

EpiMap: Towards Quantifying Contact Networks and Modelling the Spread of Infections in Developing Countries

Eiko Yoneki
University of Cambridge
Cambridge, United Kingdom
eiko.yoneki@cl.cam.ac.uk

Jon Crowcroft
University of Cambridge
Cambridge, United Kingdom
jon.crowcroft@cl.cam.ac.uk

ABSTRACT

We describe the EpiMap project, in which mobile phones and sensors record the proximity of other devices, to gather information on human interactions within the rural communities of developing countries. Collected information will be used to develop improved mathematical models of the spread of infectious diseases, such as measles, tuberculosis and pneumococcal diseases. Modelling will be complemented by the use of surveys to aid in the understanding of living conditions in these villages. EpiMap is an extension of the FluPhone project, which we carried out in 2010. FluPhone collected data on human interaction, by using mobile phones to record information such as locality and user symptoms. Delay tolerant opportunistic networks such as the Huggle framework [5] were used as a basis for communication. We introduce the EpiMap vision for a system of opportunistic networks combined with satellite communication, designed to face the challenges posed by weak power and communications infrastructure in the rural regions of developing countries in Asia, Africa and South America. We aim to use a delay-tolerant small satellite for data transfer between developing countries and Europe and North America. Data collected through EpiMap can be used to help design more efficient vaccination strategies and equitable control programmes.

Categories and Subject Descriptors

I.6 [Computing Methodology]: Simulation and Modeling; C.24 [Computer Systems]: Computer Communication Networks

General Terms

Measurement, Experimentation, Algorithms

Keywords

Contact Networks, Epidemiology, Small Satellite

1. INTRODUCTION

Many of Africa's significant diseases such as measles, tuberculosis, meningococcal and pneumococcal disease, respiratory syncy-

tial virus, and influenza are spread directly via person to person contact. These diseases are vaccine preventable and there has been a significant investment in improving vaccine coverage and introduction of new vaccines in some of the poorest countries. Funding for vaccination programmes is limited and many countries face difficult decisions to refine the effective vaccination strategies within the limited budget.

Modelling the spread of infectious disease mathematically has been a useful tool for helping to design efficient immunisation programmes. Despite this, there is a lack of such studies. In Africa, few transmission models of vaccine preventable diseases have been developed. Thus, it is desirable to develop models for specific disease by making available social contact data, thereby encouraging others to develop their own models. For example, having a model prepared in advance based on up-to-date data on contact patterns and measles transmission data will greatly aid the efficient design of measles control and eradication efforts. The lack of any such model with which to evaluate the polio endgame has certainly hampered and delayed any progress of controlling the spread.

To achieve this goal, both developing an advanced mathematical modelling and more importantly reliable and informative social contact data are essential. Mathematical models are only as reliable as the assumptions and parameters upon which they are built. Social mixing patterns are central determinants of transmission for infections which demand close contact between individuals. In reality, very little is known about contemporary social contact patterns, especially in the third-world countries. Quantifying the social network structures through which close-contact diseases are spread will greatly aid our understanding of their epidemiology and help build better models. Such models can be used to predict the course of an outbreak, help design efficient control programmes, and help us to understand and control the emergence and spread of novel pathogens.

Contact diaries have traditionally been used to record such information. The contact diary is developed through self-reporting or an interview-led processes. Participants in a number of studies were asked to keep a detailed record of the number and characteristics of their encounters over a period of time. The diary based approach has its limitations. It becomes a burden when repeated over many days or weeks; consequently it is usually applied as a single day snapshot of individuals' contact patterns.

Recently, the use of sensor devices to collect social contact information has emerged as an alternative. All use variations of wireless proximity sensors. These have primarily made use of either Bluetooth enabled mobile phones, or RFID/Bluetooth sensors. In either instance low power radio transmissions are sent and received on a regular basis, detecting similar devices within range (typically 5-10 metres for Bluetooth, whereas RFID sensors are more tun-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ACWR'11 December 18-21 2011, Amritapuri, Kollam, Kerala, India.
Copyright 2011 ACM 978-1-4503-1011-6/11/1 ...\$10.00.

able). These methods have the advantage of being minimally invasive, just requiring participants to carry or wear the device over a period of time. Records of encounters can typically be stored on the device, and uploaded periodically to a server or base-station. Currently these methods have primarily been used in small-scale studies in enclosed settings (such as schools, hospitals and conferences). However, they offer the possibility of capturing the fine structure and dynamics of social networks at different spatial and temporal scales.

We have previously demonstrated the FluPhone Project [10][2], which aims to bring together epidemiologists, sociologists, and computer scientists, with the goal of developing novel and innovative methods with which to measure and understand social encounters based in Cambridge, UK. Such information helps scientists and medical researchers to understand how close-contact infections, such as swine flu, spread between different people. The FluPhone project was mainly targeted for tackling flu-like symptoms, which was a threat in our society a couple of years ago. Human proximity information is collected using phones with Bluetooth communication from the general population to build time dependent contact networks. The project also included a ‘virtual disease’ experiment, where a specific model of disease is spread through the proximity based communication upon encountering of two devices. The spread of different stages of the disease was then mapped across the locality of the study, and fed back to the user.

