

Mobile network big data for development: Applications and policy issues

Rohan Samarajiva & Sriganesh Lokanathan
Yuan Ze University, Taiwan, 8 January 2015



Our mission

Catalyzing policy change through research to improve people's lives in the emerging Asia Pacific by facilitating their use of hard and soft infrastructures through the use of knowledge, information and technology



Where we work



Why big data? Why now?

- Proximate causes
 - Increased “datafication”: Very large sets of schema-less (unstructured, but processable) data now available
 - Advances in memory technology: No longer is it necessary to archive most data and work with small subset
 - Advances in software: Hadoop, Mapreduce

But more than that . . .

- James Beniger's Control revolution (1986) provides theoretical context
 - Speed up of transportation caused crisis of control; remedied by early ICTs (telegraph, time zones)
 - As mass production brought down unit costs, another crisis of control emerged; this time remedied by mass media & advertising
 - Now, a crisis of control over attention; data analytics & micro targeting of audiences

Private & public purposes

- Data analytics in use within companies since 1990s
 - American Express was using Cray Supercomputers in 1990s
- Big push by IBM, Cisco, etc. to use big data for smart cities in 2000s
 - Rio as prototype in 2010; but used mostly video feeds and GPS
 - These models rely on proprietary software and installed-for-purpose sensors

Is there an alternative approach?

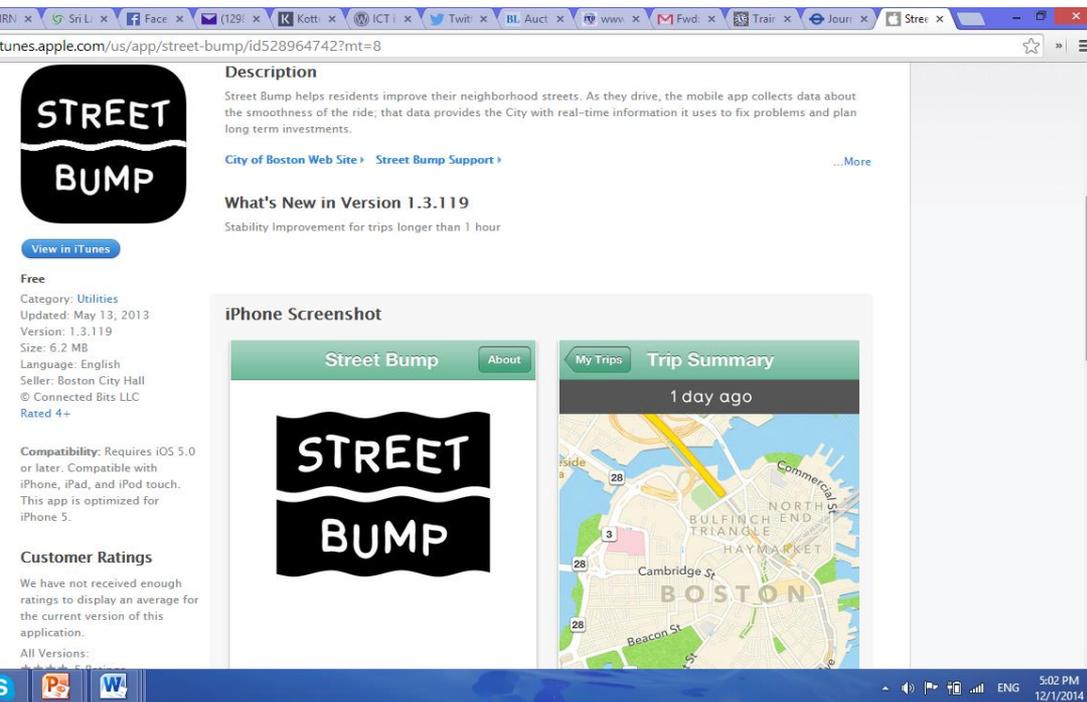
- Can we have smart cities on the cheap?
- Transaction-generated data
 - Ubiquitous mobiles can make every citizen a sensor
 - If prepaid travel cards such as Octopus & Oyster in place and public transport is popular, they too can serve as sensors
 - Open-source analytics; cheap hardware

But we need to take care

- Humans changing their behavior (Google flu trends)
- Problems of representativeness

Bias in big data → why mobile network big data in developing countries

- Streetbump is a Boston crowdsourcing + big data application that uses the natural movement of citizens to improve street maintenance
 - Data generated from an app downloaded to a smartphone “mounted” in a car



The screenshot shows the iTunes page for the Street Bump app. The page includes the app's logo, a description, a 'View in iTunes' button, and an iPhone screenshot. The iPhone screenshot shows the app's interface with a 'Street Bump' logo and a 'Trip Summary' map of Boston. The map highlights a specific trip with a yellow line. The page also lists the app's category as 'Utilities', its version as 1.3.119, and its compatibility with iOS 5.0 or later.

Can Streetbump be transplanted in Colombo at this time?

- Feature phones >> Smartphones
- “Something better than nothing” may not apply
- Bias toward roads traversed by smartphone owners → In conditions of limited resources, may skew resource allocation

Mobile network big data are more inclusive, especially in our cities

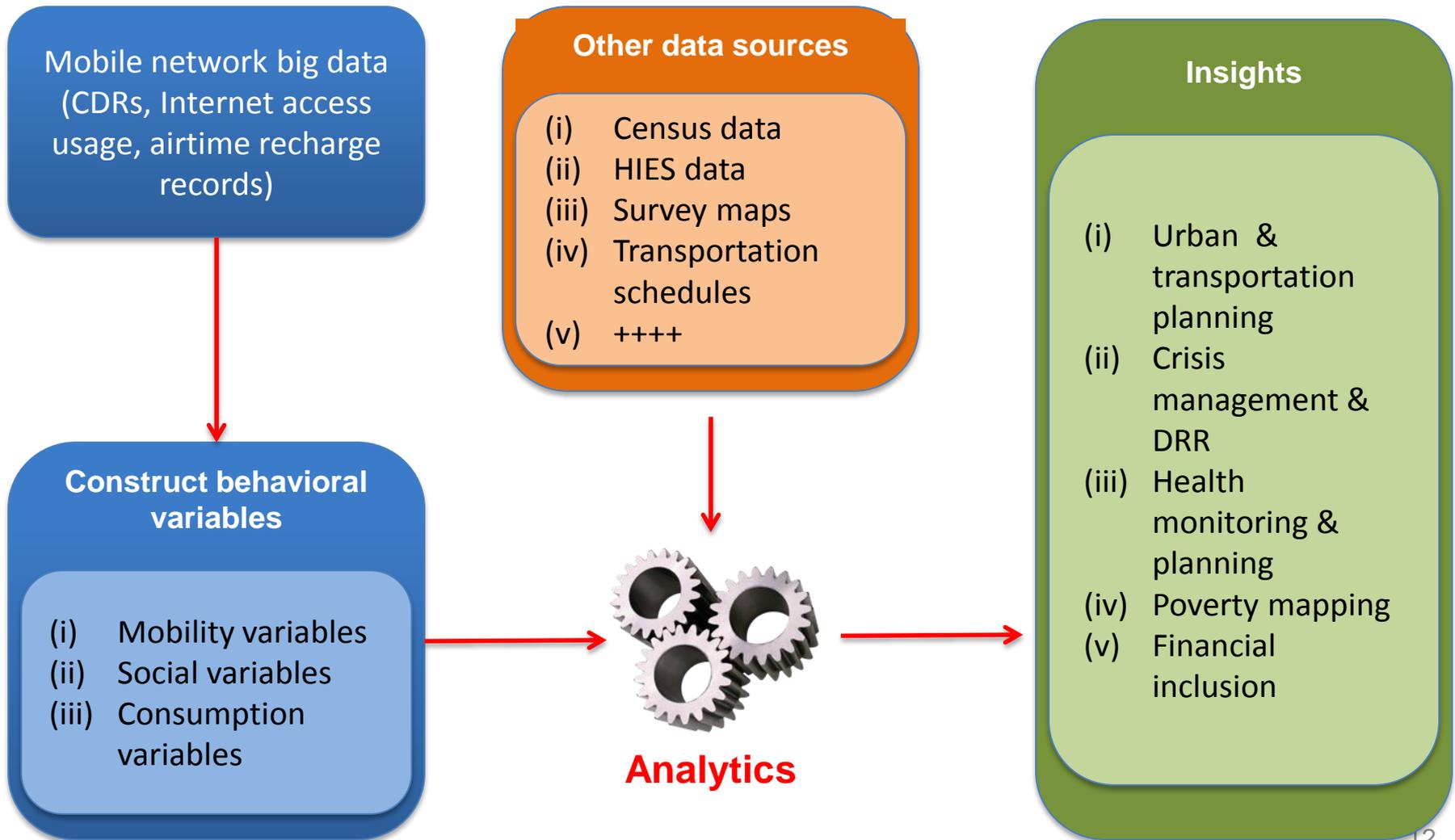
	Mobile SIMs/100	Internet users/100	Facebook users/100
Myanmar	13	1	4
Bangladesh	67	7	6
Pakistan	70	11	8
India	71	15	9
Sri Lanka	96	22	12
Philippines	105	39	41
Indonesia	122	16	29
Thailand	138	29	46

