

7.2 Hotspot Detection Experiment

To verify the approach discussed in Sect. 4 for hotspot detection, we first test it on 20 users and mark them on google map as shown in Fig. 10. All the users, the COVID +ve patients, and the suspected users (users who came in contact with COVID +ve users but still having no symptoms) are plotted on the city map. The figure also shows that at time t_0 , users ids 1, 3, 9, and 14 are COVID +ve. The Fig. 10(b) shows that user 6 who came in contact with user 9 (a COVID +ve), has been suspected and marked with the yellow color.

Similarly, user 16 becomes susceptible at time t_2 and so on. Indeed, the figure shows how the infected/suspected users *may* be moving in the city to provide real-time information of danger zones, benefiting normal users to plan accordingly. If the number of infected/suspected users in the map relative to the reference point exceeds a certain threshold, then the area around the reference point can be marked as a potential hotspot area. Alternatively, any area comprising a suspected COVID-19 positive patient can be declared as a *potential hotspot*.

One concern related to privacy is that we are marking users at the individual level. To avoid such a thing, we can draw a circle of larger radius such as 10 or 20m to conceal the COVID +ve users' identity. Alternatively, we can group nearby users and form large clusters to make the identity indistinguishable at the individual level. Additionally, the user's data remains with a single entity (government entity likely) hence the privacy is not compromised due to data sharing (assuming the government is truthful).

To show the scalability of the proposed approach, the plot of runtime with respect to the number of users considered in a hotspot region is shown in Fig. 12. The graph clearly shows that the runtime is following a *near-linear* trend. Therefore, our approach to detect hotspots in real-time may be useful.

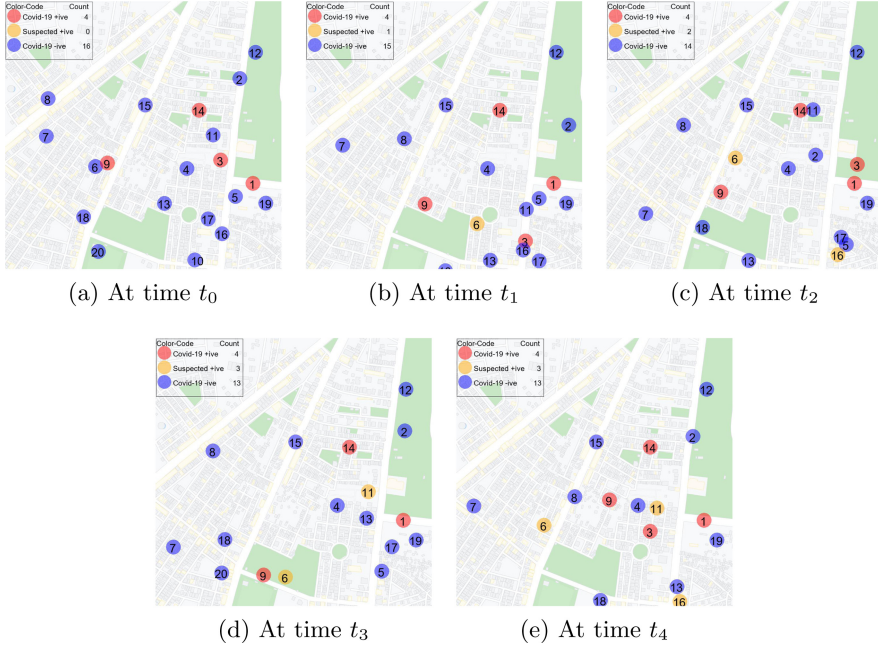
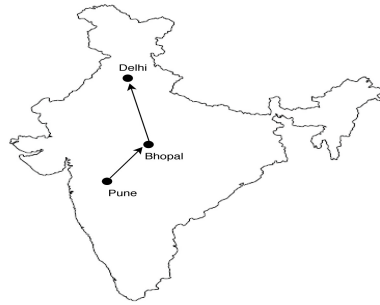


Fig. 10. Hotspot information at 5 consecutive time periods. Best viewed with zooming the image and in color.



(a) Route recommended by Google map



(b) Route recommended by our approach

Fig. 11. (a) Route recommended by Google map (b) Route recommended by our approach. We can see that Google map recommends a route from Pune to New Delhi, which passes through Indore city, a hotspot area and not allowed for travel. On the other hand, our approach recommends a route that avoids the hotspot city.

