

Big Data en el INEGI

UNA VISIÓN TECNOLÓGICA

11 de Noviembre 2015



10 -12 DE NOVIEMBRE, 2015



4 mil Empleados en Aguascalientes y 18 mil en todo el País

COLABORACIÓN INTERINSTITUCIONAL

- Nacional



Innovación con propósito de vida.



- Internacional



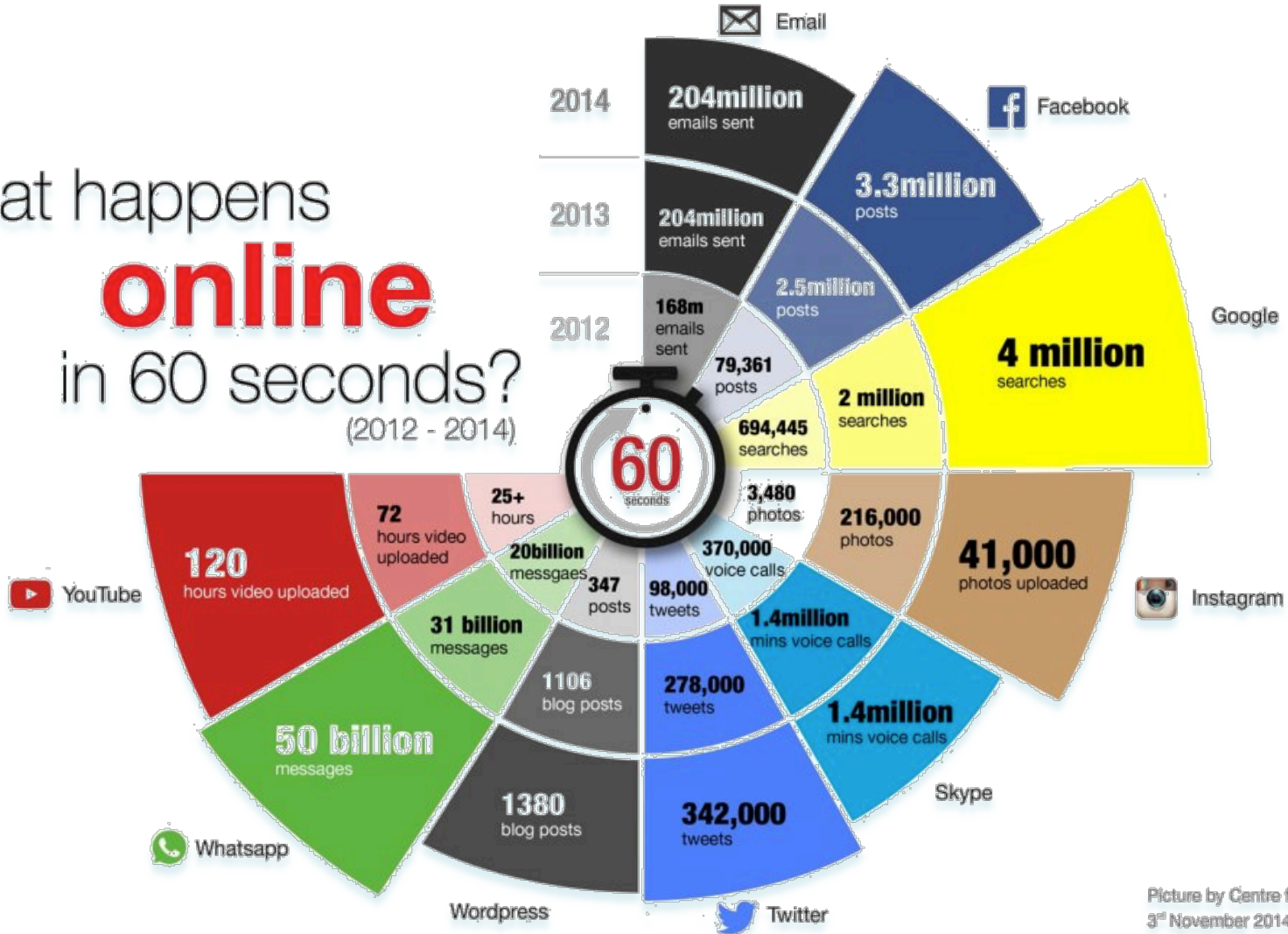
Considerar las Nuevas Fuentes de Datos para Complementar a las Tradicionales.

Fuentes sin un diseño original, en una diversidad de fines posibles, normalmente ajenos a las causas que permitieron la generación



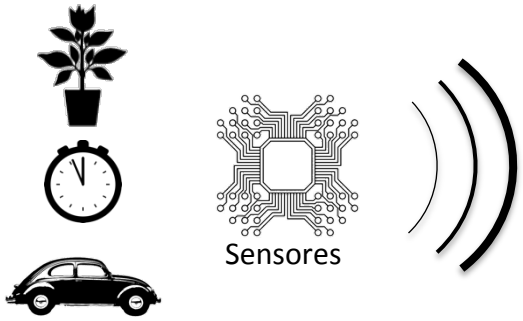
Las Fuentes de Información siguen creciendo

What happens
online
 in 60 seconds?
 (2012 - 2014)