The FluPhone is built over the Haggle framework [5], in which we introduced Pocket Switched Networks (PSNs), a type of Delay Tolerant Network (DTN), exploring proximity based communication. PSNs provide communication in highly stressed settings with intermittent connectivity, variable delays and high error rates in decentralised and distributed environments over a multitude of devices that are dynamically networked. A partitioned network can deal with disconnected operation using a store-and-forward approach to communication. In PSNs people carry devices in their pockets, which communicate directly with other devices within their range or with infrastructure. Because device mobility is reflected by the user’s movement, we have worked on understanding the social structure among the people who carry the devices. In many ways, the concept of PSNs is analogous to how infectious diseases spread. One key aspect of this is working out the numbers of social encounters and links in a chain of contacts between different people (similar to the idea of how many steps we are away from a particular person). One way to measure this is to record how often different people (who may not know each other) come into close proximity with each other, as part of their everyday lives. PSNs share many issues with epidemiological studies.

In this paper, we introduce a vision of the EpiMap project, where we extend the FluPhone project to be able to deploy large-scale social contact data collection in developing nations. We plan to include not only phones but other sensing devices such as RFID tags and to exploit various communication methods. For example, the use of satellite communication will have a great potential in overcoming the limitations of sparse power and communications networks in the targeted regions. In parallel to data collection by EpiMap, we propose to conduct surveys for validating models derived by EpiMap on the spread of respiratory infections. We also briefly demonstrate what we can be modelled using the collected social contact data. It is our intention to develop a range of mathematical models based on contact patterns, which can be used to help guide vaccine policy development over the coming decade in Africa and other developing countries.

2. FLUPHONE PROJECT

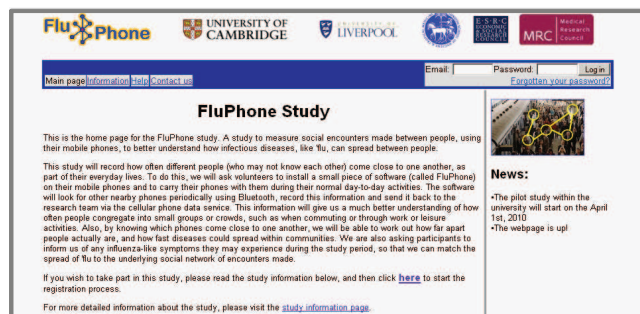


Figure 1: FluPhone Study Web Interface

The FluPhone project studies how often different people (who may not know each other) come close to one another, as part of their everyday lives. To do this, we asked volunteers to install a small piece of software on their mobile phones and to carry their phones with them during their normal day-to-day activities. The software will look for other nearby phones periodically using Bluetooth communication, record this information and send it back to the research team via a cellular phone data service or other means. This information gives us a much better understanding of how often people congregate into small groups or crowds, such as when commuting or through work or leisure activities. Also, by knowing which phones come close to one another, we will be able to work out how far apart people actually are, and how fast diseases could spread within communities. We also asked participants to inform us of any influenza-like symptoms they may experience during the study period, so that we can match the spread of flu to the underlying social network of encounters made. Participants were able to log-on to the study website and see an estimate of how many people they have encountered. Further details can be found in FluPhone study web depicted in Figure 1 (see also [10][11]).

The FluPhone application can be downloaded from the web following a registration and authentication process. Data collection can operate over three different methods such as periodic uploading via the web, real-time collection via 3G, and post-study collection from the devices.

The FluPhone Study is being carried out in Cambridge, United Kingdom, and was advertised via the channels of the University of Cambridge and various online social media such as Facebook [7] and Twitter [21]. Targeted participants include university members, their families, colleagues, friends, and people who work or live in Cambridge. Participants must be over 12 years old (under 16s require parental/carer consent), have the use of a compatible mobile phone, and permission from the owner and bill-payer of the phone to participate. The registration process requires the consent of each participant.

In our previous work, we have developed a range of ways for detecting and recording spatial proximity [5]. These include small custom built battery-powered sensor devices (i.e. Intel iMote) and mobile phones. In each case, the software has been developed to record contacts with other devices. Each device is uniquely identifiable, so a network of contacts, which includes information about which devices interact, can be built. Furthermore, the duration of interactions are logged (both the duration of a single interaction and the cumulative duration of all interactions over the study period) to enable weights to be assigned to links in this contact network. The technology has proved robust with reliable data collection at initial level. However, for tackling reliable network modelling for epidemiology requires further massive and precise experimental data from the general population. The FluPhone project aims at such



Figure 2: FluPhone: Encountering Statistics, Symptom Type and Time Entry Screens, and Brief Diagnostics

massive scale data collection and, to our knowledge, is the first attempt of using mobile phones for this purpose.

2.1 FluPhone Software

FluPhone provides software that runs on the users' mobile phones which the users carry with them during their normal day-to-day activities. FluPhone adopts a simple client-server design consisting of a mobile phone application in the phone and a receiver as a PHP (Hypertext Preprocessor) script on a web server. The mobile phone application is written in Java (J2ME), and collects Bluetooth device proximity data, Global Positioning System (GPS) coordination data, and data on self-reported flu symptoms through a user interface. The collected data is sent via GPRS/3G to the server. The user can also upload the data over a web interface. An overview of the system architecture is shown in Figure 3. Figure 2 shows a screen capture of the application. Communication between the application and the server is based on the secure Internet protocols and standard public-key infrastructure. Through the secure login to his/her account, the participant can look at their history of activities.

The collected information provides us with a much better understanding of how often people congregate into small groups or crowds, such as when commuting or through work or leisure activities. Also, by knowing which phones come into close proximity, we will be able to work out how far apart people actually are, and how fast diseases could spread within communities. Without complex analysis a simple history of encounters can be used to measure activity, environment, cyclic behaviour, and so forth. The experiment result shows an average person encountered over 1500 unique devices over a ten day period.