Myanmar mobile SIMs/100 was 22.6 by September 2014

There is a role for other sources of big data

- But for smart cities, MNBD is the best
 - Low cost, compared to fitting all vehicles with GPS or electronic toll cards/toll infrastructure
 - But can/should be complemented with GPS and other sensor data
 - Proposed two-year study to LK Ministry of Urban Development, based on first results with MNBD, after which appropriate sensors can be installed
 - Global Pulse analysis of food-related Twitter content in Jakarta shows value in social media content, even if not as “representative”
- Visitor Location Register (VLR) data is best for physical mobility, but Call Detail Records (CDR) can serve as acceptable proxy
 - Something we plan to explore in relation to infectious diseases in 2015

Mobile network big data + other data → rich, timely insights



Data used in the research

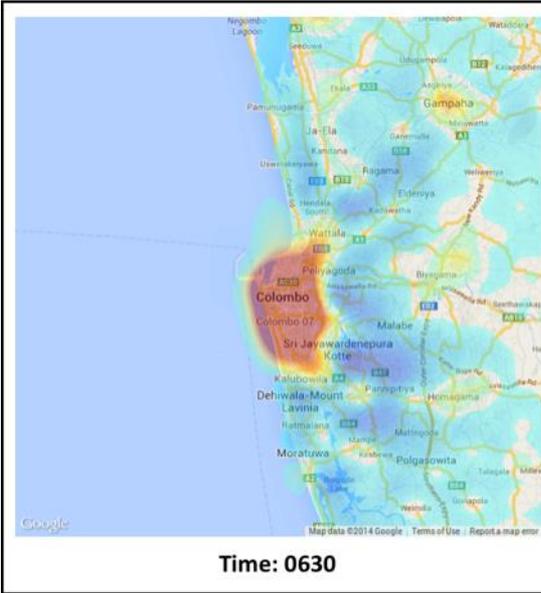
- Multiple mobile operators in Sri Lanka have provided four different types of meta-data
 - Call Detail Records (CDRs)
 - Records of calls
 - SMS
 - Internet access
 - Airtime recharge records
 - No Visitor Location Register (VLR) data
- Data sets do not include any Personally Identifiable Information
 - All phone numbers are pseudonymized
 - LIRNEasia does not maintain any mappings of identifiers to original phone numbers
- Cover 50-60% of users; very high coverage in Western (where Colombo the capital city is located) & Northern (most affected by civil conflict) Provinces, based on correlation with census data

UNDERSTANDING CHANGES IN POPULATION DENSITY

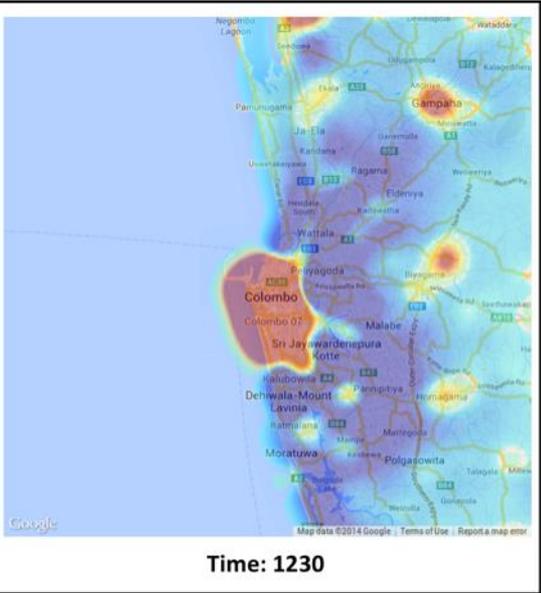
Population density changes in Colombo region: weekday/ weekend

Pictures depict the change in population density at a particular time relative to midnight