Panorama Tecnológico

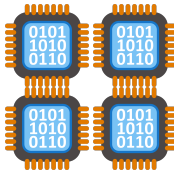
Internet de las Cosas



Infraestructura de Cómputo

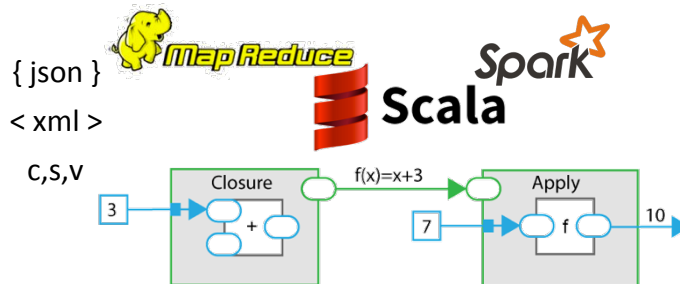


{ json }
< xml >
C,S,V



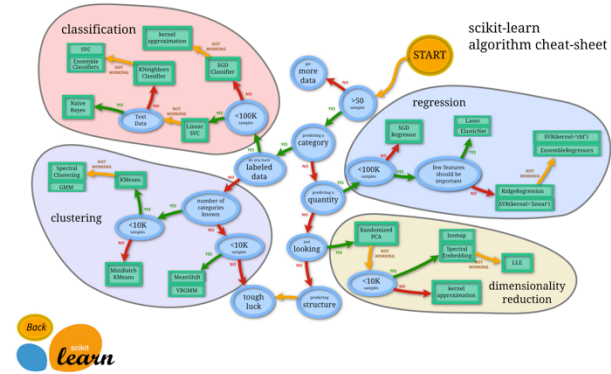
Computo Paralelo y
Concurrente

Programación Funcional



Razonamiento Algebraico

Estadística



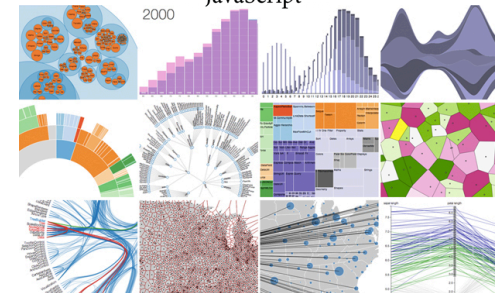
Análisis Multivariado
Machine Learning
Análisis de Interacción Espacial



Visualización



JavaScript



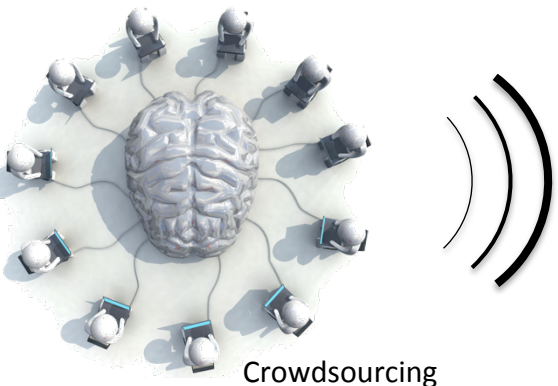
Data-Driven Documents

Internet de las Personas



{ json }
< xml >
C,S,V

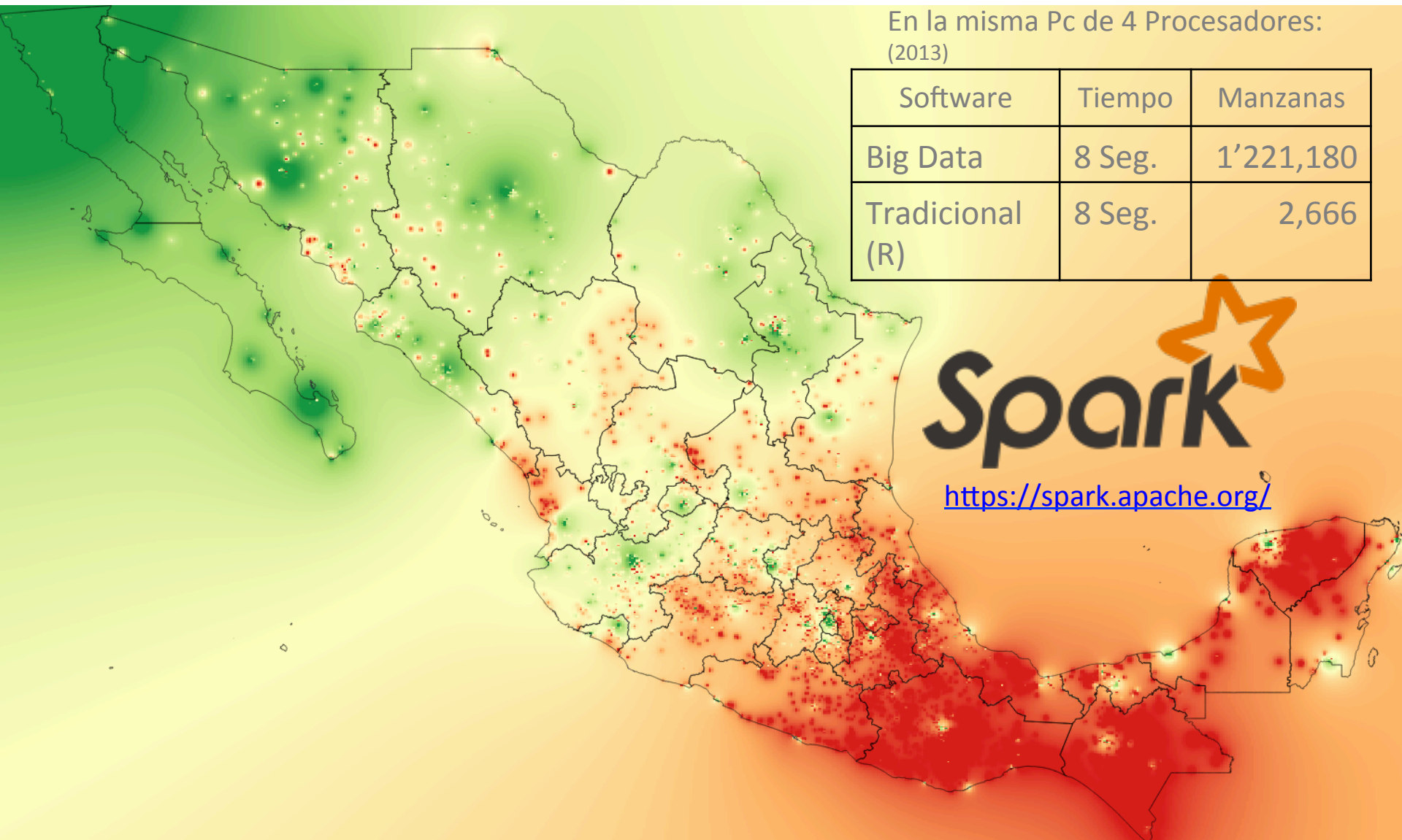
Internet de las Ideas



Estratificación de 1.2 M de Manzanas (2013)

