



Opportunities and Requirements for Leveraging Big Data for Official Statistics and the SDGs in Latin America

Emmanuel Letouzé

Director and co-Founder, Data-Pop Alliance
Visiting Scholar, MIT Media Lab

Eight meeting of the Statistical Conference of the Americas of ECLAC

Substantive seminar: “The data revolution and the 2030 Agenda for Sustainable Development: Challenges and Opportunities for National Statistical Institutes”

Quito, November 17, 2015



Thanks

- **ECLAC:** Alicia Barcena, Romain Zivy, Pilar Arturo, Paula Fuenzalida, Maria Eugenia Johnson...
- **UNFPA:** Pablo Salazar, Sabrina Juran...
- **INEC:** Carlos Riva,...
- All other organizers..



DATA-POP ALLIANCE
WHITE PAPER SERIES

Opportunities and Requirements for Leveraging Big Data for Official Statistics and the Sustainable Development Goals in Latin America

FULL DRAFT V1 FOR DISCUSSION

November 2015

Please send comments & suggestions to eletouze@datapopalliance.org

DATA-POP ALLIANCE
WHITE PAPERS SERIES

OFFICIAL STATISTICS,
BIG DATA AND
HUMAN DEVELOPMENT

March 2015

Emmanuel Letouzé
Johannes Jütting



Both a new and old topic...

- 2 years ago...feels like 10 years ago..



Big Data and Official Statistics: Perspectives From the Case of Poverty Monitoring

Emmanuel Letouzé*
eletouze@berkeley.edu

ISI World Statistical Congress
Special Technical Session 18 "Big Data"
Hong Kong, August 29th, 2013

60th Anniversary of the National Department of Statistics of Colombia (DANE)

Official Statistics in the Big Data Era

Emmanuel Letouzé
PhD Candidate, UC Berkeley
Fellow, Harvard Humanitarian Initiative
Non-Resident Adviser, International Peace Institute
eletouze@berkeley.edu

Biblioteca Luis Angel Arango
Bogotá
October 28th, 2013

'Data is the new oil'—the rush

"We are at the beginning of what I call **The Industrial Revolution of Data.**"

Joe Hellerstein,
November 19, 2008

WIRED MAGAZINE: 16.07

SCIENCE + DISCOVERY
The End of Theory: The Data Deluge Makes the Scientific Method Obsolete



Big data: The next frontier for innovation, competition, and productivity



WORLD ECONOMIC FORUM
COMMITTED TO IMPROVING THE STATE OF THE WORLD

Big Data, Big Impact:
New Possibilities for International Development

The Economist



SOCIAL SCIENCE

Computational Social Science

David Lazear,¹ Alex Pentland,² Lada Adamic,³ Sinan Aral,^{2,4} Albert-László Barabási,⁵ Devon Brewer,⁶ Nicholas Christakis,¹ Noshir Contractor,⁷ James Fowler,⁸ Myron Gutmann,⁹ Tony Jebara,¹ Gary King,¹ Michael Macy,¹⁰ Deb Roy,⁷ Marshall Van Alstyne^{2,11}

Statistics 2.0
The next level

By Enrico Giovannini

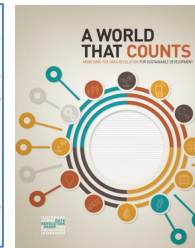
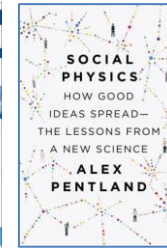


Africa's statistical tragedy

Big Data for Development:
Challenges & Opportunities

May 2012

6 provocations



Opportunities and Requirements for Leveraging Big Data for Official Statistics and the Sustainable Development Goals in Latin America
DATA-POP ALLIANCE WHITE PAPER STUDY
FULL DRAFT V1 FOR DISCUSSION
November 2015

DATA-POP ALLIANCE WHITE PAPER STUDY
OFFICIAL STATISTICS, BIG DATA AND HUMAN DEVELOPMENT
March 2015

Data and development
Off the map

Rich countries are deluged with data; developing ones are suffering from drought

Nov 15th 2014 | NEW YORK AND LONDON | From the print edition

Returns for « Big Data » over time on Google

2008

2009

2010

2011

2012

2013

2014



Main message:

*“Big Data as a new **ecosystem** presents the **official statistical community** with a **historical opportunity** and a **political obligation** to **engage with it** to retain or regain its role as the **legitimate custodian of statistical knowledge** and creator of a **deliberative public space** to discuss and drive human development in and about societies on the basis of **sound democratic and statistical principles**.”*

DATA-POP ALLIANCE WHITE PAPER SERIES

Opportunities and Requirements
for Leveraging Big Data for
Official Statistics and the
Sustainable Development Goals
in Latin America

FULL DRAFT V1 FOR DISCUSSION

November 2015

DATA-POP ALLIANCE WHITE PAPERS SERIES

OFFICIAL STATISTICS,
BIG DATA AND
HUMAN DEVELOPMENT

March 2015

60th Anniversary of the National Department of Statistics of Colombia (DANE)