3. VIRTUAL DISEASE EXPERIMENT

There is also a function called 'virtual disease', which models epidemics on participants' phones, giving a real-time picture of the social network between participants from the perspective of infectious disease. The virtual disease application is implemented as an

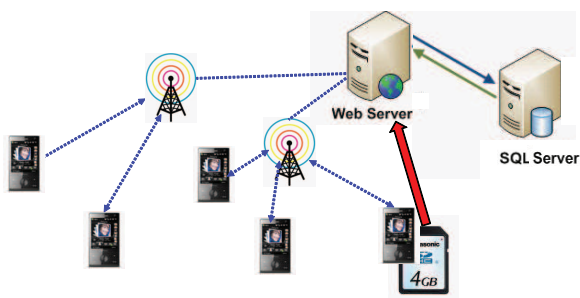


Figure 3: Overview of FluPhone System

Android application built to run on devices that support the haggel architecture [5]. The application broadcasts information about virtual diseases. It is currently infected to all devices in range, and receiving devices are infected by these virtual diseases based on a simple probability calculation. The application logs all incoming diseases and stores information regarding how they are processed. It also regularly scans for Bluetooth and GPS based location data. In Virtual Disease experiment, the spreading of diseases

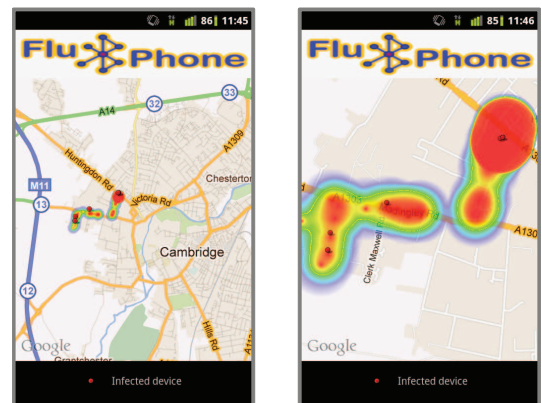


Figure 4: Virtual Disease showing Infection Map

is simulated with a simple SEIR model (S: Susceptible, E: Exposed, I: Infectious, and R: Recovered). In this model, each device is originally susceptible to a disease. Once it is infected by another device, it becomes exposed for a specified time. Whilst it is exposed, a device has the disease but cannot yet infect other devices. Once the exposed duration has run out, the device becomes infectious for a specific time. Whilst it is infectious, the device can infect other devices. Each disease has an associated infection probability which indicates the likelihood that another device will be infected. Once the infectious duration has run out, the device has recovered from the disease and cannot be reinfected. The application shows the state of the infection of each disease including who/when passed the virtual disease. The screenshot of the heatmap of virtual disease statistics is shown in Figure 4. The locations of the infectious nodes are depicted in the map shown in the left side of Figure 4, which can be zoomed in showing more details in the area in the right side of Figure 4. We deployed three different diseases as follows, where β =transmission probability, α =incubation time, and IP =infection period.

- SARS: fast ($\beta=0.8$; $\alpha=24$ hours; $IP=30$ hours)
- FLU: normal ($\beta=0.4$, $\alpha=48$; $IP=60$ hours)
- COLD: slow ($\beta=0.2$, $\alpha=72$ hours; $IP=120$ hours)

4. DISEASE SPREAD SIMULATION

In our previous work, we have worked on uncovering community structures, centrality node, weighted networks and so forth [15] [24] in the context of PSNs. In [23] and [12], we looked into inter contact time, meeting time, and epidemic spread patterns. In this section we briefly demonstrate multi-modal spread modes extracted from the contact networks. In this analysis, we defined spanning trees on various time points and aggregated them using ‘joint diagonalisation’ (JD), which is a technique to estimate an average eigenspace of a set of matrices. Using JD on matrices of spanning trees of a network is especially useful in the case of real-world contact networks in which a single underlying static graph does not exist. The average eigenspace may be used to construct a graph which represents the average spanning tree of the network or a representation of the most common propagation paths. We then examine the distribution of deviations from the average and find that this distribution in real-world contact networks is multi-modal; thus indicating several modes in the underlying network. These modes are identified and are found to correspond to particular times. See [9] for further detail of the analysis.

This technique was applied to a contact network trace with 36 students in Cambridge, and revealed five spread modes corresponding to different times of day. Figure 5 depicts one of the modes,

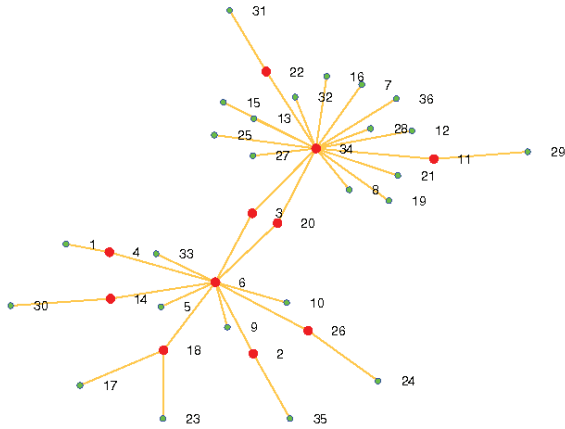


Figure 5: Average Graph in Mode 5

where a highly structured network corresponding to the day when the groups are well defined by *class year* (i.e. year 1 and year 2). There is an obvious bridge formed by nodes 3 and 20. Using the average graph as an indicator, this implies that a disease spread at this time from nodes 3 and 20 should have the fastest infection rate.