En la misma Pc de 4 Procesadores:
(2013)

Software	Tiempo	Manzanas
Big Data	8 Seg.	1'221,180
Tradicional (R)	8 Seg.	2,666



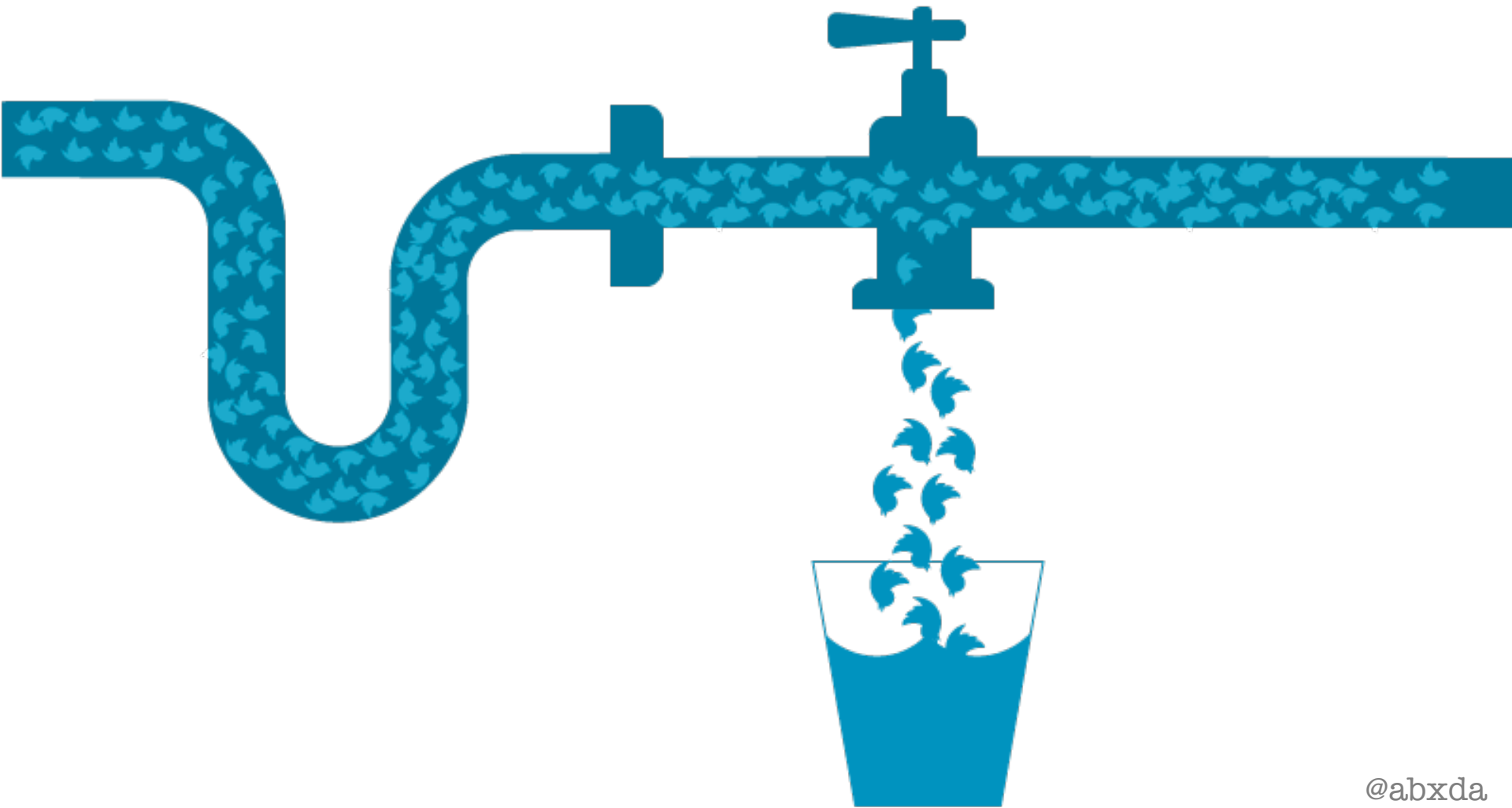
Spark

<https://spark.apache.org/>

TWITTER COMO FUENTE DE BIG DATA

Para medir el pulso emotivo de México

...y mucho más ...



OBJETIVO DEL PROYECTO

Generar indicadores experimentales, nuevos o que complementen los generados por métodos tradicionales, utilizando técnicas de Big Data para la extracción, almacenamiento, procesamiento, análisis y visualización de los datos.



Colaboración



INSTITUTO NACIONAL
DE ESTADÍSTICA Y GEOGRAFÍA

Dr. Enrique Ordaz

enrique.ordaz@inegi.org.mx

Dr. Gerardo Leyva

gerardo.leyva@inegi.org.mx

Dr. Alfredo Bustos

alfredo.bustos@inegi.org.mx

Dr. Juan Muñoz López

Juan.munoz@inegi.org.mx

Ing. Silvia Fraustro

Silvia.fraustro@inegi.org.mx

Mtro. Abel Coronado

abel.coronado@inegi.org.mx

Ing. Ricardo Olvera

Ricardo.olvera@inegi.org.mx

Lic. Marco Ibarra

Marco.ibarra@inegi.org.mx



Dr. Elio Villaseñor

elio.villaseñor@infotec.com.mx

Dr. Mario Graff

mario.graff@infotec.com.mx

Dr. Eric Tellez

eric.tellez@infotec.com.mx

Dr. Sabino Miranda

sabino.miranda@infotec.com.mx



Dr. Oscar S. Siordia

osanchez@centrogeo.edu.mx

Dra. Daniela Moctezuma

dmoctezuma@centrogeo.edu.mx

Y el apoyo de:



Innovación con propósito de vida.

Todos los tuits están disponibles para su recolección en tiempo real.