Official Statistics in the Big Data Era

Emmanuel Letouzé
PhD Candidate, UC Berkeley
Fellow, Harvard Humanitarian Initiative
Non-Resident Advisor, International Peace Institute
elouz@berkeley.edu

Biblioteca Luis Angel Arango
Bogotá
October 28th, 2013

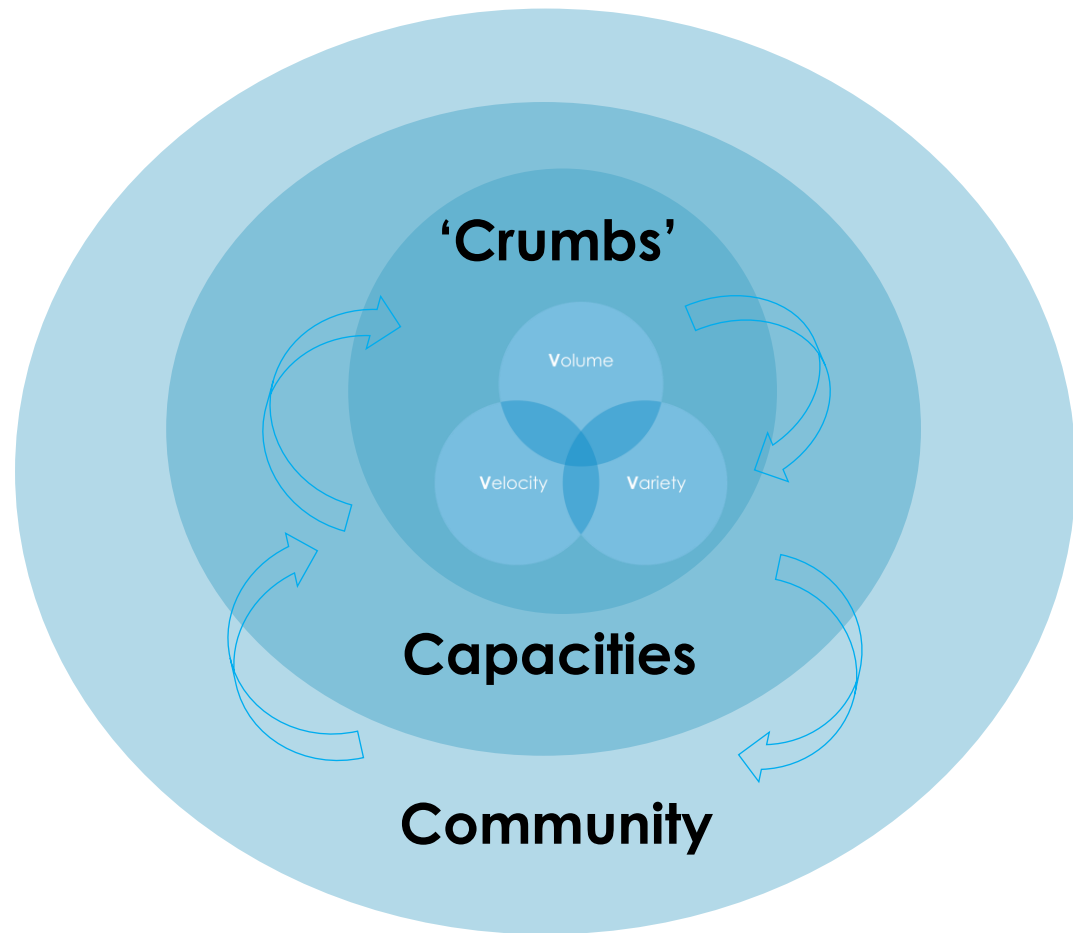
Big Data and Official Statistics:
Perspectives From the Case of Poverty Monitoring

Emmanuel Letouzé*
elouz@berkeley.edu

ISI World Statistical Congress
Special Technical Session 18 "Big Data"
Hong Kong, August 29th, 2013