To test the infection rate, an SEIR model is constructed, setting the probability of infection to 0.5; infection time Poisson distributed with mean time of 800 minutes. A disease is spread through the contact network starting at time index 250. The simulation is repeated 30 times for each node and the results bootstrapped to give estimates of the mean number of people susceptible (those that have not received the disease) at time, t , $S(t)$. Figure 6 shows the results of these simulations and as can be seen the number of susceptible people falls most rapidly for infections started at nodes 3 and 20, as expected.

5. FLUPHONE TO EPIMAP

The use of mobile phones and sensors for quantitative measurement of societal mixing patterns to underpin mathematical models of the spread of close-contact diseases has distinct advantages over other methods of collecting contact data (such as diaries and

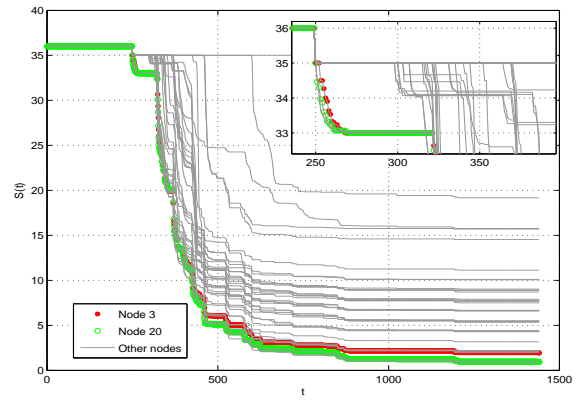


Figure 6: Mean number of nodes susceptible to disease after time t (Root infection starts at time 250. Inset focuses on the start of the infection)

interviews). Such devices can be programmed to gather proximity data automatically, allowing detailed longitudinal studies to be conducted with no possibilities of re-call bias, no barriers due to problems of literacy or understanding, and minimal disruption to the participants in the survey. They therefore offer an unparalleled opportunity to collect information on social contact patterns that would allow a step-change in our understanding of the patterns of disease spread. Despite the emergence of such technology, the use of mobile phones for this purpose has not been explored in this setting before.

The environments in developing countries vary and we will develop an application (named EpiMap) for various types of devices: mobile phones with or without GPS functionality, sensor boards, and RFID tags. Thus, we need to consider a hybrid model of data collection, allowing operation in the absence of reliable electricity supply, lack of Internet access, and so forth. We will use Bluetooth communication and RFID based communication for detecting devices within radio proximity, along with sensing (e.g. movement, light, humidity). Various sensors are embedded in phones and devices, and can be used to capture changes in surrounding context information in order to infer the behaviour and patterns of the device carriers. Accelerometers are especially useful for this purpose. The EpiMap phones and devices can be carried by the study participants and by health workers. The EpiMap application will also feature additional functions such as a display of nearby mobile phones, a history of encounters and can also support logging capability for events, activities of participants, and any useful observations of the environment. Battery life is an important issue, and the application will be carefully designed to minimise I/O and CPU usage. The study will be carried out initially in a rural African community, and we will therefore combine the use of spare batteries and power charging efficiency for the duration of study.

We will also augment the electronically-collected social contact data with several weeks of self-reported surveys. Data analysis will compare the different social contact and survey tools, and mathematical modelling will validate the mixing patterns recorded against known age-based zero-conversion rates for common respiratory infectious diseases in the population.

The information collected by EpiMap will be kept within the device until the end of the study or uplifted via automated data collection, which could be deployed using delay tolerant networks over Bluetooth communication or with WiFi base stations. The collected data will comprise a time series of encountering pairs of phones within physical proximity. A time-dependent graph of connectivity

among phones can be constructed, detailing the duration and frequency of interactions made between participants. This will form the basis of analysis with other tools. Apart from building a contact map, various types of contextual information can be recorded by EpiMap such as movement and light, which will be analysed together with the diary and interview survey tools, where those tools will record the number of unique individuals encountered by that participant during a day, as well as some contextual information (whether at home, work, health facility, congregation, school, market, socialising). The data collected by survey-based methods will include some information on individuals not participating in the study, so at the individual level, the phone-collected data will intersect with the diary and interviews information. Direct comparison between tools will consider participant-participant encounters only, and will focus on accuracy of recording, reciprocity, duration, and frequency of encounters thought to occur between a participant's contacts. This comparison will enable correction factors, by age and sex, to be estimated for future phone-only studies.

The dynamic network of connections between participants will be used to investigate the topology of the social network, including: 1) duration-weighted pairs: time spent in close-proximity is a powerful determinant of infection risk and these can be considered as a weighted link between individuals with location and context associations; 2) number of encounters per person: to determine whether some individuals are responsible for a disproportionate number of contacts 3) social distances: betweenness and centrality measures describe how far apart individuals are in a network, and strongly impact disease dynamics; 4) community structure: identifying individuals that form bridging links between otherwise distinct groups offers efficient targeted interventions. We expect all of these measures and the network to change rapidly over time, and we will consider the implications for disease transmission in the modelling stage. Added value comes from the creation of a well-defined dynamic network of real human interactions, pertinent to respiratory infectious diseases, which will be a valuable contribution to the development of a new generation of analytical methods and measures that cope explicitly with dynamic networks. While analysis of the survey data can inform greatly on the potential for disease transmission, using mathematical models of disease transmission enables all the features to be incorporated at once to consider the implications of contact patterns and network structure on transmission, prevalence and possible intervention regimes. Final validation of the survey tools will be provided by using the models to simulate the spread of respiratory infections, and compare the age-based prevalence generated by the models to age-based incidence and longitudinal household prevalence of pneumococcal disease carriage information for the region. In the initial step, we aim to deploy EpiMap in the scale of 100-300 participants. We hope to obtain donations of devices from mobile phone and sensor vendors to increase the scale of our study.