 <https://dev.twitter.com/docs/api/streaming>



Developers

API Health

Blog

Discussions

Documentation

Search

[Home](#) → [Documentation](#)

The Streaming APIs

View

[What links here](#)

Updated on Mon, 2012-09-24 14:47

API version 1

API version 1.1

Overview

The set of streaming APIs offered by Twitter give developers low latency access to Twitter's global stream of Tweet data. A proper implementation of a streaming client will be pushed messages indicating Tweets and other events have occurred, without any of the overhead associated with polling a REST endpoint.

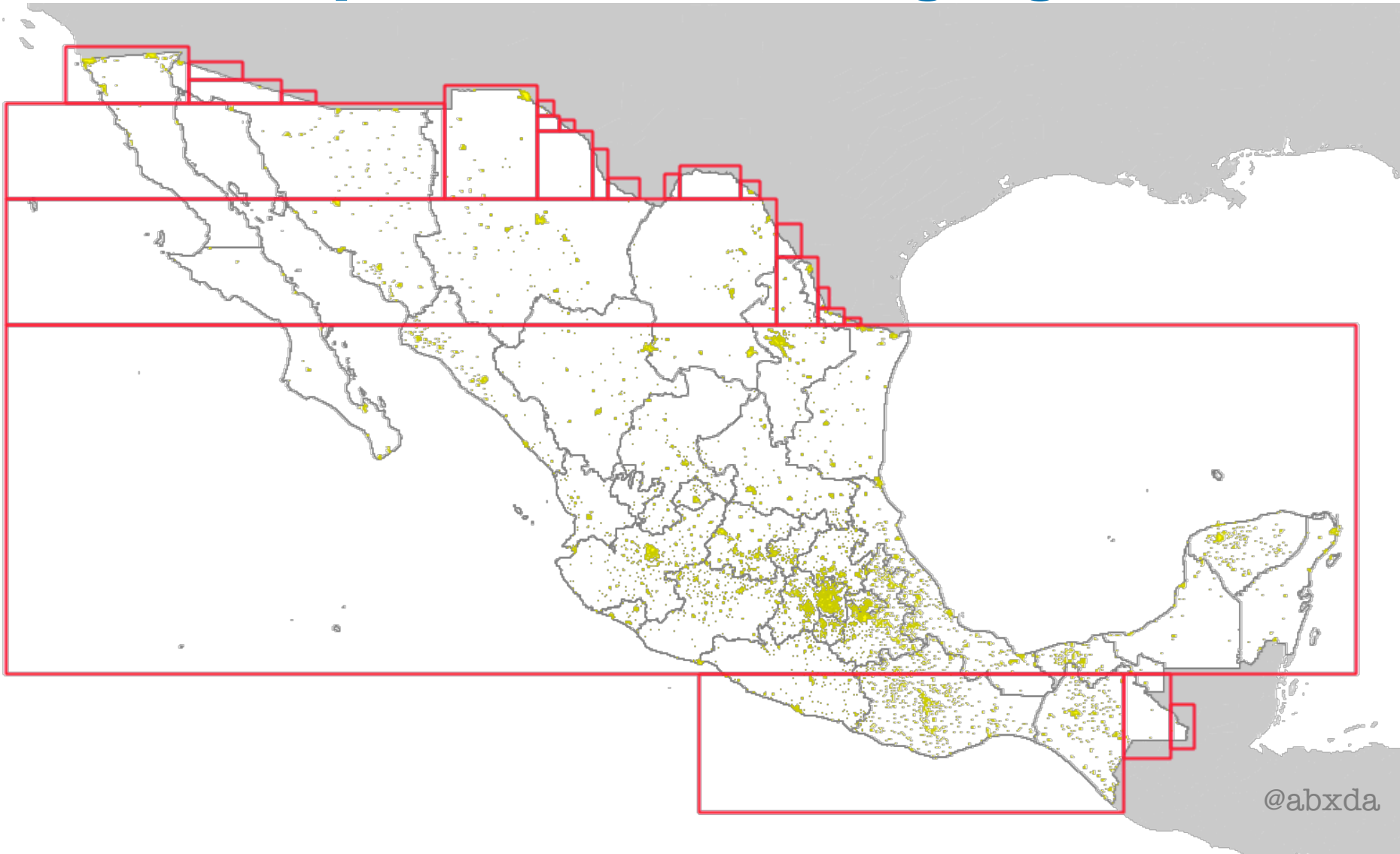
Twitter offers several streaming endpoints, each customized to certain use cases.

Public streams

Streams of the public data flowing through Twitter. Suitable for following specific users or topics, and data mining.

@abxda

Incluso permite consultas geográficas



Hydra

1 año 8 meses



< ESCALABILIDAD HORIZONTAL >



elasticsearch

<https://www.elastic.co>

@abxda

<http://cienciadedatos.inegi.org.mx/pioanalysis>



0 de 20 - nivel 0

Y ahí... entre todos tus gustos raros estaba yo.

¿El tema del tuit  es?

- Política
- Cultural / Entretenimiento
- Deporte
- Escolar / Laboral
- Personal
- Ni idea

¿El tuitero se sentía?



@hbcolectivo

@ricardoalvera

@abxda

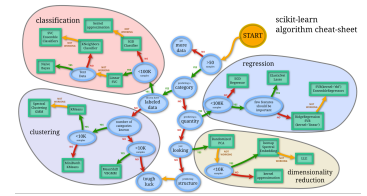
Entrenamiento



$f(x)$

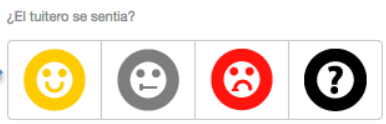


Representación numérica

$$\begin{pmatrix} 0 & 3 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 7 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 & 0 & 0 & 0 & 7 \\ 0 & 0 & 0 & 5 & 0 & 0 & 0 & 8 & 0 & 0 \\ 2 & 0 & 0 & 0 & 0 & 6 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 9 & 0 & 0 & 0 & 4 & 0 \\ 0 & 1 & 0 & 7 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 4 \\ 0 & 0 & 7 & 0 & 0 & 0 & 0 & 5 & 0 & 0 \end{pmatrix}$$