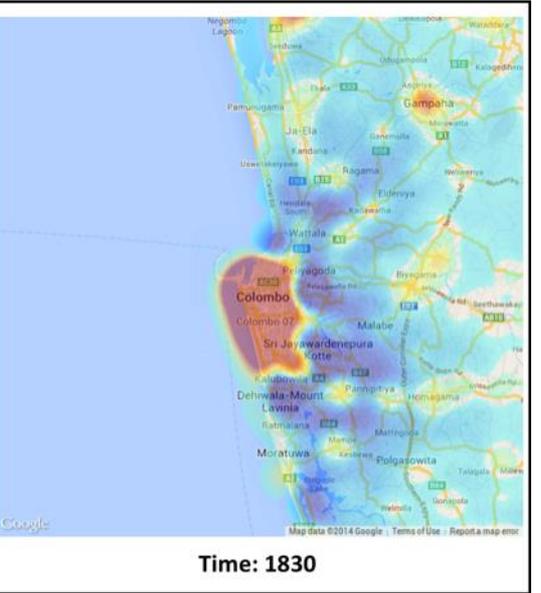
Weekday



Time: 0630

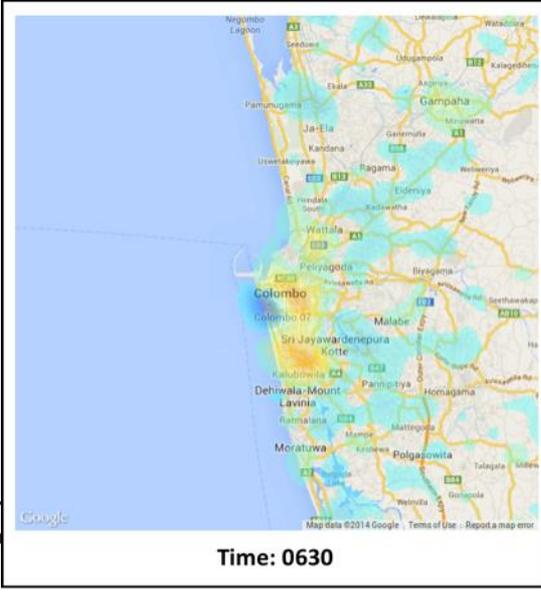


Time: 1230

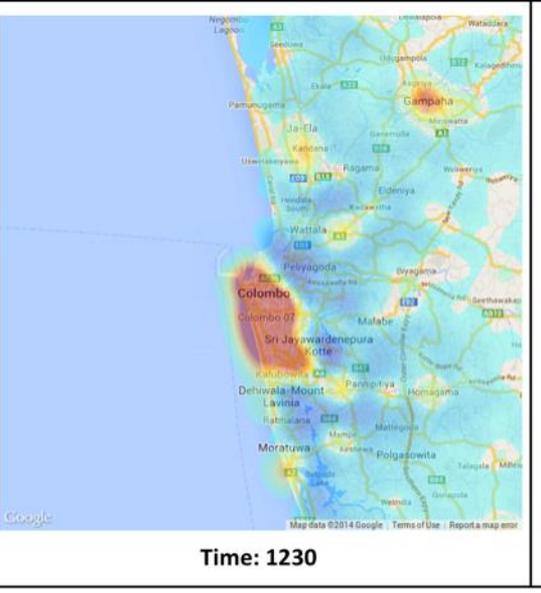


Time: 1830

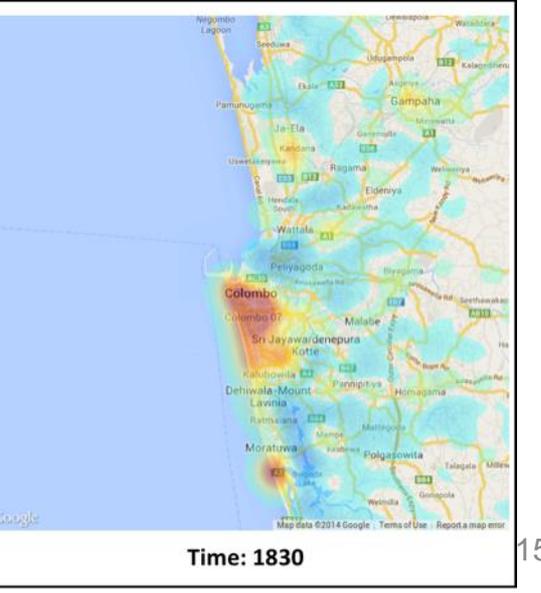
Sunday



Time: 0630



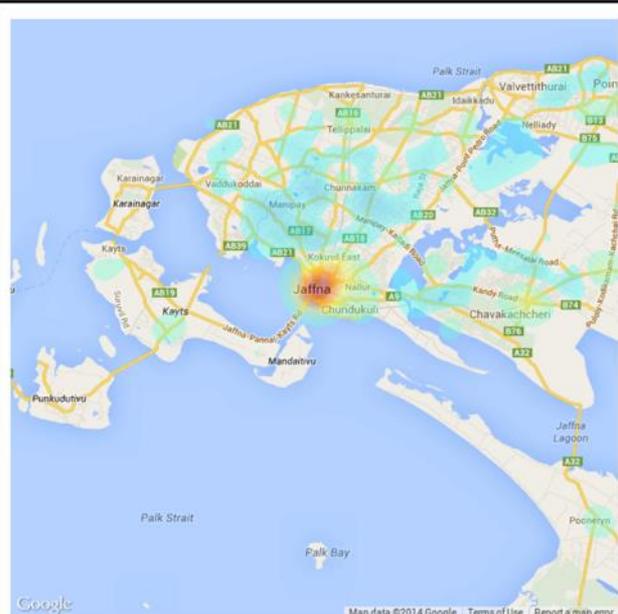
Time: 1230



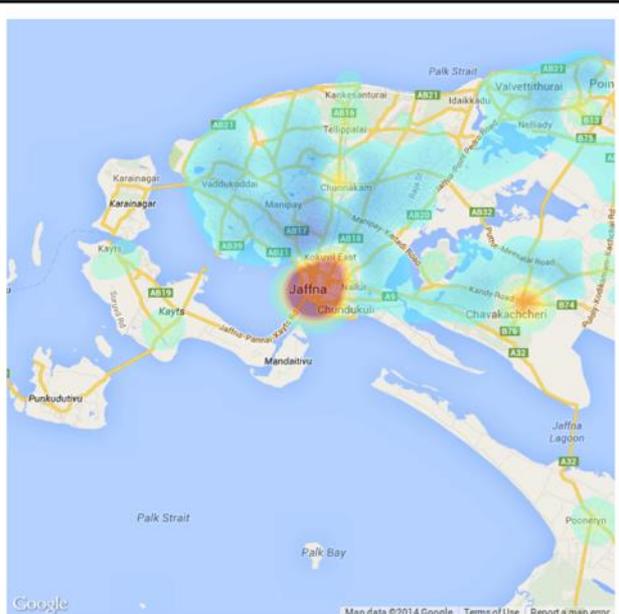
Time: 1830

Population density changes in Jaffna region on a normal weekday

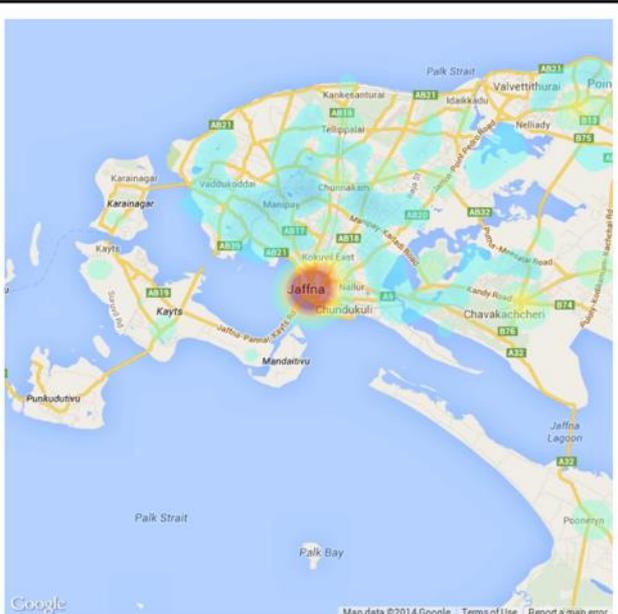
Pictures depict the change in population density at a particular time relative to midnight



Time: 0630

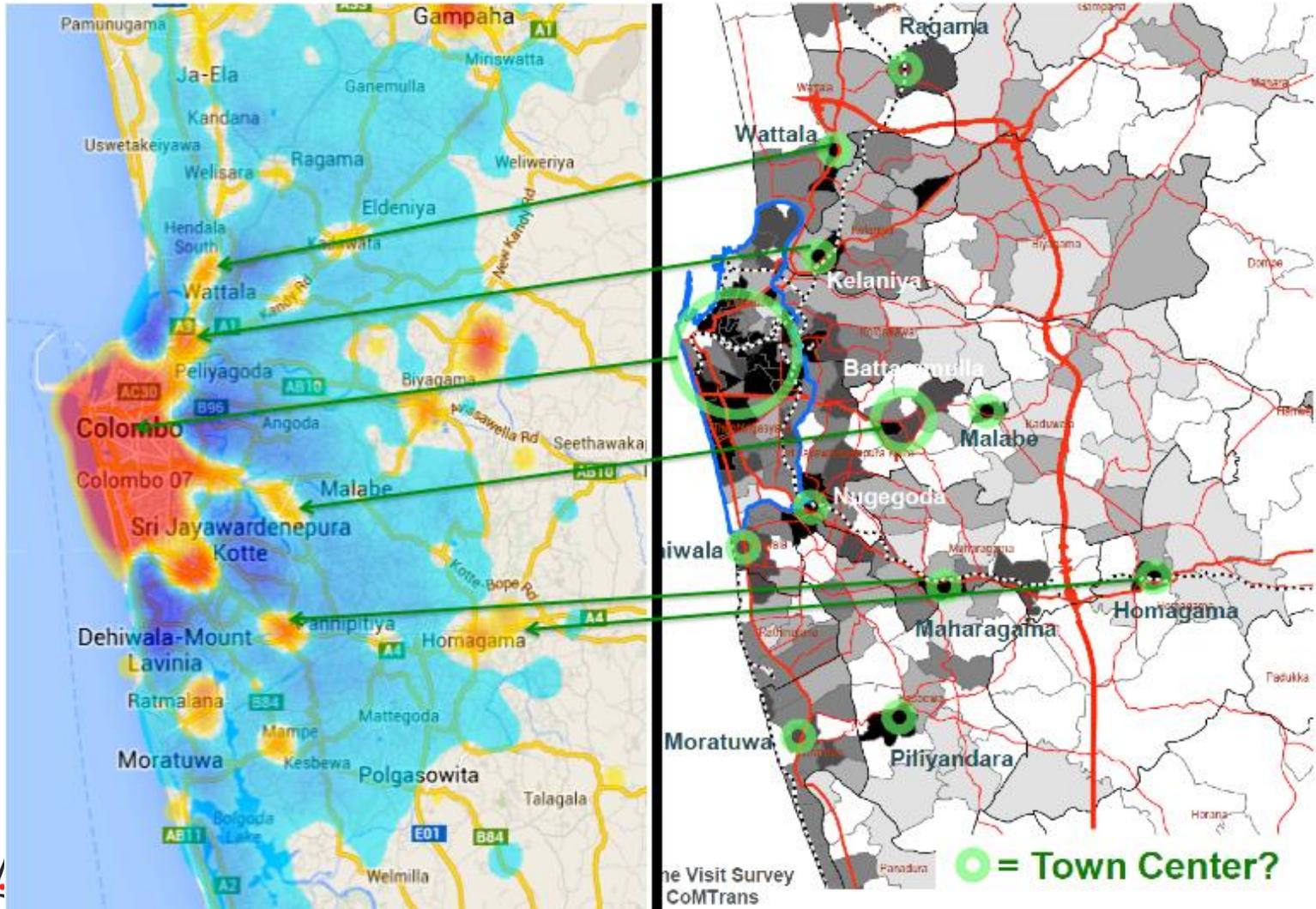


Time: 1230



Time: 1830

Our findings closely match results from expensive & infrequent transportation surveys

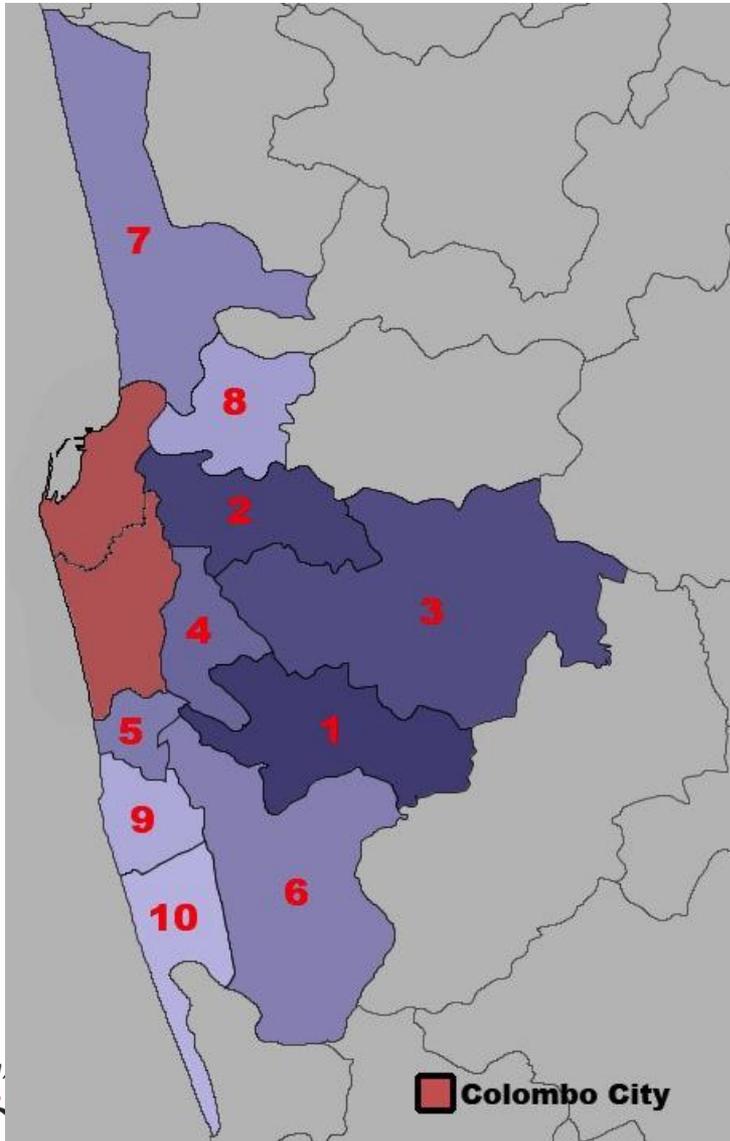


UNDERSTANDING WHERE PEOPLE LIVE AND WORK

Methodology

- Based on extracted average diurnal mobility pattern for population, choose time frames for home and work
 - Home time: 2100 to 0500
 - Work time: 1000 to 1500
- Calculate a home and work location for each SIM:
 - Match cell towers to Divisional Secretariat Division (DSD)
 - Count each DSD at most once per *day*.
 - Pick the DSD with the largest number of “hits”
 - For work consider only weekdays that are not public holidays