Emmanuel Letouzé
Johannes Jüttgen

What is Big Data? *An ecosystem*



Made up of the **3Cs of Big Data**

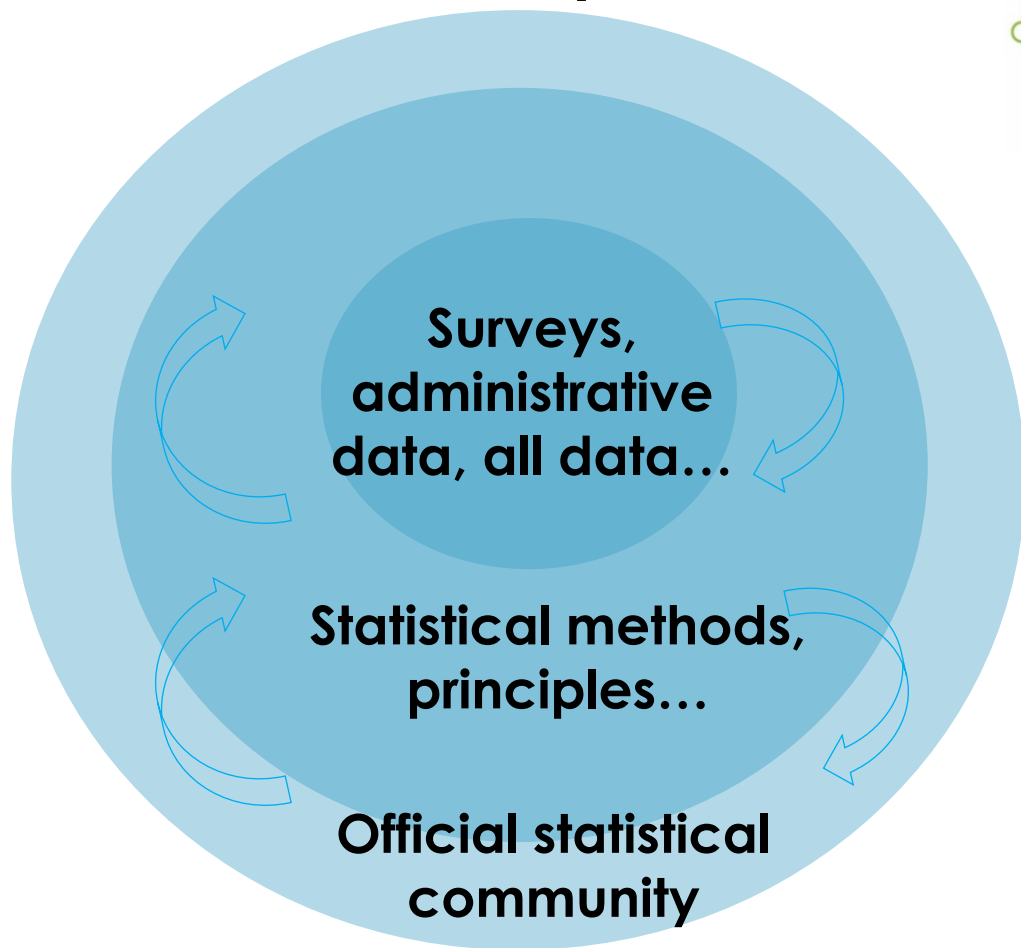
Serving **4 main applications**:

1. **Descriptive** (e.g. maps);
2. **Predictive**
 1. now-casting' or inference
 2. forecasting
3. **Prescriptive**, by establishing causal relations
4. **Discursive**; by spurring and shaping dialogue within and between communities

What is Official Statistics? also an ecosystem



An ecosystem that turns **data into knowledge** and provides a **space to drive and discuss** development outcomes **according to sound democratic and statistical principles**



Statistics 2.0
The next level

By Enrico Giovannini

What is Big Data for Official Statistics?...



*“For the Official Statistical Community, **Big Data is a wake up call**”*

*John Pullinger
UK National Statistician, Chair of the 46th
UN Statistical Commission Session,
Data-Pop Alliance—ODI-RSS event,
January 19th, 2015*



Another way to say it: ***a historical opportunity and political obligation***

DATA-POP ALLIANCE
WHITE PAPER SERIES

Opportunities and Requirements
for Leveraging Big Data for
Official Statistics and the
Sustainable Development Goals
in Latin America