<http://scikit-learn.org/>
<http://www.r-project.org/>

Machine Learning

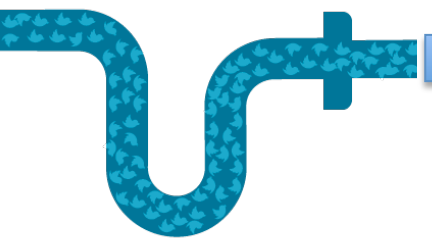


Etiquetado Manual



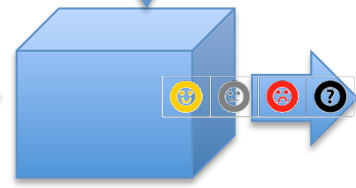
Muestra de Tuits

Producción

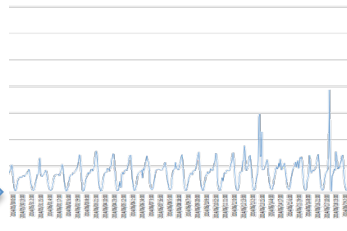


$f(x)$



$$\begin{pmatrix} 0 & 3 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 7 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 & 0 & 0 & 0 & 7 \\ 0 & 0 & 0 & 5 & 0 & 0 & 0 & 8 & 0 & 0 \\ 2 & 0 & 0 & 0 & 0 & 6 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 9 & 0 & 0 & 0 & 4 & 0 \\ 0 & 1 & 0 & 7 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 4 \\ 0 & 0 & 7 & 0 & 0 & 0 & 0 & 5 & 0 & 0 \end{pmatrix}$$


Clasificador



Indicador de sentimiento

Tuits en Tiempo Real



¿Feliz o triste?, el INEGI analiza el ánimo de los tuiteros

noticieros.televisa.com/programas-primero-noticias/1510/primero-futuro-estado-animo-tuiteros-me

Top 10 Emojis on Twitter

Least Popular Emojis on Twitter

Telesión Deportes Noticieros Espectáculos

Noticieros Televisa

México DF Estados Mundo Secciones



Primero Futuro: Estado de ánimo de tuiteros de México

Temas Relacionados
Recomendados, Primero Noticias, Carlos Loreto de Mola

Aura López comenta que gracias a un estudio que hizo el INEGI, ya es posible saber el estado de ánimo de los usuarios de Twitter en México



INEGI usará Twitter para un ambicioso proyecto de estadística

Pretende dar seguimiento en tiempo real al estado de ánimo de los tuiteros en México

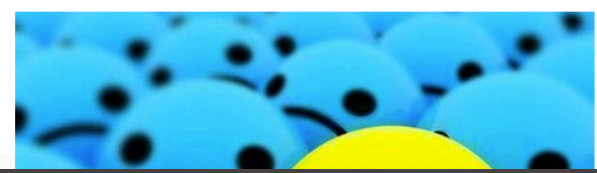
Twitter 8 G+ 0 Share Comentarios 0



¿Cuáles son los estados más felices, según sus tuits?

Por millones de tuits, el Inegi concluyó que Guerrero, Oaxaca y Michoacán son los estados más felices; Sonora y Coahuila son de los más tristes.

114 Me gusta 301 1 Compartir 416 G+ 0

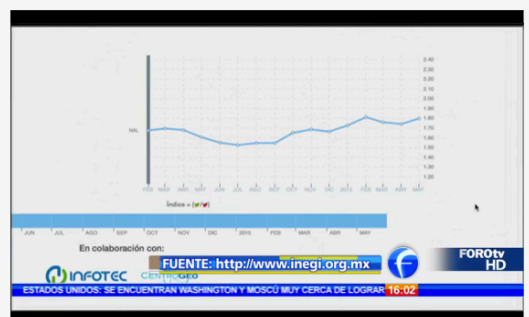
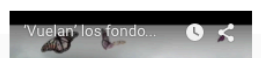