46.9% of Colombo city's daytime population comes from the surrounding regions



Colombo city is made up of Colombo and Thimbirigasyaya DSDs

Home DSD	%age of Colombo's daytime population
Colombo city	53.1
1. Maharagama	3.7
2. Kolonnawa	3.5
3. Kaduwela	3.3
4. Sri Jayawardanapura Kotte	2.9
5. Dehiwala	2.6
6. Kesbewa	2.5
7. Wattala	2.5
8. Kelaniya	2.1
9. Ratmalana	2.0
10. Moratuwa	1.8

Implications for public policy

Urban Planning

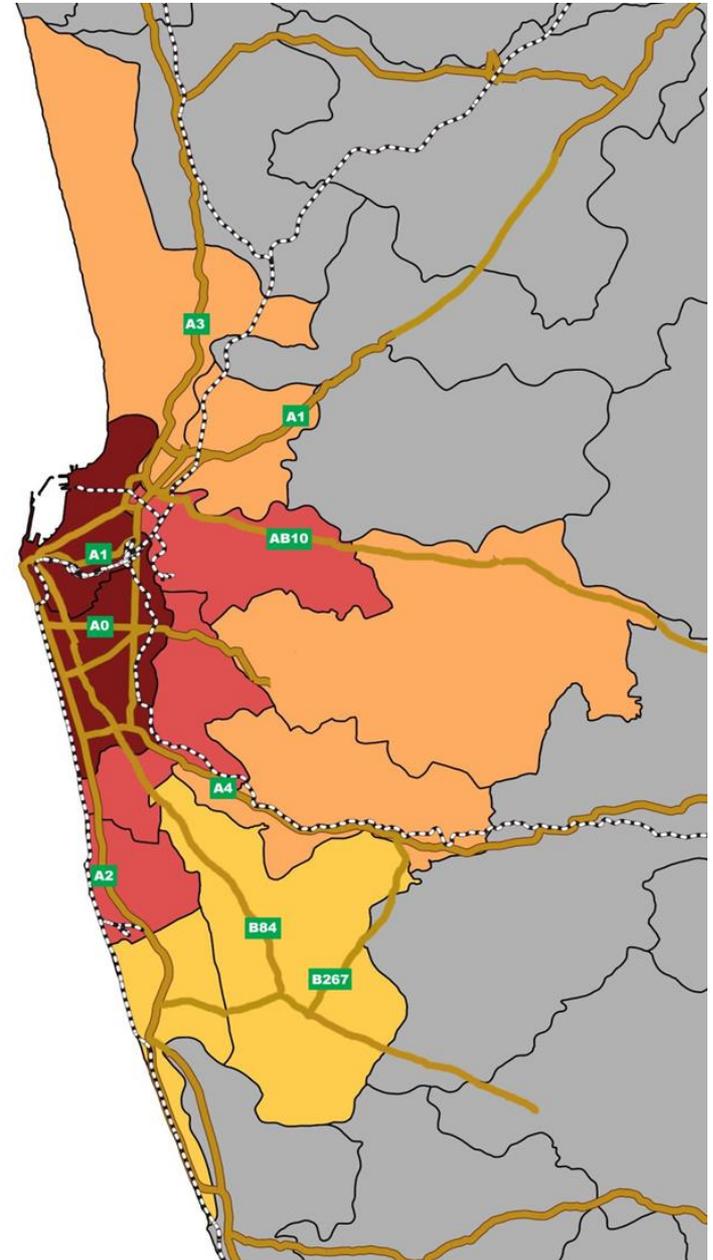
- Current municipal boundaries are obsolete; those from outside city limits cause costs but do not contribute adequate revenues; our data indicate logical boundaries of metro regions

Transportation Policy

- High volume transport corridors suitable for provision of mass transit
 - Kaduwela DSD (now served by AB 10 & A0) (3) already identified
 - High Level Road (now served by A4 & rail line) to Maharagama DSD (1)

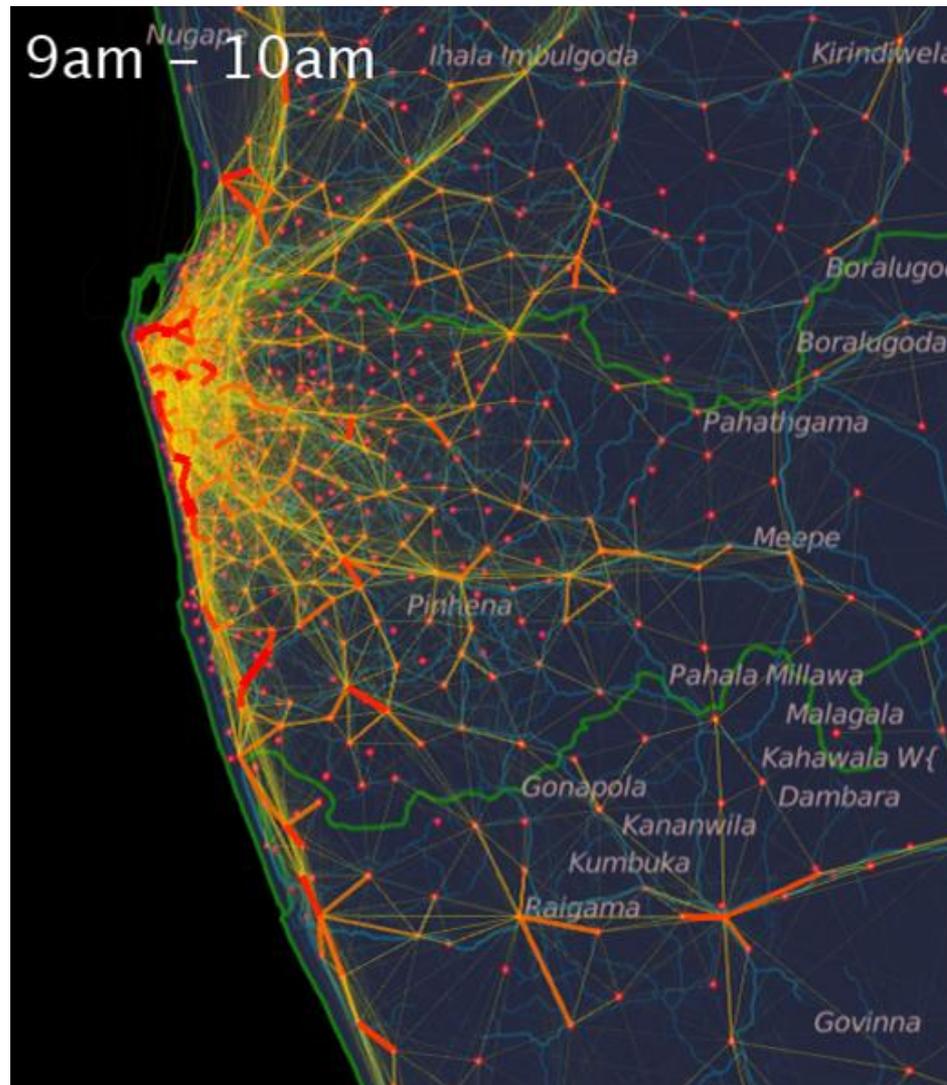
Health Policy

- Understanding people's regular mobility patterns can help model spread of infectious diseases (e.g. dengue)