FULL DRAFT V1 FOR DISCUSSION

November 2015

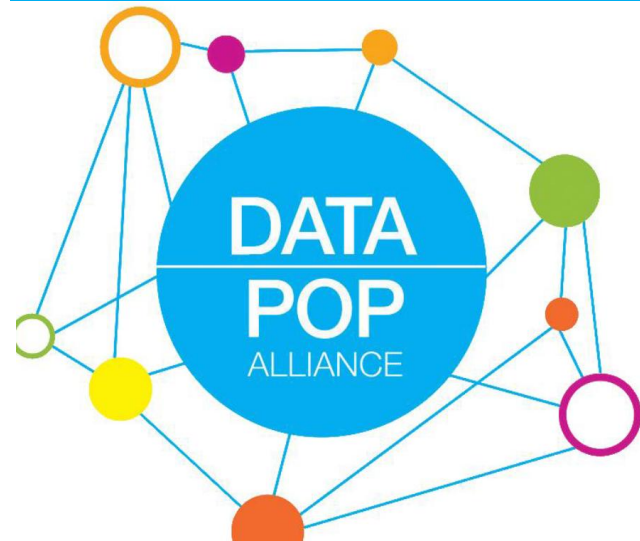


Table of Contents

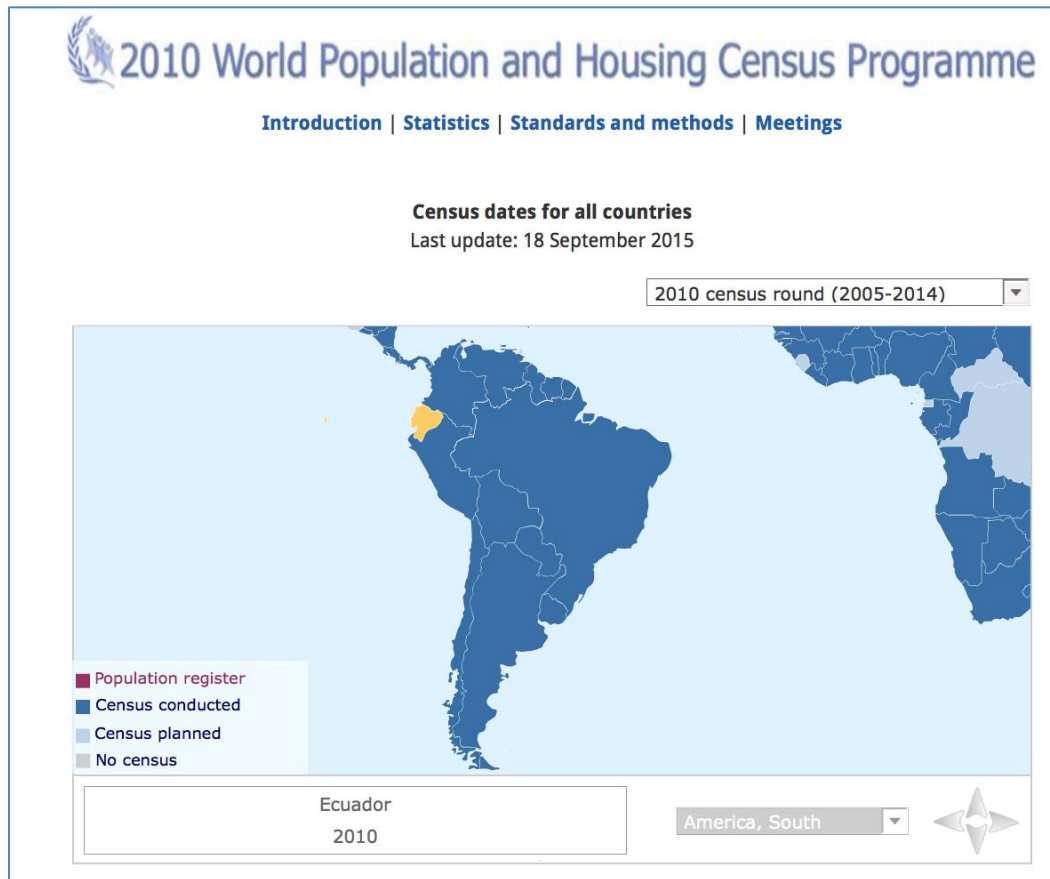
Foreword.....	5
Introduction.....	1
1 The state of Latin American NSOs: overall context and concepts	5
1.1. The role of national statistical offices in Latin America and the Caribbean	5
1.2. The state of NSOs in LAC: ongoing challenges	6
1.3. Defining “Big Data” for Official Statistics and the SDGs	8
1) Big Data provides new sources of data	9
2) Big Data provides greater diversity of data sources	9
3) Big Data has the potential to complement and improve ongoing statistical activities through its four functions	9
2 Engaging, Innovating, and Discovering Big Data in Latin America	12
2.1 Setting the stage: the burgeoning ecosystem of Big Data	12
2.2 NSOs and Big Data: trends in Latin America	15
2.3 Big Data for SDGs across the wider ecosystem of actors	19
1) Big Data Research Projects	19
1) Governments and international agencies	19
2) Private sector approaches	23
3) Civic Technology Movement	23
2.4 International attempts to use Big Data for official statistics and development	24
3 Challenges and requirements for NSOs engaging Big Data for SDGs	25
3.1. Institutional barriers to innovation and change management	25
3.2. Constraints to data access and completeness	26
3.3. Technical challenges	27
3.4 Human capacity gaps	28
3.5. Methodological challenges	30
3.6. Ethical and political risks	31
4 Towards a Regional Multi-Partner Roadmap for Leveraging Big Data for Official Statistics and the SDGs	34
4.1. Five regional trends promoting Big Data use in Latin America	34
4.2. Towards a regional multi-partner roadmap for Big Data: building on existing regional strengths and opportunities	36
Annex	41
Endnotes	53



The State of Official statistics in LAC: B+?