Indigo



MÁS SOBRE EL TEMA VIDEOS

os encontrados en



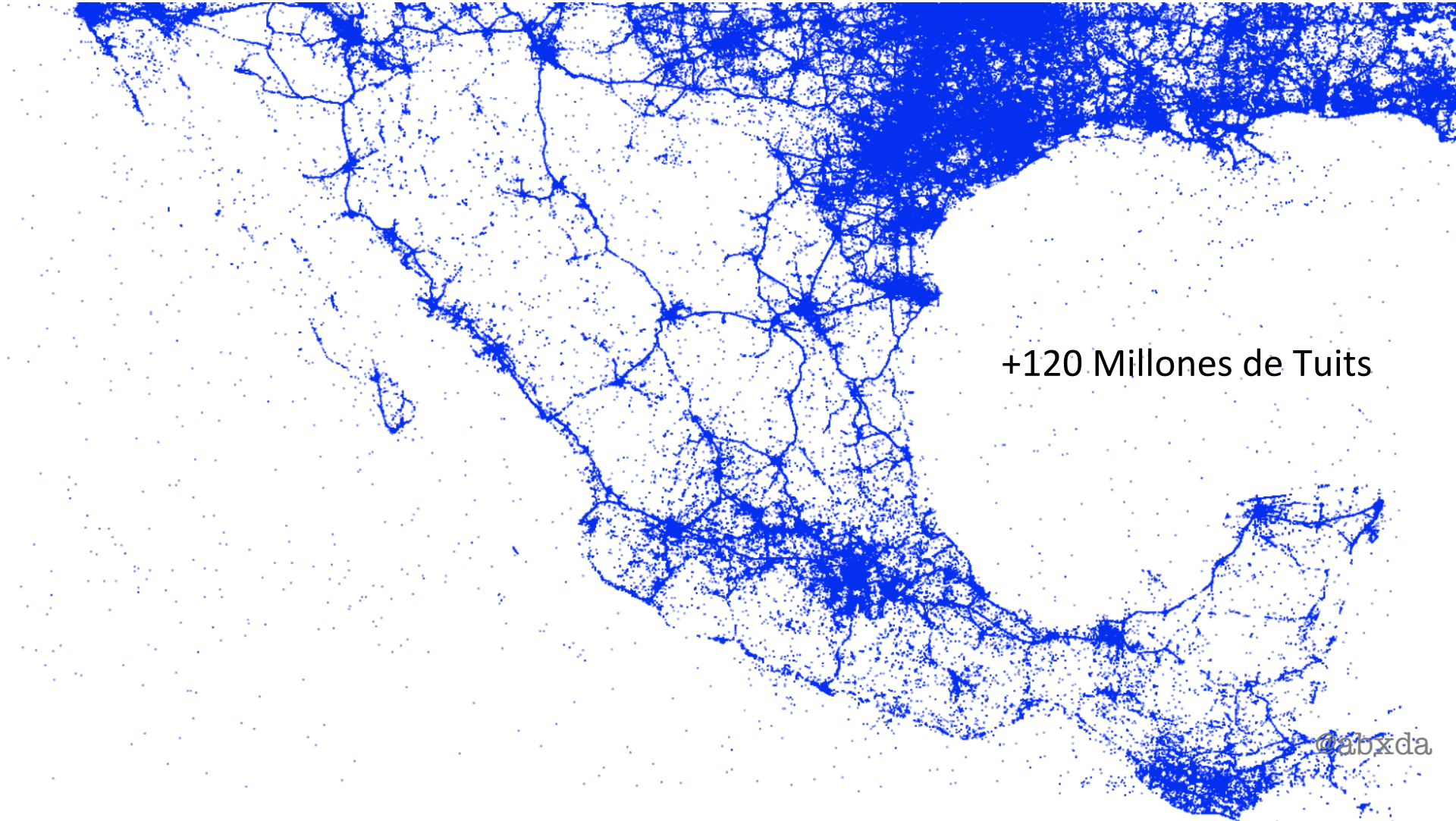
MAS ALLÁ DEL ANÁLISIS DEL SENTIMIENTO



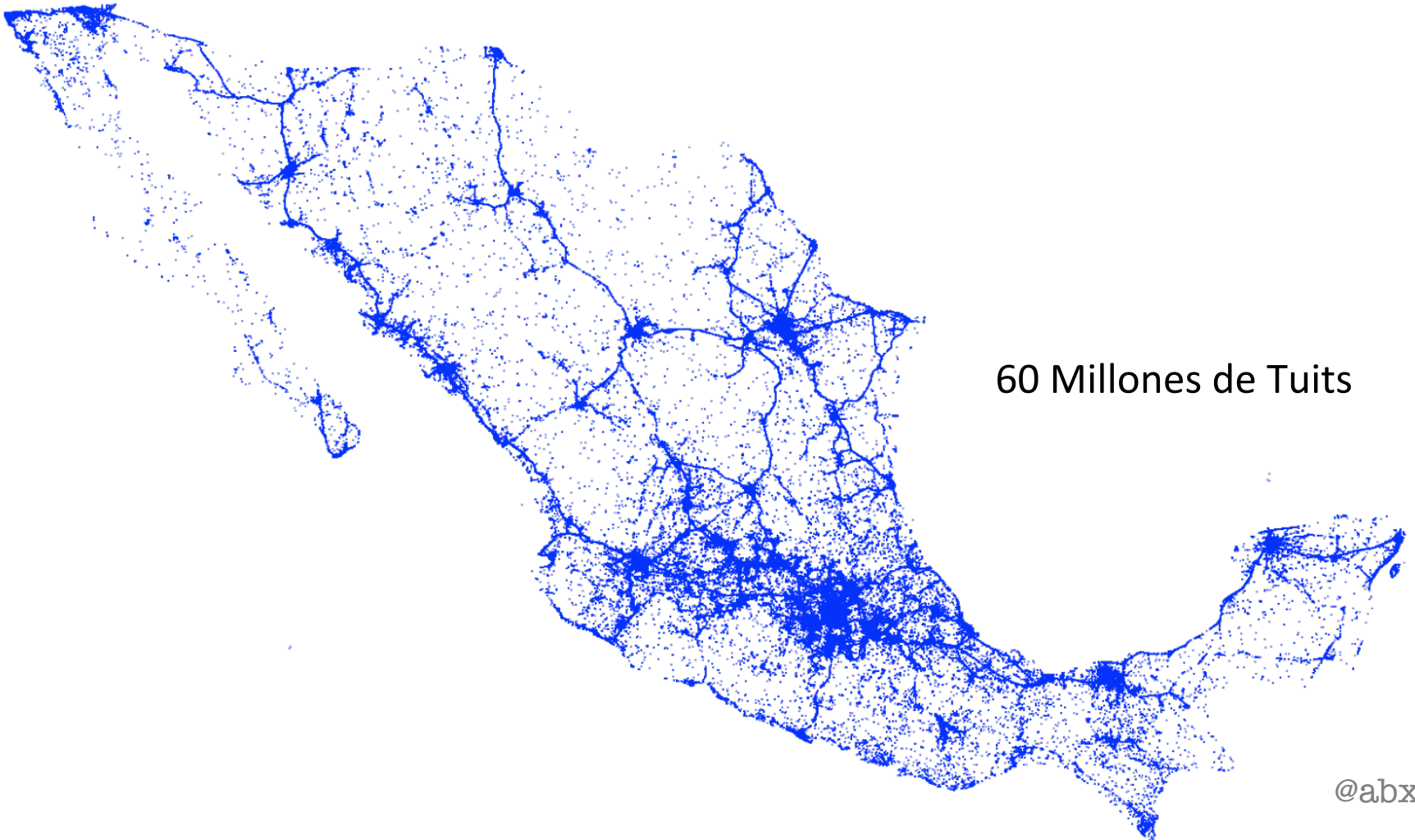