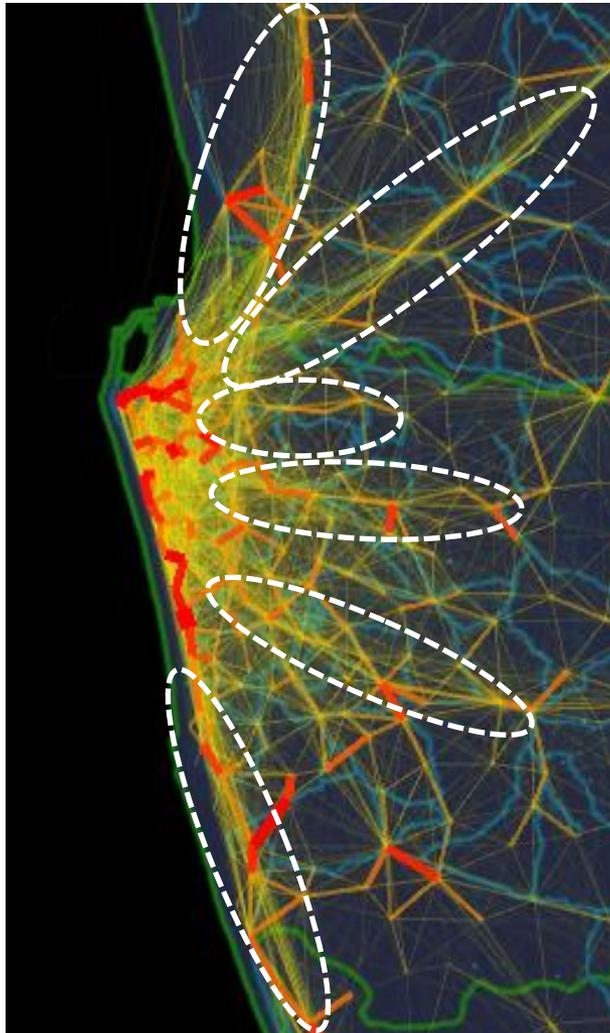


MODELING TRAVEL

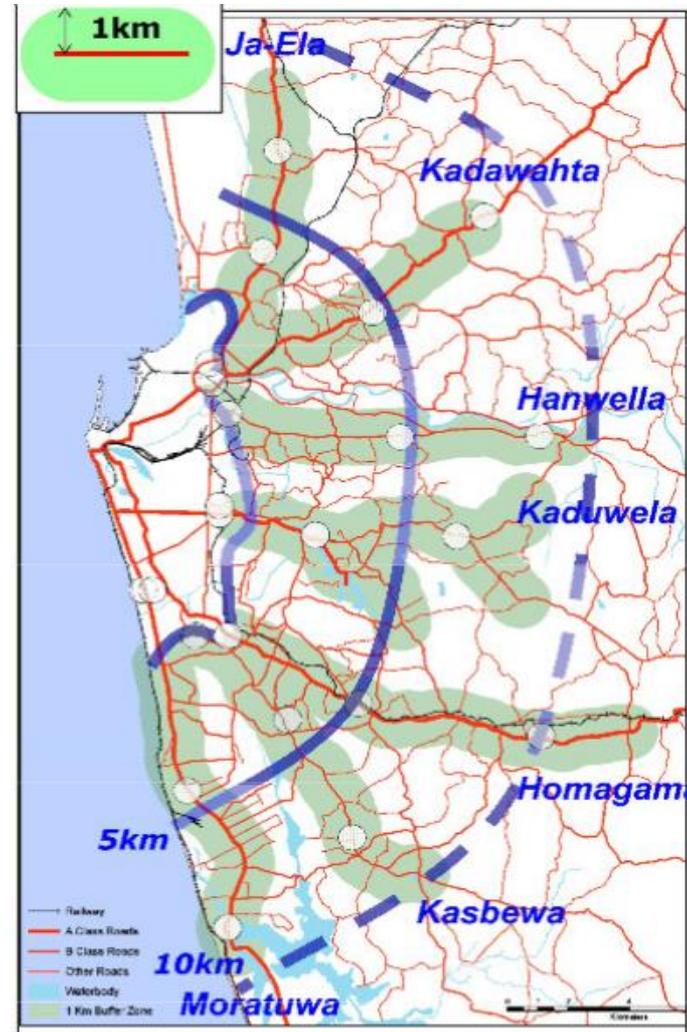
Understanding temporal variations in trips



Mobility visualization for Colombo District identifies transport corridors



Low  High
Volume of People



Source: COMTRANS report, 2013, Ministry of Transport

Implications for public policy

Transportation & Urban Policy

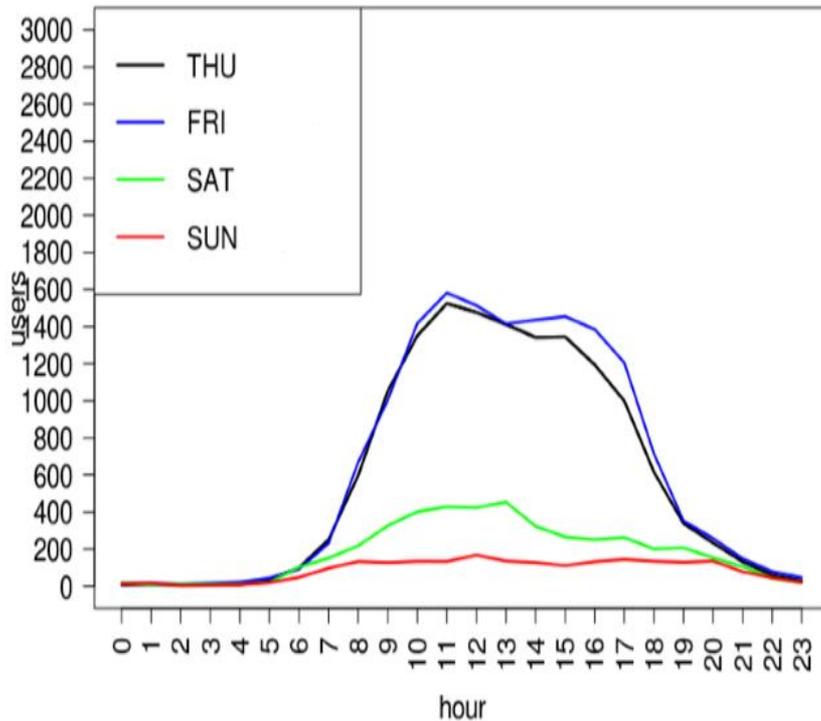
- CDR analysis can give us rough insights on principal transport corridors
- Then, with cooperation of mobile operators and additional computing power, we can zoom in on priority corridors to do detailed analysis using Visitor Location Register (VLR) data

Health Policy

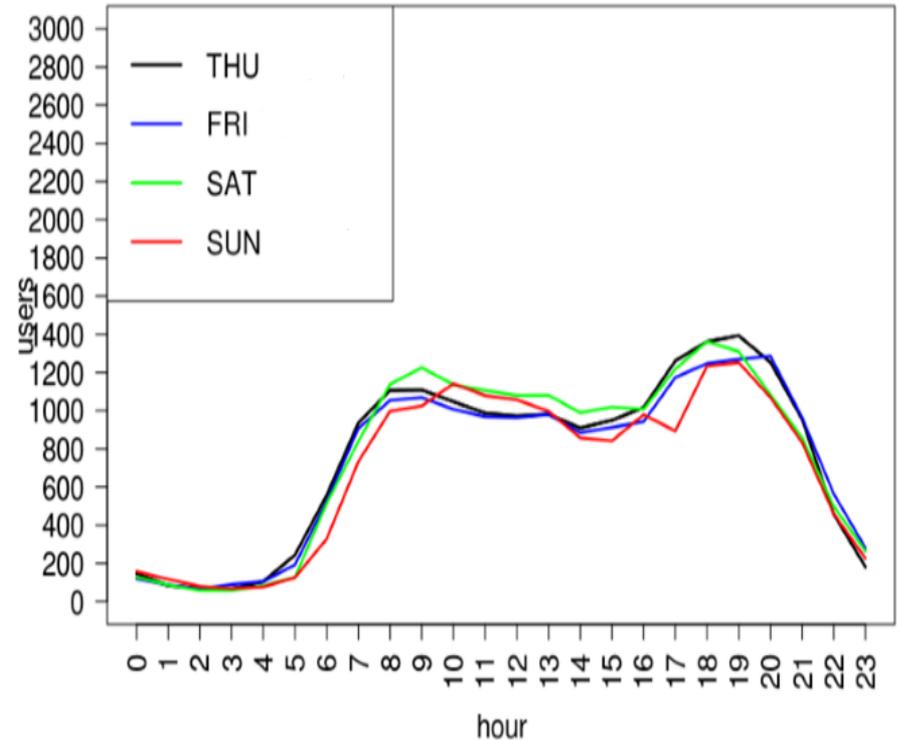
- Understanding people's regular mobility patterns can help model spread of infectious diseases (e.g. dengue)

UNDERSTANDING LAND USE CHARACTERISTICS

Hourly loading of base stations reveals distinct patterns



Type X: ?

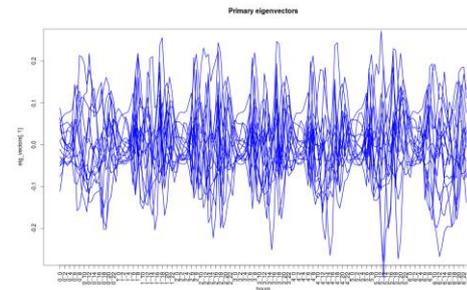
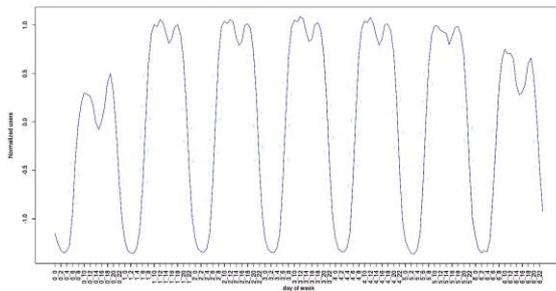


Type Y: ?

- We can use this insight to group base stations into different groups, using unsupervised machine learning techniques