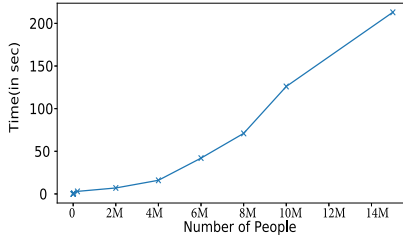


Fig. 12. The runtime of the hotspot detection algorithm wrt different number of users

7.3 Safe Route Recommendation Experiment

The evaluation of the safe route recommendation approach is done in the following way: 1) The proposed method is compared with an existing route Recommendation system, that is, the Google Map, in terms of the total cost of the path. The total distance from the source to destination constitutes the total cost and whether it is a Hotspot Zone free path or not. The efficacy of the proposed approach is analyzed on a real World dataset, which is extracted from Google Map. For the experiment, we have considered major cities in India (Indore, Bhopal, Chennai, Mumbai, Delhi, Bangalore, Lucknow, Hyderabad, Pune, Kolkata).

We have compared our proposed approach with the current Google map recommendation, considering Indore as a Hotspot City, Source as Pune, and Destination as New Delhi. Google maps recommend the path which passes through the hotspot city (Indore) as shown in Fig. 11(a) with a total cost of 1455 Km whereas, the path displayed by the proposed algorithm is an alternate route where the intermediate city is Bhopal (non-hotspot city) instead of Indore with a total cost of 1540 Km which is displayed in Fig. 11(b). Indeed, the path recommended by our approach is not optimal. Still, since our objective is to avoid a hotspot city, our approach is suitable for choosing a safe route rather than an optimal path.

8 Conclusion and Future Work

The proposed approach for contact tracing seems plausible based on the initial experiments on simulated data. Safe route recommendation and potential hotspot information further add new features to our method which is absent in the apps available in the market. Our approach scales well for large data and hence can be used for deploying over big cities. We are planning to release the app in the future after verification on the real-user study.

References

1. Ascione, G.: On the construction of some deterministic and stochastic non-local sir models. *Mathematics* **8**(12), 2103 (2020)
2. Christaki, E.: New technologies in predicting, preventing and controlling emerging infectious diseases. *Virulence* **6**(6), 558–565 (2015)
3. Dagum, L., Menon, R.: Openmp: an industry standard api for shared-memory programming. *IEEE Comput. Sci. Eng.* **5**(1), 46–55 (1998)
4. Dijkstra, E.W., et al.: A note on two problems in connexion with graphs. *Numerische mathematik* **1**(1), 269–271 (1959)
5. Farrahi, K., Emonet, R., Cebrian, M.: Epidemic contact tracing via communication traces. *PloS one* **9**(5), e95133 (2014)
6. Ferretti, L., et al.: Quantifying sars-cov-2 transmission suggests epidemic control with digital contact tracing. *Science* **368**(6491) (2020)
7. Hellewell, J., et al.: Feasibility of controlling covid-19 outbreaks by isolation of cases and contacts. *Lancet Global Health* **8**(4), e488–e496(2020)
8. Huerta, R., Tsimring, L.S.: Contact tracing and epidemics control in social networks. *Physic. Rev. E* **66**(5), 056115 (2002)
9. Lessler, J., Azman, A.S., McKay, H.S., Moore, S.M.: What is a hotspot anyway?, *Am. J. Tropical Med. Hygiene* **96**(6), 1270–1273 (2017)
10. Mahapatra, G., Pradhan, P., Chattaraj, R., Banerjee, S.: Dynamic graph streaming algorithm for digital contact tracing (2020)
11. Martin, T., Karopoulos, G., Ramos, J.L.H., Kambourakis, G., Fovino, I.N.: Demystifying COVID-19 digital contact tracing: a survey on frameworks and mobile apps. CoRR abs/2007.11687 (2020). <https://arxiv.org/abs/2007.11687>
12. Müller, J., Kretzschmar, M., Dietz, K.: Contact tracing in stochastic and deterministic epidemic models. *Math. Biosci.* **164**(1), 39–64 (2000)
13. Ojagh, S., Saeedi, S., Liang, S.H.: A person-to-person and person-to-place covid-19 contact tracing system based on ogc indoorgml. *ISPRS Int. J. Geo-Inf.* **10**(1), 2 (2021)
14. Sanampudi, A.: Indore navigation mobile application using indore positioning system (ips). *Int. J. Basic Sci. Appl. Comput. (IJBSAC)* **2** (2020)
15. Team, P.P.: Pan-european privacy-preserving proximity tracing (2020)
16. Vaughan, A.: The problems with contact-tracing apps. *New Sci.* **246**(3279), 9 (2020)
17. Wiki, O.: Precision of coordinates – openstreetmap wiki, (2019). wiki.openstreetmap.org/w/index.php/title/Precision/of/coordinates/oldid/1872063. Accessed 20 May 2020
18. Yoneki, E.: Fluphone study: virtual disease spread using haggie. In: *Proceedings of the 6th ACM Workshop on Challenged Networks*, pp. 65–66 (2011)