Apache Spark

<http://spark.apache.org/>

Visualización de la Base de Datos

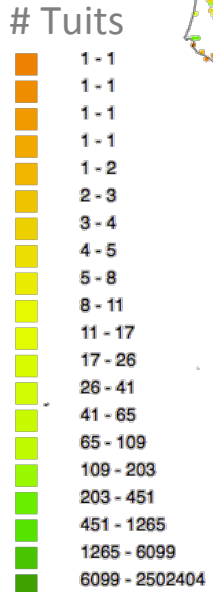
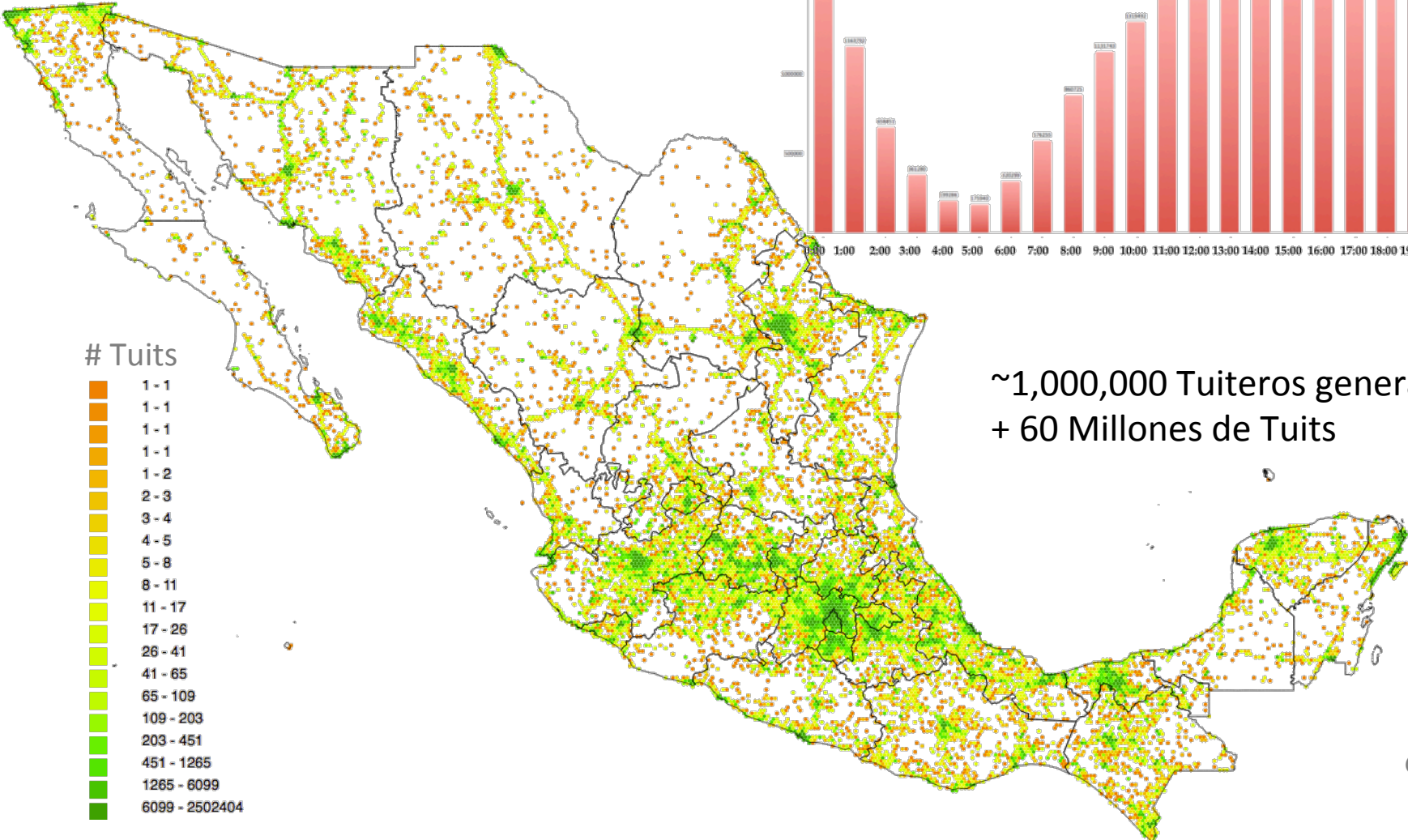