Understanding land use characteristics: methodology

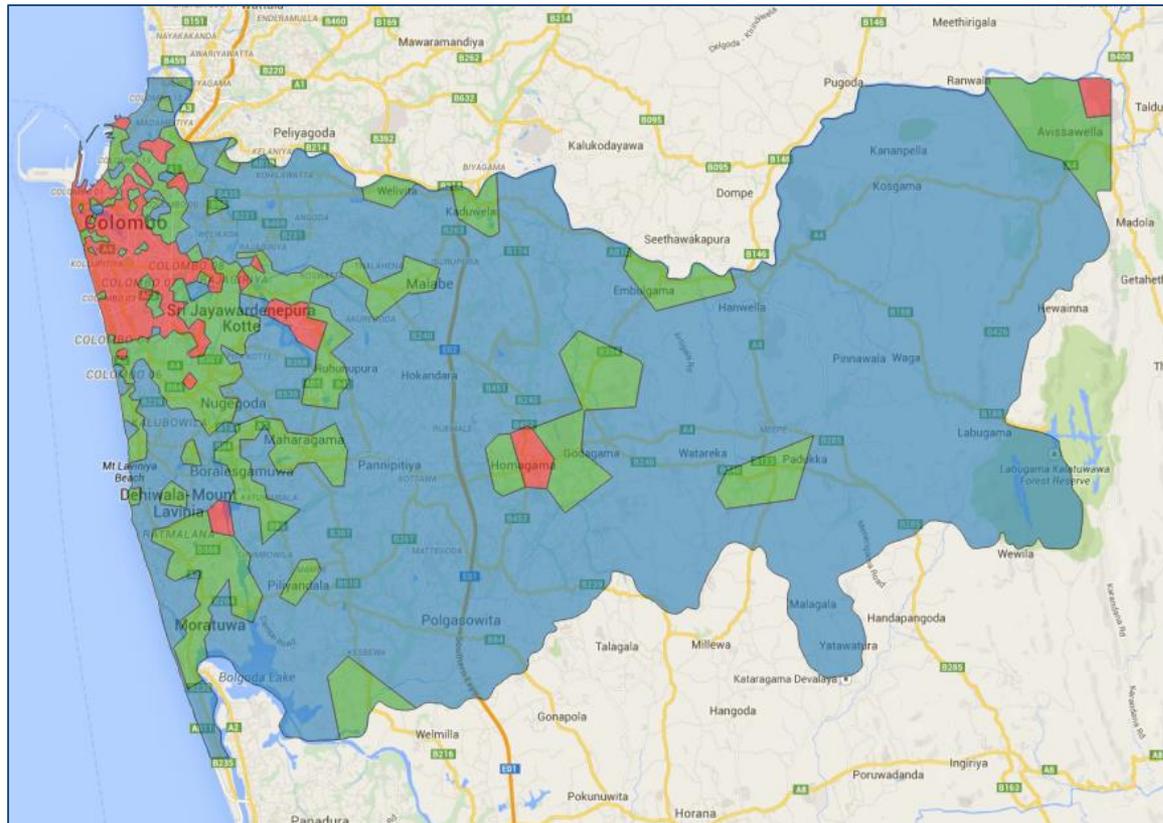
- The time series of users connected at a base station contains variations, that can be grouped by similar characteristics
- A month of data is collapsed into an indicative week (Sunday to Saturday), with the time series normalized by the z-score
- Principal Component Analysis(PCA) is used to identify the discriminant patterns from noisy time series data



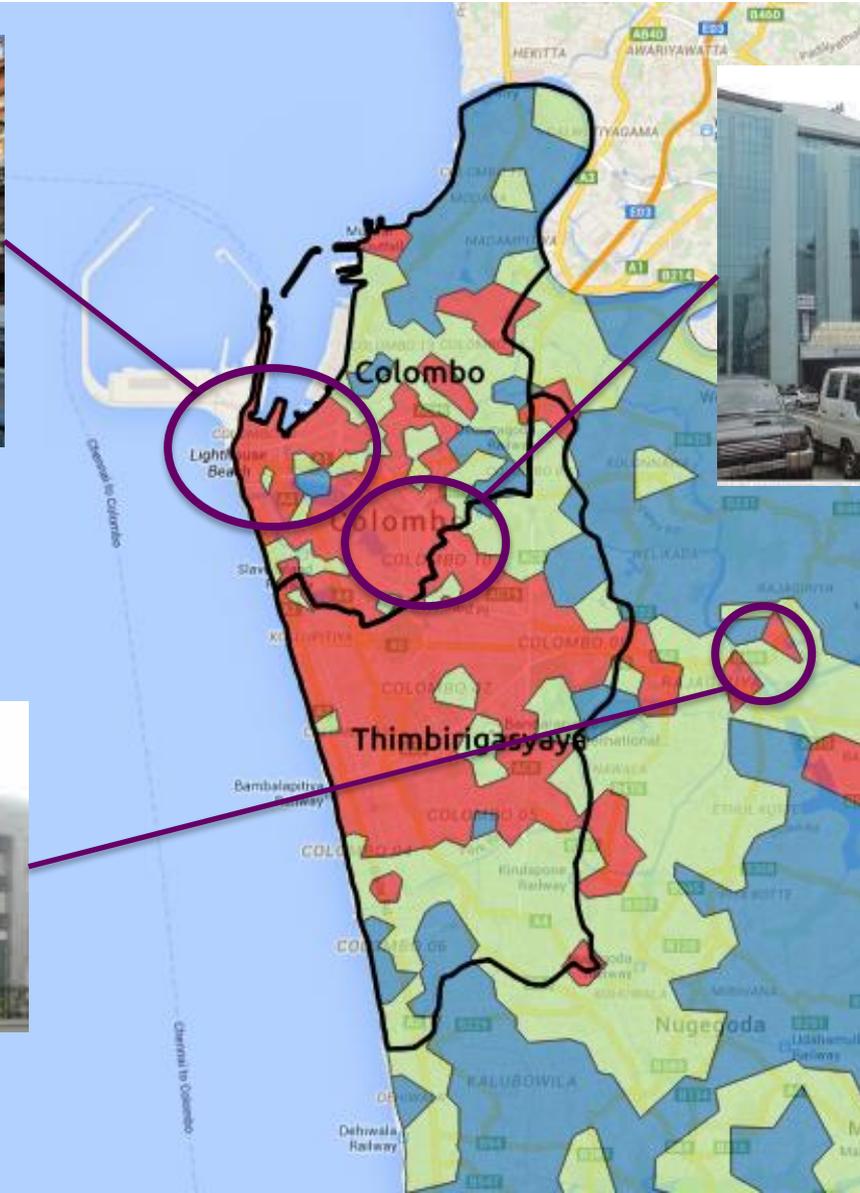
- Each base station's pattern is filtered into 15 principal components (covering 95% of the data for that base station)
- Using the 15 principal components, we cluster all the base stations into 3 clusters in an unsupervised manner using k-means algorithm

Three spatial clusters in Colombo District

- **Cluster-1 exhibits patterns consistent with commercial area**
- **Cluster-3 exhibits patterns consistent with residential area**
- **Cluster-2 exhibits patterns more consistent with mixed-use**



Our results show Central Business District (CBD) in Colombo city has expanded



Small area in NE corner of Colombo District classified as belonging to Cluster 1?

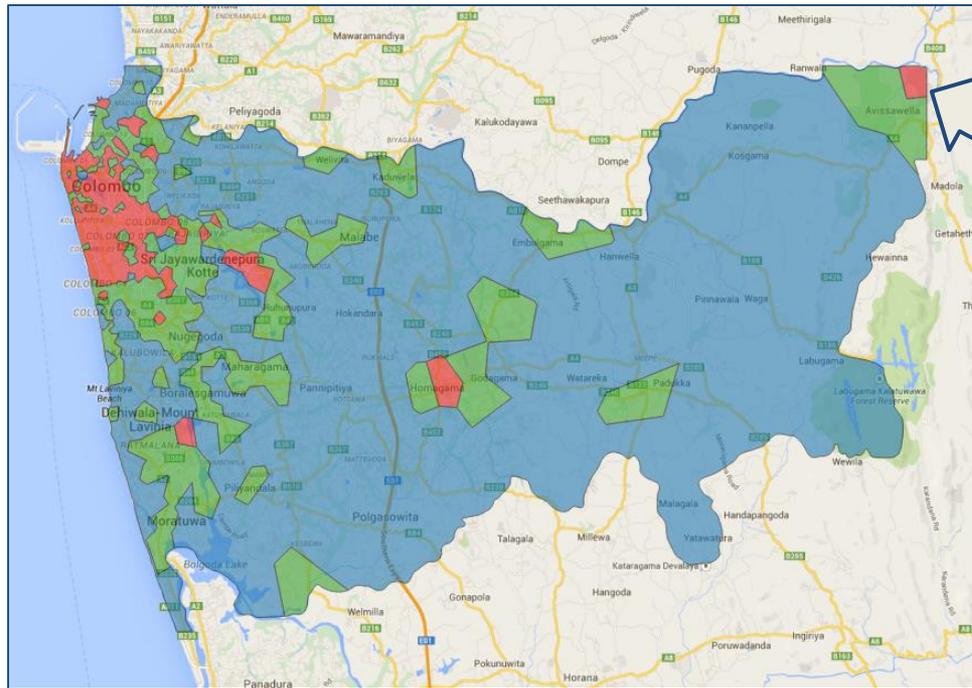
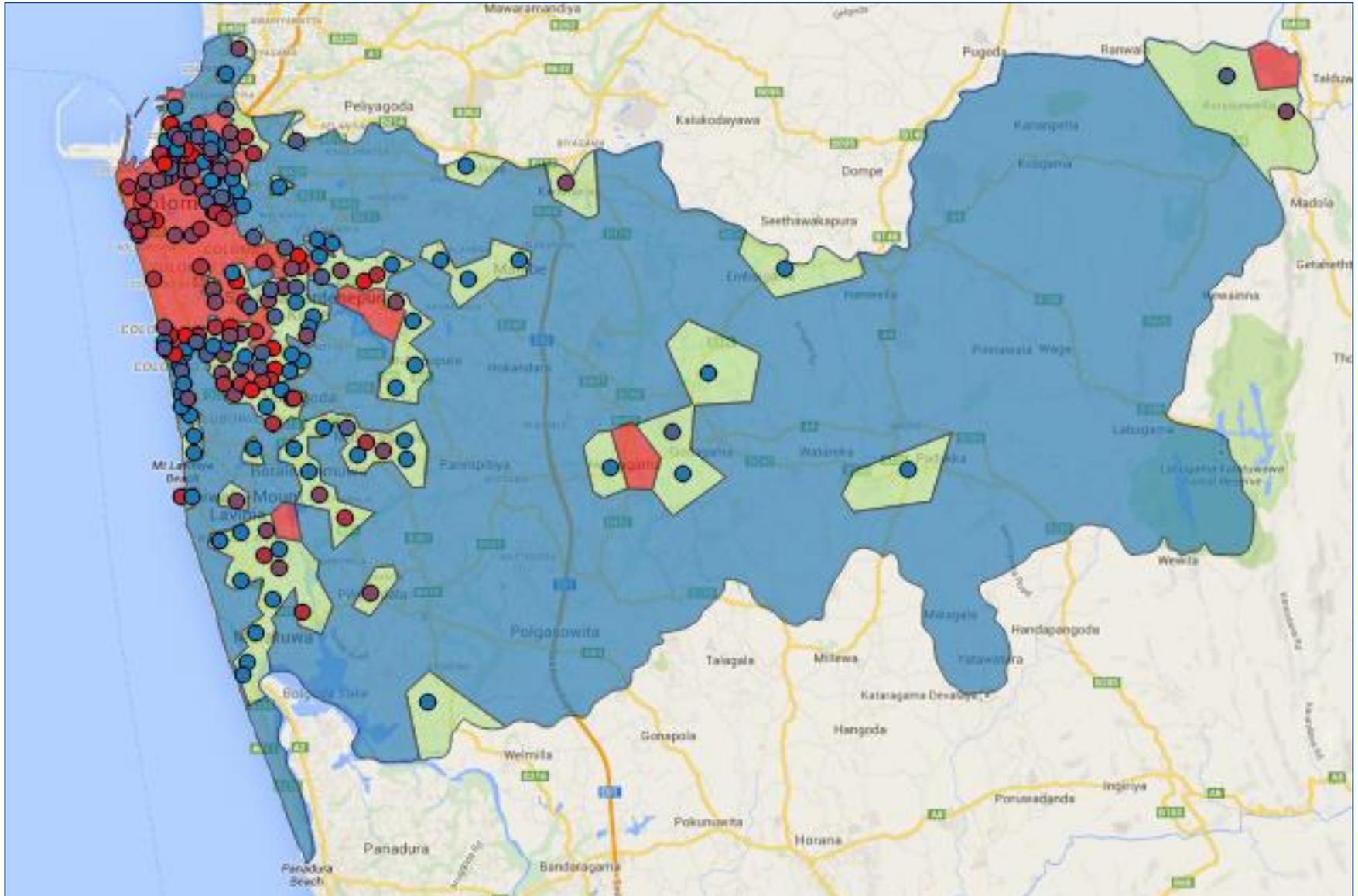