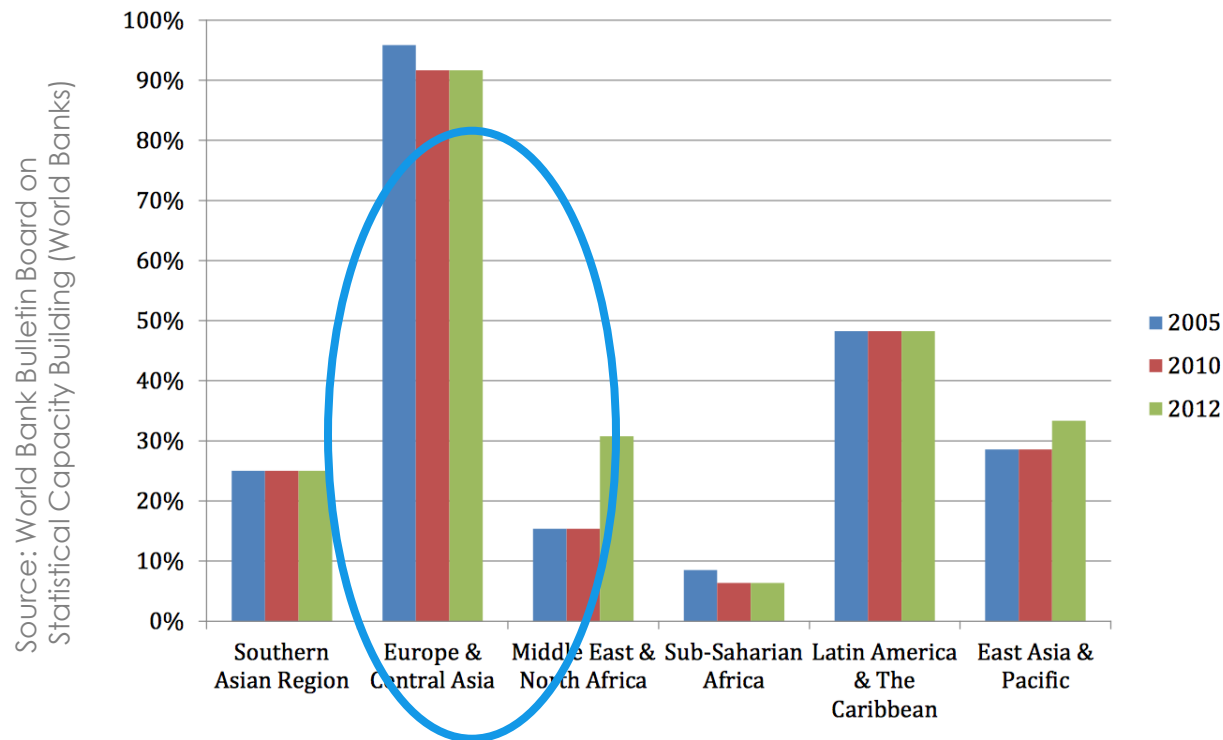
Censuses: A-



The State of Official statistics in LAC: B+?



**Other
aspects: B**



→ CEPAL project: hard to estimate maternal mortality because of lack of certification or registry in areas inhabited by indigenous populations or remote

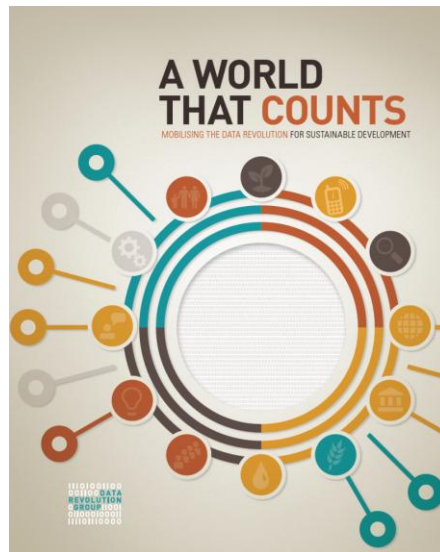
Main obstacles for LAC NSOs to leverage Big Data



1. **Institutional and cultural barriers to innovation and change**, e.g. lack of internal digital culture, scepticism vs. on Big Data, political will & coordination
2. **Constraints to data access and completeness**, particularly in access and continued use of data held by the private sector, limited ownership rights involving people and their relationships with data;
3. **Technical challenges**, including infrastructure for capturing, cleaning, processing, analyzing and visualizing structured and unstructured data
4. **Human capacity gaps**, including talent discovery, data literacy, limited data science training programs and involvement of academia
5. **Methodological challenges**, including challenges in data representativeness, biases, and the lack of standards and guidelines;
6. **Ethical and political risks**, including risks to privacy and weak legal frameworks

But huge window of opportunities and assets too

Post-2015 development / SDG / Big and Open Data agendas place NSOs at their core



CYBERSECURITY Saving Big Data from Itself

A three-step plan for using data right in an age of government overreach

By Alex "Sandy" Pentland

For the first few decades of its existence, the National Security Agency was a quiet department with one primary job: keeping an eye on the Soviet Union. Its enemy was well defined and monolithic. Its principal tools were phone taps, spy planes and hidden microphones. After the attacks of September 11, all of that changed. The NSA's chief enemy became a diffuse network of individual terrorists. Anyone in the world could be a legitimate target for spying. The nature of spying itself changed

as new digital communication channels proliferated. The exponential growth of Internet-connected mobile devices was just beginning. The NSA's old tools apparently no longer seemed sufficient.

In response, the agency adopted a new strategy: collect everything. As former NSA director Keith Alexander once put it, when you are looking for a needle in a haystack, you need the whole haystack. The NSA began collecting bulk phone call rec-



GLOBAL PARTNERSHIP
FOR SUSTAINABLE DEVELOPMENT DATA



Latin American NSOs are active!