After the initial prototype stage, we will widen the study to a range of settings, including urban and rural sites across different parts of the developing world. The application could be extended to web-based data collection. We will also expand the temporal window over which data are collected, to capture seasonal differences in contact patterns and assess how these may affect the spread of disease. We will develop mathematical models based on these new data, which will help us gain valuable insight into the spread and control of diseases. Examples of diseases to be modelled are tuberculosis, pneumococcal disease, meningococcal disease, measles, and disease associated with *Hemophilus influenzae*.

6. DATA COLLECTION AND COMMUNICATION

Building an effective and reliable human proximity detection system raises various issues. Particularly, optimal exploitation of technologies available across the hardware and software is necessary. Current detection mechanisms in the FluPhone using WiFi access points or Bluetooth expect high failure, communication protocol limitation and complex statistics. Without in-depth understanding of the data collection mechanism, modelling networks will not be reliable. For example, the symmetry of edge detection is extremely low according to our experiments using Bluetooth. Missing edges from device detection leads to inaccurate clustering coefficient calculation. This noise hampers our ability to infer deep knowledge from this data. We need to understand at least the scale of missing edges. However, such important information is entirely missing in current research efforts.

6.1 Proximity Detection

Bluetooth is a low-power open standard for Personal Area Networks (PANs) and has gained its popularity due to its emphasis on short-range, low-power and easy integration into devices. The platform used in the experiment in the Huggle project [5] is the Intel Mote ISN100-BA (known as the 'iMote'). The iMote runs TinyOS and is equipped with an ARM7TDMI processor operating at 12MHz, with 64kB of SRAM, 512kB of flash storage, and a multi-colored LED, and a Bluetooth 1.1 radio. The specifications lists the radio range to be 30 meters.

It is a complex task to collect accurate connectivity traces using Bluetooth communication, as the device discovery protocol may limit detection of all the devices nearby. Bluetooth uses a special physical channel for devices to discover each other. A device becomes discoverable by entering the inquiry substate where it can respond to inquiry requests. The inquiry scan substate is used to discover other devices. The discovering device iterates (hops) through all possible inquiry scan physical channel frequencies in a pseudo-random fashion. For each frequency, it sends an inquiry request and listens for responses. Therefore, a Bluetooth device cannot scan for other devices and be discoverable at the same time. Bluetooth inquiry can only happen in 1.28 second intervals. An interval of $4 \times 1.28 = 5.12$ seconds gives a more than 90% chance of finding a device. However, there is no data available when there are many devices and many human bodies around. The Bluetooth standard recommends being in the inquiry scan substate for 10.24 seconds in order to collect all responses in an error-free environment. The power consumption of Bluetooth also limits the scanning interval, if devices have limited recharging capability. The iMote connectivity traces in Huggle use a scanning interval of approximately 2 minutes, while the Reality Mining project uses 5 minutes.

Bluetooth for proximity detection is widely available and a lot of people carry a Bluetooth enabled mobile phone with them. Thus, it is possible to detect a certain amount of peoples' phones without handing a special device to each of them, which makes Bluetooth appealing for experiments involving a large quantity of people. The range of Bluetooth varies between 10m and 100m, depending on the device class. In mobile phones, the range is usually 10m. We have observed the devices can be detected in 20m range if there is no obstacles, while if there is any obstacles such as a thick wall it limits to 5m range.

We plan to extend to include audio recording when two devices are in proximity range so that the type of interaction can be inferred. Note that the audio is used for determining the patterns and tone of interaction and it does not examine the content of audio. as a prior work, Wyatt et al [22] have shown effective analysis of privacy

sensitive Audio.

It is increasingly popular to use radio frequency identification (RFID) for identification and security application. RFID transmits the identity (in the form of a unique serial number) of an object or person wirelessly, using radio waves. It's grouped under the broad category of automatic identification technologies. RFID is in use all around us. If you have ever chipped your pet with an ID tag, used EZPass through a toll booth, or paid for gas using Speed-Pass, you've used RFID. Unlike ubiquitous UPC bar-code technology, RFID technology does not require contact or line of sight for communication. RFID data can be read through the human body, clothing and non-metallic materials. RFID requires three components are depicted; an antenna, a transceiver (with decoder), and a transponder (RFID tag) electronically programmed with unique information. The antenna emits radio signals to activate the tag and to read and write data to it. then the reader emits radio waves in ranges of anywhere from one inch to 100 feet or more, depending upon its power output and the radio frequency used. When an RFID tag passes through the electromagnetic zone, it detects the reader's activation signal. The reader decodes the data encoded in the tag's integrated circuit (silicon chip) and the data is passed to the host computer for processing.

RFID quickly gained attention because of its ability to track moving objects. As the technology is refined, it will be more pervasive and invasive. A typical RFID tag consists of a microchip attached to a radio antenna mounted on a substrate. Normally the chip can store as much as 2 kilobytes of data.