Photo ©Senanayaka Bandara - [Panoramio](#)

Seethawaka Export Processing Zone

Internal variations in mixed use regions: More commercial or more residential?



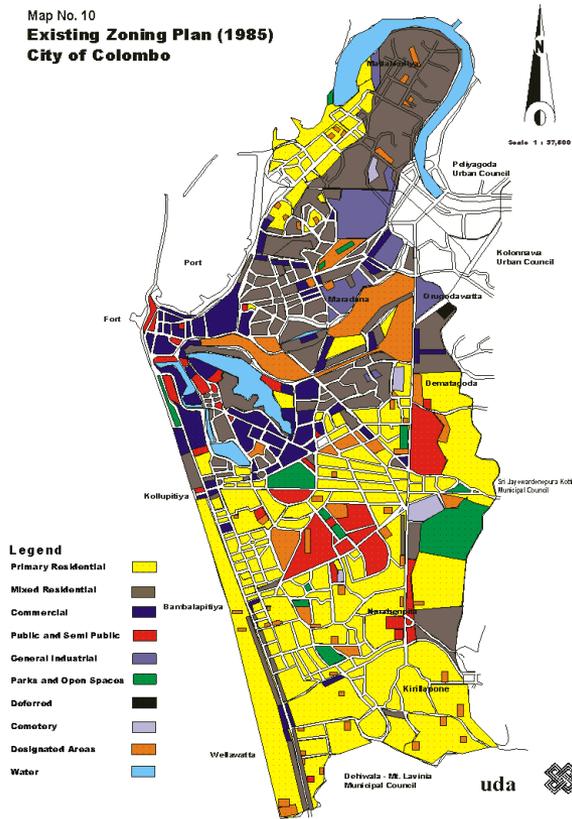
Blue dots: more residential than commercial

Red dots: more commercial than residential

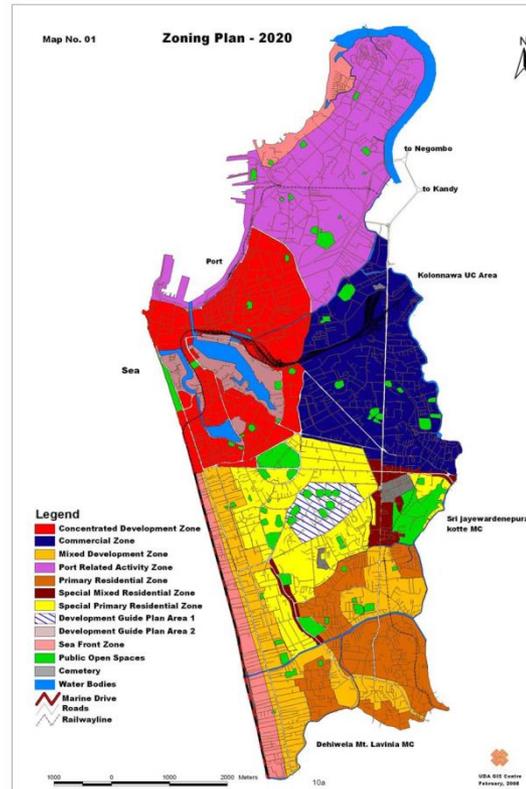
Plans & reality

1985 Plan

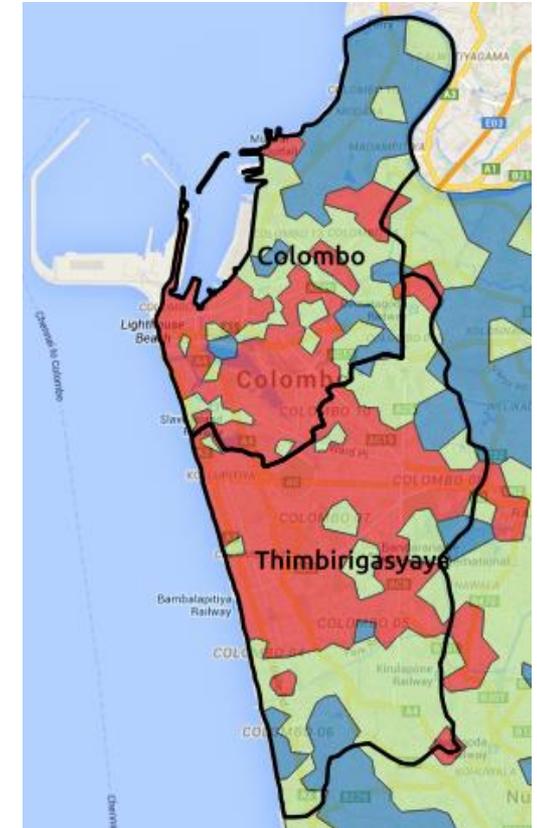
Map No. 10
Existing Zoning Plan (1985)
City of Colombo



2020 UDA Plan



2013 reality



Implications for urban policy

- Almost real-time monitoring of urban land use
 - We are currently working on understanding temporal variations in zone characteristics (especially the mixed-use areas)
- Can dispense with surveys & align master plan to reality
- LIRNEasia is working to unpack the identified categories further, e.g.,
 - Entertainment zones that show evening activity

UNDERSTANDING COMMUNITIES

Identifying communities: Methodology

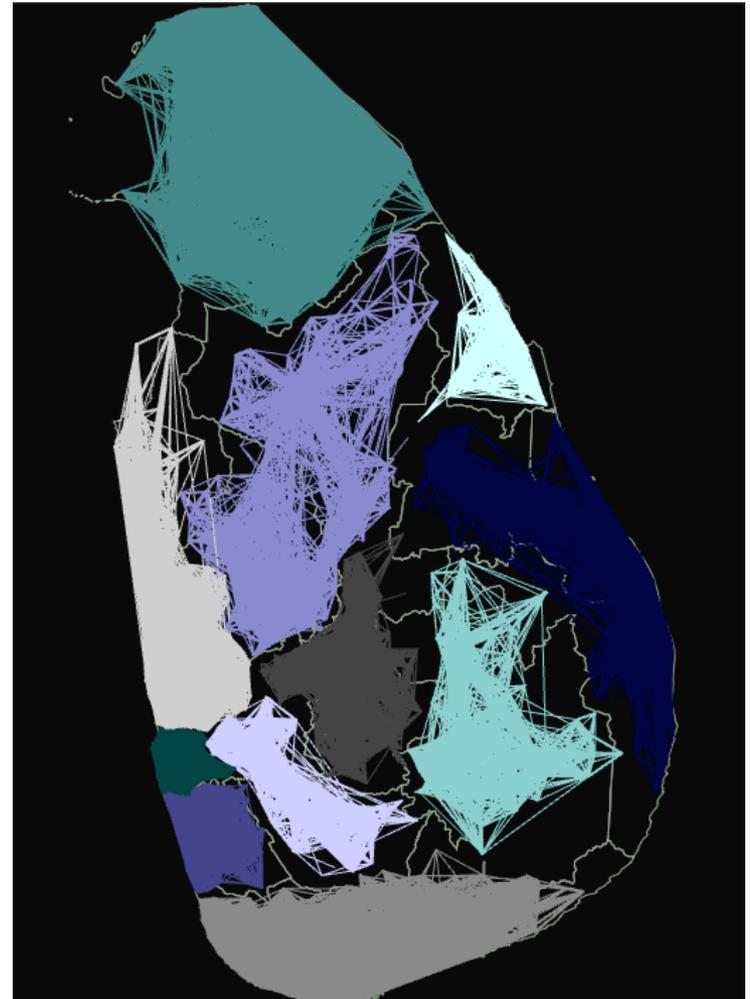
- Social networks segregated so overlapping connections between communities are minimized
- Strength of a community is determined by *modularity*
 - Modularity Q = (edges inside the community) –
(expected number of edges inside the community)

$$Q = \frac{1}{2m} \sum_{a,b} \left(A_{a,b} - \frac{k_a k_b}{2m} \right) \delta(c_a, c_b)$$

M. E. J.-Newman, Michele-Girvan, "Finding and evaluating community structure in networks", Physical Review E, APS, Vol. 69, No. 2, p. 1-16, 2004.