Overview of selected Big Data projects involving regional NSOs

Table 2: Overview of Big Data projects in selected LAC NSOs

Type of Big Data	Data Used in Current NSO Activities	Projects	Other Organizations Involved	Project Status
Argentina (INDEC)				
<i>Exhaust data</i>	Web scraping	IPC Online		Planned
Brazil (IBGE)				
<i>Digital content</i>	Google Maps CDRs	Developing Water Accounts Tourism Monitoring	National Water Agency IBGE	Implemented /On-Gong Planned
Colombia (DANE)				
<i>Exhaust data</i>	Web scraping	IPC Online SIPSA		Planned Implemented/On-gong
<i>Digital content</i>	CDRs	Monitoring crime activities Socio-economic levels and networks	World Bank Data-Pop Alliance TransMilenio	Pilot stage Pilot stage Pilot stage
<i>Sensing data</i>	Satellite	Complementing the National Agriculture Census		Completed
Ecuador (INEC)				
<i>Exhaust data</i>	Web scraping	IPC Online		Pilot stage
<i>Digital content</i>	Twitter CDRs	Measuring Subjective Wellbeing Daytime Migration		Pilot stage Planned
Guatemala (INE)				
<i>Digital content</i>	CDRs	Monitoring Poverty Levels	World Bank Telefónica	Pilot stage
Mexico (INEGI)				
<i>Digital content</i>	Twitter	Subjective Wellbeing Subjective Wellbeing of Women Tourism Monitoring Movements Across Borders	InfoTec and Tec Monterrey Data2x, the University of Pennsylvania Ministry of Tourism	Completed Pilot Completed Planned



Box 4: Twitter for Tourism Monitoring in Mexico

In 2014, a working group on Big Data at INEGI conducted a pilot study to track domestic tourism from Twitter data, in order to contribute to the empirical modelling of individual tourist behavior. The objective of this pilot program was to identify the characteristics of an average Tweeting tourist in order to identify how many people travelled to Puebla and Guanajuato during the holiday weekend of February 1-3, 2014. The team of researchers from INEGI, in collaboration with the Mexican Ministry of Tourism, analysed 60 million Tweets from January to July 2014, from the continuous 1% georeferenced sample that Twitter makes available for free.⁶⁵ From this data, INEGI collected Tweets from the 7,955 Twitter users who Tweeted in Guanajuato (48%) and Puebla (52%) during the holiday. They then gathered all the Tweets sent by those users in the remainder of the target period (amounting to 827,424 total Tweets), and identified which users Tweeted from another state (presumably their homestate) after being in Guanajuato or Puebla, in order to map the origin of domestic tourism to those two areas during the holiday.⁶⁶ The resulting estimates of domestic tourism to Guanajuato and Puebla were compared to estimates made by the respective offices of tourism of those two states.⁶⁷