The SocioPatterns project [19][18] use active RFID devices, embedded in unobtrusive wearable badges. Detailed information on how this technology is used to monitor social interactions and to identify contact patterns. Individuals are asked to wear the devices on their chests, so that badges can exchange radio packets only when the individuals wearing them face each other at close range (about 1 to 1.5 m). This range was chosen as a proxy of a close-range encounter during which a communicable disease infection can be transmitted, for example, either by cough or sneeze, or directly by hands contact. The infrastructure parameters are tuned so that the proximity of two individuals wearing the RFID badges can be assessed with a probability in excess of 99% over an interval of 20 seconds. The problem of RFID tag is that the reader is required, which is normally expensive and the tag normally does not have enough storage to keep the data for long time. Thus, frequent transferring the logged data to the reader device is necessary.

None of technique for the proximity detection is perfect currently and we will plan a hybrid system to tailor to the environment in developing countries.

Existing trace data typically lacks geographical information. We are experimenting GPS equipped mobile phones, small computers, and embedded Linux boards to design tracking and localisation mechanisms in an efficient and inexpensive way. Software that detects proximate devices for the phone will be developed based on our previous work, extending to GPS tracking and capturing image/sound as necessary to record contexts in the environments. The proximity networks represent pair relationships, proximity based modularity, and social interactions, while online based interactions such as email, instant messaging, and social network services (e.g. Facebook, LinkedIn, Orkut) represent another type of social interactions. We collect data from online social networks to be used for network analysis, including correlation between two types of social networks.

6.2 Satellite Communication

What about moving data to the location where data analysis can be

performed? In general, experimental data collection may need to be repeated many times with different configurations and it would be desirable to move data from the developing country to where we can work on the data. There may not be any infrastructure for the Internet or electricity in rural area. The condition is different from remote sensing. In remote sensing, the small sensor board can be left for a long time and later on those sensors would be collected.

In the TIER (Technology and Infrastructure in Emerging Regions) project [20], the concept of the Village Base Station (VBTS) is proposed, where the users in the village receive various services ranging from distributed caches and uploading user-generated content as data services [13]. VBTS provides an outdoor PC with a software radio that implements a low-power low-capacity GSM base station including power supply via solar/wind.

There may not be GSM access capability in rural area, and we attempt to integrate satellite communication using Delay Tolerant Networks (DTNs) [8]. DTN research started with Vint Cerf and the Interplanetary Internet initiative [4], which proposed a new architecture that could work over both terrestrial and interplanetary links. This architecture could enable applications such as the remote operation of scientific experiments on other planets, controlled using TCP/IP from Earth.

In contrast, in the Haggie Project [5] we exploited a Pocket-Switched Networks (PSNs) variant of DTN, where people carry devices in their pockets, which communicate directly with other devices within their range or with infrastructure. As people move around, they can exchange messages with nearby devices, carrying a message until it is close to another device. We will build a hybrid DTN using the satellite based communication and PSNs. All the collected data will be gathered up to a node that can communicate with the CubeSat by PSNs and the CubeSat communication node exchange the data when the satellite communication is available.

Aspirations of use of microsats and smallsats for remote sensing applications have been increasing in developing countries. Small startups in the United Kingdom and now South Africa have taken the lead in providing end-to-end, innovative solutions for countries with limited economies. CubeSat type platforms are considered to be well suited for building an early space technology capacity. Low Earth Orbits small satellites, positioned on inclined polar orbit will also be used.

Uploading collected data and receiving instruction do not require the same level of end-to-end steady connection as email or Voice over IP (VoIP) do. High-delay connectivity is acceptable as far as the data collection can be performed within required timeline. We vision that DTN could be one solution to bring communication in the developing countries. For data collection in EpiMap, the DTN concept is ideal. The important point here is that CubeSat needs to implement a DTN protocol. Figure 7 shows a system overview, where CubeSat moves around the earth and connects to the node, when the CubeSat is reachable by the node. The black-patterned and green(gray) nodes are in a rural area, where a black-patterned node is able to access a CubeSat, while a green(gray) node can only connect locally. The black-patterned nodes can upload/download data to the CubeSat and transfer them to the red(dark) nodes in Europe/USA. The red(dark) nodes will form an overlay to share data. Krupiarz et al [14] proposed using small SmallSats and DTN for Communication in Developing Countries. As the technology has been progressing, CubeSat will be increasingly becoming popular in the near future and CubeSat with DTN will be beneficial for the communication for the applications of 'disaster warning', 'volunteer aid worker consultation' and 'reporting', 'search and rescue', 'disaster control information', 'weather forecasts and education', and remote sensing data upload/download [3].

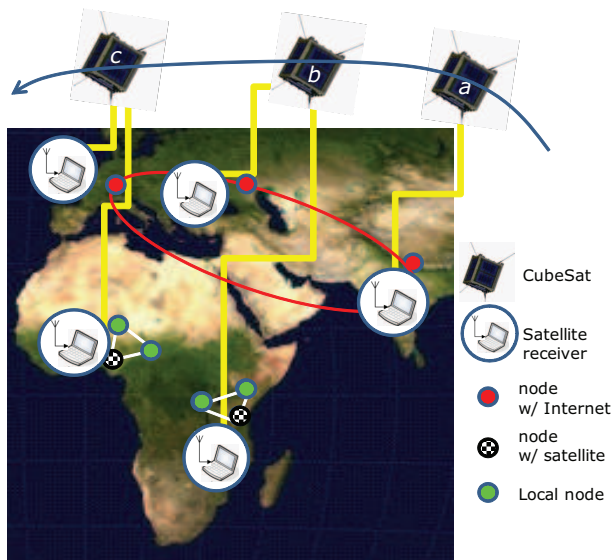


Figure 7: Use of CubeSat for Data Collection from Africa

SPICE project [17] is planned to build an overlay called ‘Space-Data Routers’ (SDR) [16], where space-data generated by a single or multiple missions can be shared among Space Agencies, Academic Institutes and Research Centres in a natural, flexible, secure and automated manners. A communication overlay modelled will be developed according to a thematic context of missions, Ground Segment topological distribution, Agency policies and Application restrictions and requirements. A DTN-enabled device that incorporates the Space Agency administrative instructions and policies for data dissemination and resource utilization and integrates the DTN protocol stack with application, network and link layer protocols is essential.