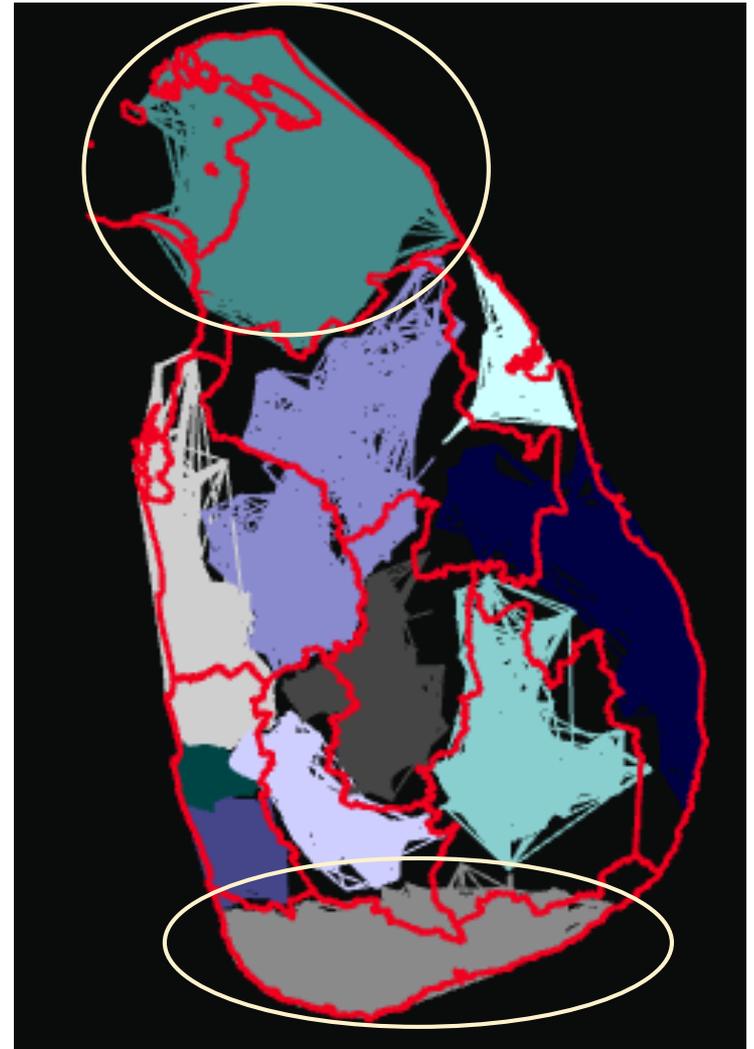
Resultant communities

- The optimal number of communities discovered by the algorithm was 11



How do these communities mesh with existing administrative boundaries?

- Southern & Northern provinces mesh the best
- Surprisingly, also Uva, hitherto thought to have a district aligned with Central plantations & other with Southern Province
- Eastern Province most intriguing
 - Trinco on its own
 - Polonnaruwa in NCP (predominantly Sinhala) tied to Batticaloa (Tamil/Muslim) and Ampara (Muslim/Sinhala) districts through rice economy
- Rest suggests administrative boundaries have been transcended



Work performed by collaborative inter-disciplinary teams

- LIRNEasia
 - Sriganesh Lokanathan
 - Kaushalya Madhawa
 - Danaja Maldeniya
 - Prof. Rohan Samarajiva
 - Dedunu Dhananjaya (lost to industry in Nov)
 - Nisansa de Silva (moved on to U of Oregon)
- LIRNEasia/ MIT
 - Gabriel Kreindler (Economics)
 - Yuhei Miyauchi (Economics)
- Technical partners:
 - WSO2 (Dr. Srinath Perera)
 - Auton Lab at Carnegie Mellon University

- University of Moratuwa
 - Prof. Amal Kumarage (Transport & Logistics Management)
 - Transport
 - Dr. Amal Shehan Perera (Computer Science & Engineering)
 - Data Mining
 - Undergraduates working on projects
- Other US Universities
 - Prof. Joshua Blumenstock (U Washington, School of Information)
 - Data Science
 - Saad Gulzar (NYU Poli Sci)
 - Political Science

Addressing challenges

Challenge	Solution(s)
Negotiating access to data	<ul style="list-style-type: none">• Win-win; insights/ techniques for public policy outputs can be leveraged for operator's business interests• Pro-active action by operator(s) rather than reactive to growing government interest in using such data
Minimizing harms from data sharing	<ul style="list-style-type: none">• Development of self-regulatory guidelines for operators
Skills	<ul style="list-style-type: none">• Assemble interdisciplinary teams that are superior to what consultants can offer
Research → policy	<ul style="list-style-type: none">• Policy enlightenment as step 1

**DRAFT GUIDELINES FOR THIRD-PARTY
USE OF MOBILE NETWORK BIG DATA**

Purpose

- Reduce transaction costs of releasing mobile network big data (MNBD) to third parties for public and commercial purposes

First step in a process that will hopefully lead to the adoption of a voluntary code of conduct by the region's mobile network operators (MNOs) that will be the most effective in minimizing possible harms

Method

- Potential harms have been identified through
 - the literature (Annex 1) and
 - engagement with ongoing analysis of MNBD at LIRNEasia

Anchored on my work on utility transaction-generated data since 1991

Privacy and other harms, from the ground up

- Guidelines address harms that have emerged in society and recognized as worthy of remedy in the Common Law and not on abstract principles.
- Solove (2008: 174) argues that privacy as an abstract concept is difficult to pin down, since it “involves a cluster of protections against a group of different but related problems.”
- He identifies 16 privacy problems, grouped into four general types:
 - Information collection;
 - Information processing;
 - dissemination; and
 - invasion

Considered harms

- Privacy (9 out of 16 recognized in the Common Law in multiple countries)
 - Surveillance
 - Aggregation
 - Identification, individual and group
 - Insecurity
 - Secondary use
 - Exclusion
 - Breach of confidentiality
 - Disclosure
 - Increased accessibility
- Anti-competitive effects
- Marginalization

Remedy	Identified potential harm	Include in agreements transferring identifiable MNBD	Include in agreements transferring anonymized MNBD
<p>Mobile Network Operators (MNOs) will not engage in active surveillance of their customers, except as required by applicable law. MNOs will desist from collecting more data than are needed for the efficient operation of the networks and the supply of good service to customers. To the extent feasible, data collection practices will be transparent.</p>	<p><i>Active surveillance</i></p>	<p>No. Applying only to MNOs, this need not be included in agreements. However, active surveillance is a root cause of problems that could be manifested in other forms at the subsequent information processing and dissemination phases.</p>	<p>No. Applying only to MNOs, this need not be included in agreements. However, active surveillance is a root cause of problems that could be manifested in other forms at the subsequent information processing and dissemination phases.</p>

Remedy	Identified potential harm	Include in agreements transferring identifiable MNBD	Include in agreements transferring anonymized MNBD
<p>Best efforts will be made to prevent de-anonymization.</p> <p>Working groups may be formed with data users to monitor the state of knowledge in techniques of anonymization and de-anonymization.</p>	<p><i>De-anonymization</i></p>	<p>No</p>	<p>Yes</p>

Remedy	Identified potential harm	Include in agreements transferring identifiable MNBD	Include in agreements transferring anonymized MNBD
<p>Individually identifiable data will not be released to third parties, unless the purposes are specified in the agreement and have been approved by an ethics review committee, if one is available. If an ethics review committee is not available, an equivalent third-party review should be sought.</p>	<p><i>Individual identification</i></p>	<p>Yes</p>	<p>No</p>

Remedy	Identified potential harm	Include in agreements transferring identifiable MNBD	Include in agreements transferring anonymized MNBD
Any agreement transferring identifiable data to a third party will also transfer responsibility to maintain safeguards to ensure security of individually identifiable data.	<i>Insecurity</i>	Yes	No

Remedy	Identified potential harm	Include in agreements transferring identifiable MNBD	Include in agreements transferring anonymized MNBD
The agreement governing the transfer will include provisions to minimize risks posed by increased accessibility when data are released to third parties.	<i>Increased accessibility</i>	Yes	Yes

Remedy	Identified potential harm	Include in agreements transferring identifiable MNBD	Include in agreements transferring anonymized MNBD
<p>The principle of non-discrimination shall govern the release of MNBD to third parties who do not compete with the MNO. Those in the same class will be treated equally, subject to reasonable accommodation for resource constraints.</p>	<p><i>Anti-competitive effects</i></p>	<p>Yes</p>	<p>Yes</p>