Overview of selected Big Data projects involving regional NSOs

Box 5: Maternal Morbidity and Remote Sensing of Malaria in Brazil

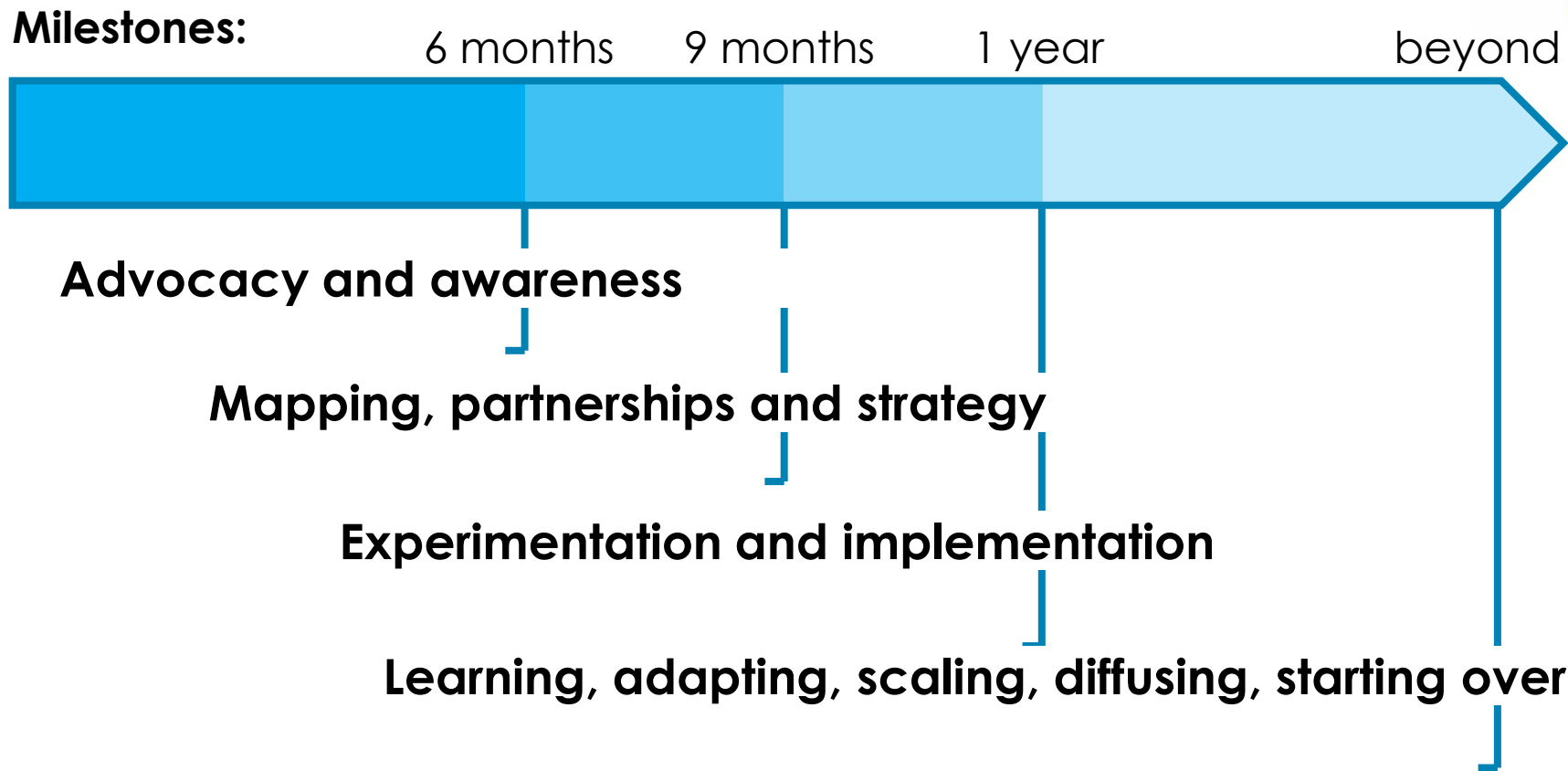
Remote sensing satellite data on vegetation density, soil moisture, population density, and spatial pattern of human infrastructure have long been used to predict levels of malaria risk. Advances in computing now allow more powerful use of these big datasets, including analysis of extreme spatial and temporal heterogeneity and inclusion of greater numbers of explanatory variables. This project seeks to create malaria risk maps for the Amazon Basin, focusing first on urban and peri-urban zones along the Brazil-Guyana border, which are areas with highly variable vector habitats and elevated incidences of illness. At least two vector distribution-mapping studies in this region exist, but to our knowledge there is no high-resolution dynamic mapping of malaria risk. The first phase of the project will use remote sensing data and existing health records, in combination with information about the economic, cultural, and health system, to estimate a spatial regression model that predicts morbidity burden in pregnant women, using DALYs (Disability-Adjusted Life Year) as the principal metric. The second phase will then test the accuracy of this model using data collected in real-time. UN Women and IBGE and the leading institutions piloting this study, drawing support from partner institutions The Vargas Foundation and the Amazon Malaria Initiative.

LAC countries and their NSOs have many assets too



1. An **urban, relatively young, innovative and technology hungry population** that share 4 major languages (Sp., Port., Fr., Eng.)
2. A long experience of the **Open Data movement**
3. The presence of strong **region-wide civil-society groups, institutions, academic networks, private sector, and working groups**
4. The emergence of **several pilots and public-private partnerships on Big Data**
5. The active involvement of **governments and public actors** that devise new **strategies and policies** and participate in **global forums**

Strategies must be *multi-partner, multi-sector, multi-year* and *regional*





The paper was part of multi-partner research, training and strategic assistance program funded by the World Bank

1. **Technical assistance:** strategy paper + strategic support
2. **Research:** 2 pilots on crime and poverty
3. **Regional training program** seed-funded by the Hewlett Foundation with other partners (Paris21, CEPAL..UNFPA...)



WORLD BANK



Telefonica



Universidad
de los Andes



DANE
Para tomar decisiones

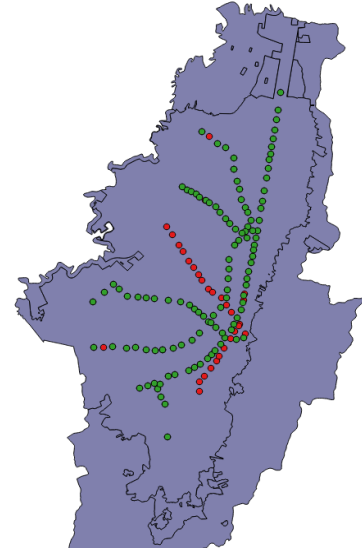


Piloto 2: predicting crime trends and patterns using aggregate data on people's dynamics



Research questions:

- Which areas are more likely to experience crime based on their digital traces?
- Can a public transportation system act as a public safety system?
- How do social networks in CDRs and crime interact?