Potentially we would extend the SDR overlay by integrating CubeSat. In an overlay, some nodes may not have access to the CubeSat but locally it can transmit the data to the node that has access to the CubeSat. This enables the developing countries to upload the collected data also receive information/result of analysis from the data processing unit in Europe/USA. The challenges to realising such a communication platform include storing Tera byte of data, rule based routing, and load balancing.

7. DATA ANALYSIS STRATEGY

Human proximity networks display extremely dynamic topology on a spatial and temporal dimensions. We refer to such networks as ‘time-dependent networks’. Generalisation of the measurements of complex networks is a recent active research topic. However, modelling dynamic temporal and spatial series of sub-networks (e.g. trees or motif) in time-dependent networks in a discrete form is a future challenge.

Several researchers have worked on predictive models for epidemics such as an influenza pandemic [6]. Such models require precise information about mobility, interaction, and behavioral assumption of the population. However, interactions between individuals are assumed to follow existing contact models, that do not take into the changeable behaviour of human movement. We aim to build a model incorporating spatial and temporal information for improving the predictions.

Apart from confirming previously known results, such as that degree distributions with high variance of occurrence of high-degree individuals can be associated with an accelerated course of the epi-

demio [1], a little work has been done in this area. Many other network characteristics (e.g. population size, geographical location) can be uncovered. Clustering will be an important factor to drive the epidemic, and looking into causal contact patterns of the epidemics will give additional insight. The patterns of interactions between individuals are key to understanding how infectious diseases spread. Only considering one-dimensional pair relationships may not be sufficient and consideration of the strength and regularity of connections will be necessary.

In EpiMap, we set several goals of data analysis below:

- Estimate social contact parameters relevant to the spread of close-contact infectious diseases in a number of contexts in a number of different African countries including determination of age-related contact patterns, and comparison between settings (e.g. East versus West Africa, urban versus rural and wet versus dry season).
- Estimate patterns of contact with domestic animals and link with patterns of contact between humans
- Elucidate risk behaviours by linking epidemiological and social contact data
- Develop and parameterise mathematical models of a range of infections based on observed social contact data
- Provide a quantitative description of human-human and human-animal contact patterns within Africa for a wide range of researchers and scientific purposes.

Regression analyses will be used to assess factors underlying contact patterns, including age, gender, household size and composition, region, country, season etc, based on previously derived methods. The model can be used to help explain differences in contact patterns observed between different settings or in different seasons and generalise the results to other settings. Age and context-specific contact matrices will be generated, using previously developed methods. These information will be used to estimate the basic reproduction number with knowledge of the population distribution from the dominant eigenvalue of the next generation operator. Bootstrap confidence intervals for the reproduction number will be generated by re-sampling from the empirical contact data. The relative change in the reproduction number in different settings, and for diseases spread through different routes (e.g. airborne versus physical contact) will be calculated as will the effect of season on the reproduction number.

Data on contact with animals will be combined with data on contacts amongst humans to derive a series of animal-human-human contact matrices by age group and other key variables. These will be used to develop mathematical models designed to assess the emergence of novel pathogens into the human population and how they may spread during the critical first steps. Confidence intervals for relevant contact patterns will be derived by bootstrapping from empirical distributions. In addition, the data on travel and contact patterns within and between localities will be used to parameterise spatially-explicit models.

8. ETHICAL/PRIVACY/ANONYMITY ISSUES

We are well aware of ethical and privacy issues for the collected data, and the data will never be used to identify individuals. The collected data will be anonymised before analysis. Software developed for sensor devices and mobile phones may involve collaboration between ad hoc groups of members. When new encounters occur, there are complex issues in knowing what entities to trust. Based on predefined trust, recommendations, risk evaluation and

experience from past interactions, an entity may derive new trust metrics to use as the basis for authorisation policies for access control.

There will be various concerns about privacy, surveillance and freedom of action. For example, while providing location information can clearly be a one-way system where the location providing tools do not track the receivers, once a device receives information its location is potentially available to others. We will exploit various methodologies to protect the participants' privacy in EpiMap.

9. SUMMARY AND FUTURE WORKS

We describe our vision of the EpiMap project, where mobile phones and sensors record proximity to other devices. This project will gather information on human interactions in rural communities of developing countries in Asia, Africa and South America. The EpiMap project evolved from the FluPhone project in which we deployed data collection of human contact, flu-like symptoms and virtual disease spread using various phones. In EpiMap, Delay Tolerant Networking takes an important role to collect data and transmit them from developing countries to Europe. We envision to use CubeSat to move data in a delay tolerant manner and build an overlay network to share the data among the institutions. Within the developing countries, we will build a system for hybrid networking of opportunistic and infrastructure-based networks. Main challenges are the computer networks and power supply in rural villages of developing countries. The collected information will be used to develop improved mathematical models for the spread of infectious diseases, such as measles, TB and pneumococcal diseases. The modelling is complemented by surveys to understand the characteristics of living conditions in such rural villages. The outcome of EpiMap can be used to help design more efficient vaccination strategies and equitable control programmes.