Regional professional training program in Big Data en development 2015-2017



THE WILLIAM AND FLORA
HEWLETT
FOUNDATION

\$500k



**TODOS POR UN
NUEVO PAÍS**
PAZ EQUIDAD EDUCACIÓN



MIT Connection Science
the technology of innovation



Universidad
de los Andes

Knowledge platform with
literature reviews, articles,
toolkits, pilots, codes...

Training workshops and regional seasonal 'Big Data and development' schools in various countries (Bogotá, Dakar, Kigali, MIT Media Lab, Singapore,...) covering technical, institutional, legal aspects....

Regional Data Spaces

Brokering and fostering regional synergies to foster a people-centered Big Data revolution



IGARAPÉ INSTITUTE
think connect transform



MIT Connection Science
the technology of innovation



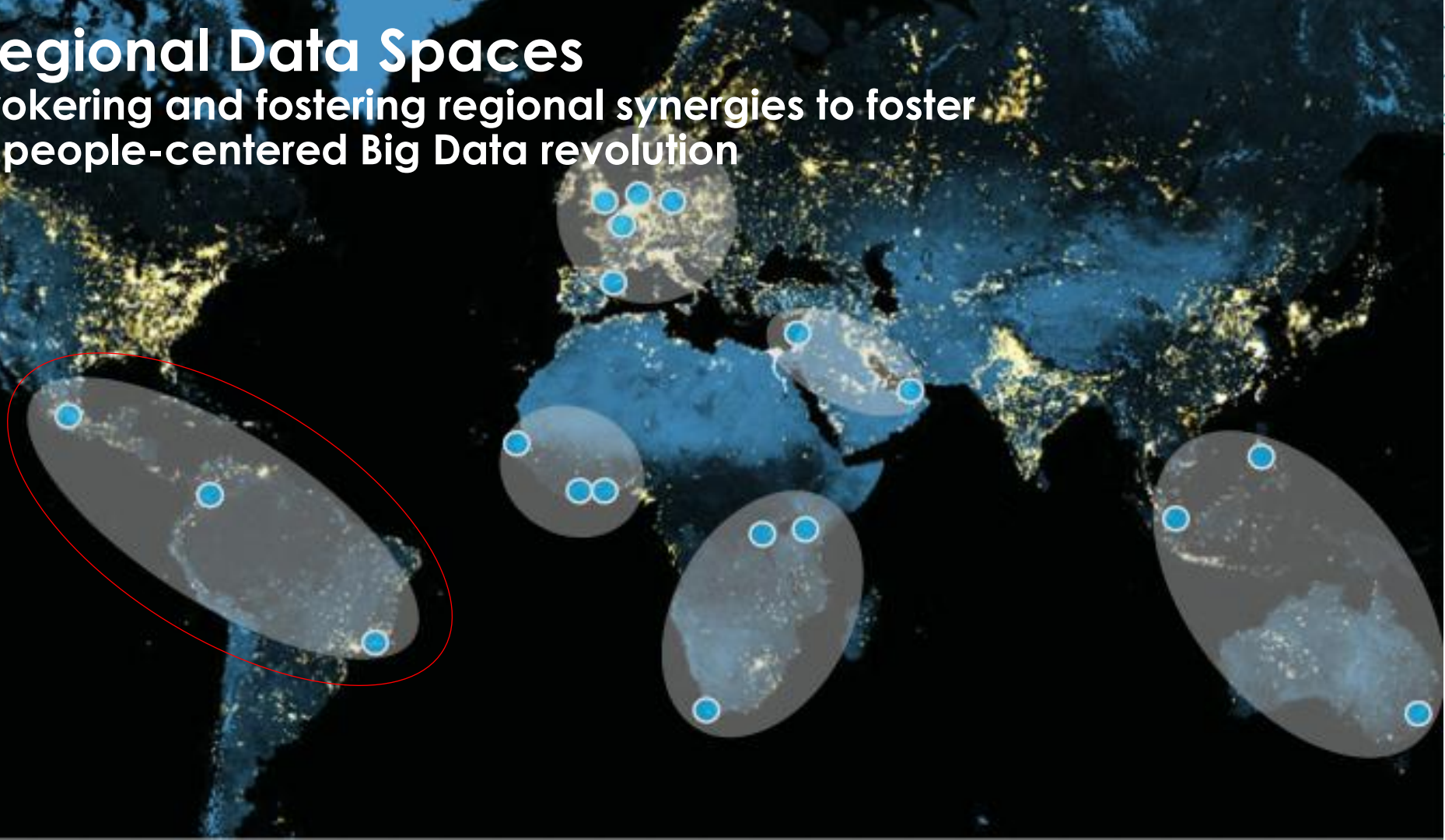
**SUSTAINABLE DEVELOPMENT
SOLUTIONS NETWORK**
A GLOBAL INITIATIVE FOR THE UNITED NATIONS



OXFAM
México

Regional Data Spaces

Brokering and fostering regional synergies to foster
a people-centered Big Data revolution





Gracias!

www.datapopalliance.org

eletouze@datapopalliance.org