In EpiMap, our study will be extended towards understanding human behaviour. Individuals may change their behaviour for several reasons: by being ill themselves, by caring for others who are ill, or by changing their normal habits in the belief it will minimise their risk of infection. A recent study suggests that public transport usage may decline in the event of an influenza pandemic and that people may stay at home rather than go to work and risk infection. If such precautionary behaviours were to be adopted by a large number of individuals, the economic implications would be profound.

Acknowledgment

This research is part-funded by the EPSRC DDEPI Project, EP/H003959. The data analysis and visualisation are carried out partially with Damien Fay and Fehmi Ben Abdesslem. The EpiMap Project will be a potential collaboration with John Edmunds and Ken Eames at London School of Hygiene and Tropical Medicine, and with Jon Reed at the University of Liverpool. We would like to thank Scott Burleigh at JPL for sharing the insight of DTN CubeSat, and Steven Smith and Karthik Nilakant for valuable comments.

10. REFERENCES

- [1] M. Barthlemy, A. Barrat, R. Pastor-Satorras, and A. Vespignani. Velocity and hierarchical spread of epidemic outbreaks in scale-free networks. *Physical Review Letters*, 92:178701, 2004.
- [2] BBC-News. <http://www.bbc.co.uk/news/uk-england-cambridgeshire-13281131>.
- [3] S. Burleigh and E. Birrane. Toward a Communications Satellite Network for Humanitarian Relief. In *International Conference on Wireless Technologies for Humanitarian Relief (ACWR)*, 2011.
- [4] S. Burleigh, A. Hooke, L. Torgerson, K. Fall, V. Cerf, B. Durst, K. Scott, and H. Weiss. Delay-tolerant networking: an approach to interplanetary internet. *IEEE Communications Magazine*, 41(6):128–136, 2003.
- [5] EU FP6 Huggle Project. <http://www.huggleproject.org>, 2010.
- [6] S. Eubank, H. Guclu, V. Kumar, M. Marathe, a. Srinivasan, z. Toroczkai, and N. Wang. Modelling disease outbreaks in realistic urban social networks. *Nature*, 429, 2004.
- [7] Facebook. <http://www.facebook.com>.
- [8] K. Fall. A delay-tolerant network architecture for challenged internets. In *Proc. SIGCOMM*, 2003.
- [9] D. Fay, J. Kunegis, and E. Yoneki. Uncovering multi-modal spread modes using joint diagonalisation in contact networks. Technical report, UCAM-CL-TR-806, University of Cambridge, 2011.
- [10] FluPhone-project. <http://www.cl.cam.ac.uk/research/srg/netos/-fluphone2/>.
- [11] FluPhone-Study. <https://www.fluphone.org>.
- [12] M. Freeman, N. Watkins, E. Yoneki, and J. Crowcroft. Rhythm and randomness in human contact. In *Proc. International Conference on Advances in Social Networks Analysis and Mining*, 2010.
- [13] K. Heimerl and E. Brewer. The Village Base Station. In *ACM Workshop on Networked Systems for Developing Regions (NSDR)*, 2010.
- [14] C. Krupiarz, C. Belleme, D. Gherardi, and E. Birrane. Using smallsats and dtn for communication in developing countries. In *Proc. International Astronautical Congress (IAC-08.B4.1.8)*, 2008.
- [15] F. Lynch and M. Zapp. Bubble Rap: Forwarding in Small World DTNs in Ever Decreasing Circles. Technical Report UR-CDL-TR-684, Care of David Lodge University of Rummidge, Cyber Science Lab Euphoric State University, January 2007.
- [16] SDR-Project. *SDR: Space-Data Routers*. <http://www.spacedatarouters.eu/>, 2011.
- [17] SPICE-Project. *SPICE:Space Internetworking*. <http://www.spice-center.org/description/>, 2011.
- [18] J. Stehle, N. Voirin, A. Barrat, C. Cattuto, L. Isella, J.-F. Pinton, M. Quaggiotto, W. V. den Broeck, C. Regis, B. Lina, and P. Vanhems. Simulation of an seir infectious disease model on the dynamic contact network of conference attendees. *BMC Medicine*, 9(87), 2011.
- [19] J. Stehle, N. Voirin, A. Barrat, C. Cattuto, and Others. High-resolution measurements of face-to-face contact patterns in a primary school. *PLoS ONE*, 6(8), 2011.
- [20] TIER. <http://tier.cs.berkeley.edu/drupal/>.
- [21] Twitter. <http://www.twitter.com>.
- [22] D. Wyatt, T. Choudhury, J. Tanzeem, and J. Kitts. Inferring colocation and conversational networks using privacy-sensitive audio. *ACM Transactions on Intelligent Systems and Technology*, 2(1), 2011.
- [23] E. Yoneki and J. Crowcroft. Wireless Epidemic Spread in Dynamic Human Networks. *Bio-Inspired Computing and Communication*, LNCS(5151), 2008.
- [24] E. Yoneki, P. Hui, S. Chan, and J. Crowcroft. A socio-aware overlay for multi-point asynchronous communication in delay tolerant networks. In *Proc. MSWiM*, 2007.