Visualización de la Base de Datos



Frecuencia de Tuiteo

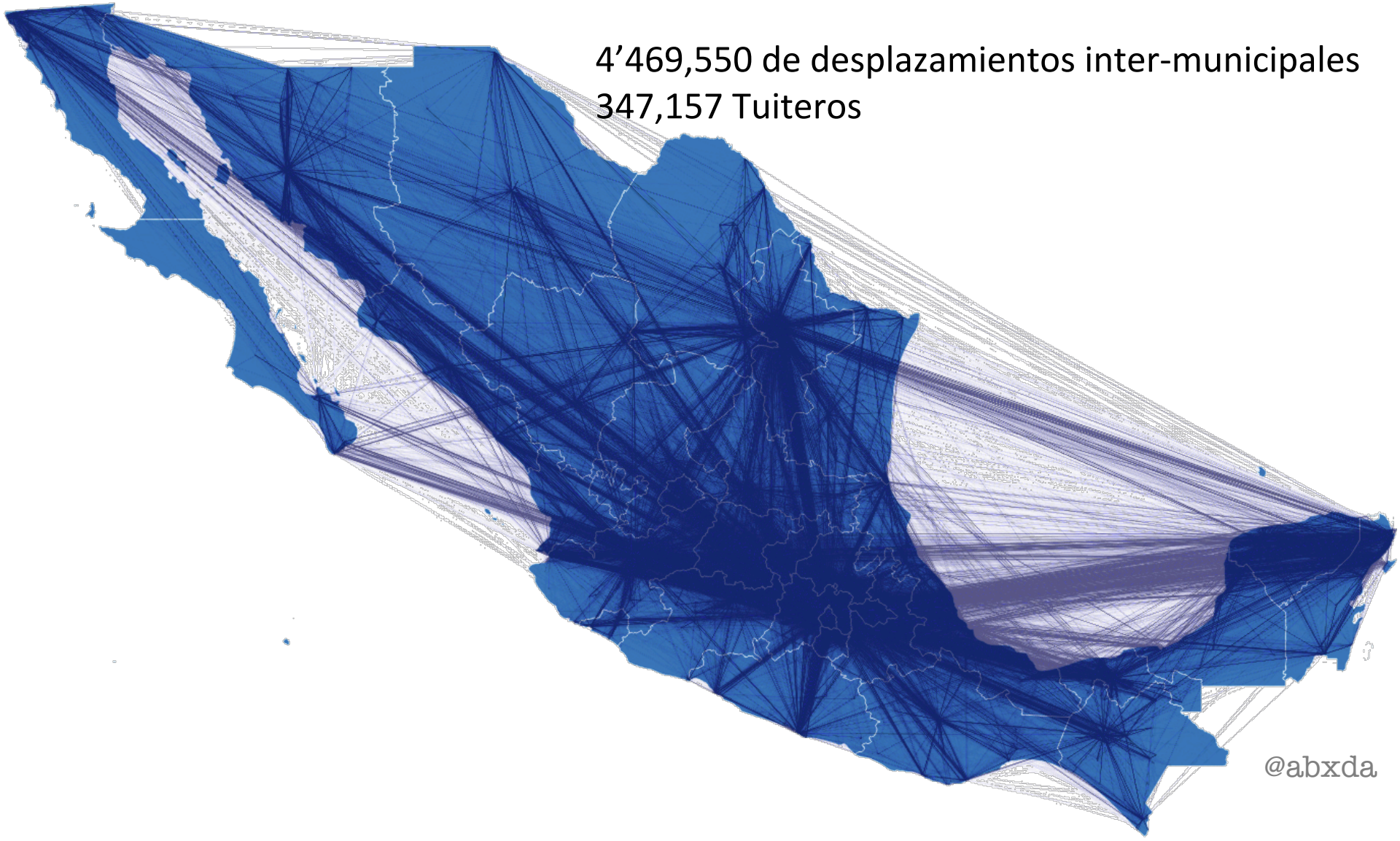
Frecuencia por hora del día



~1,000,000 Tuiteros generaron + 60 Millones de Tuits

Movilidad de los Tuiteros

4'469,550 de desplazamientos inter-municipales
347,157 Tuiteros



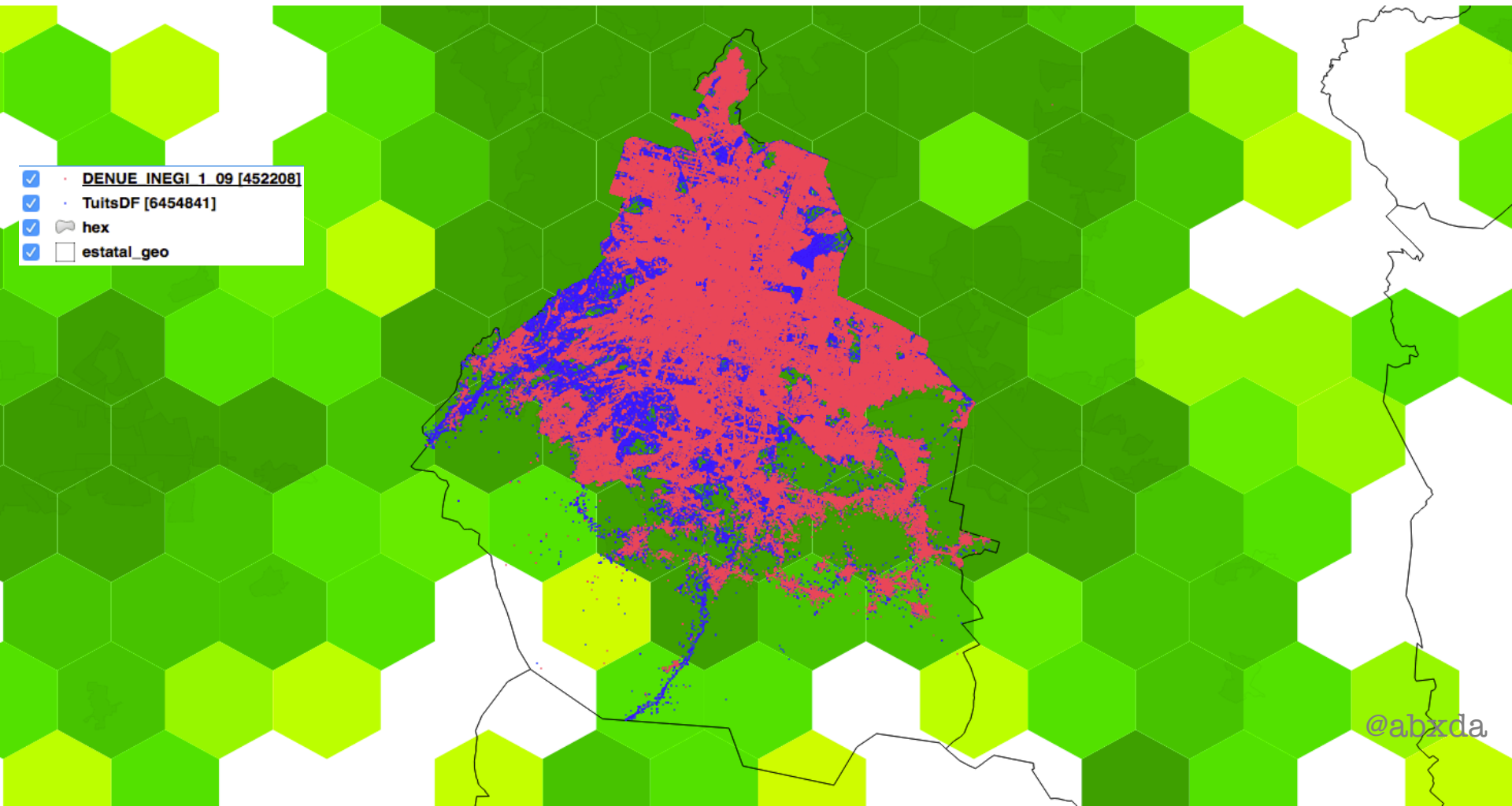
Red Nacional de Caminos (Open Data) y **INEGI** **Twitter**



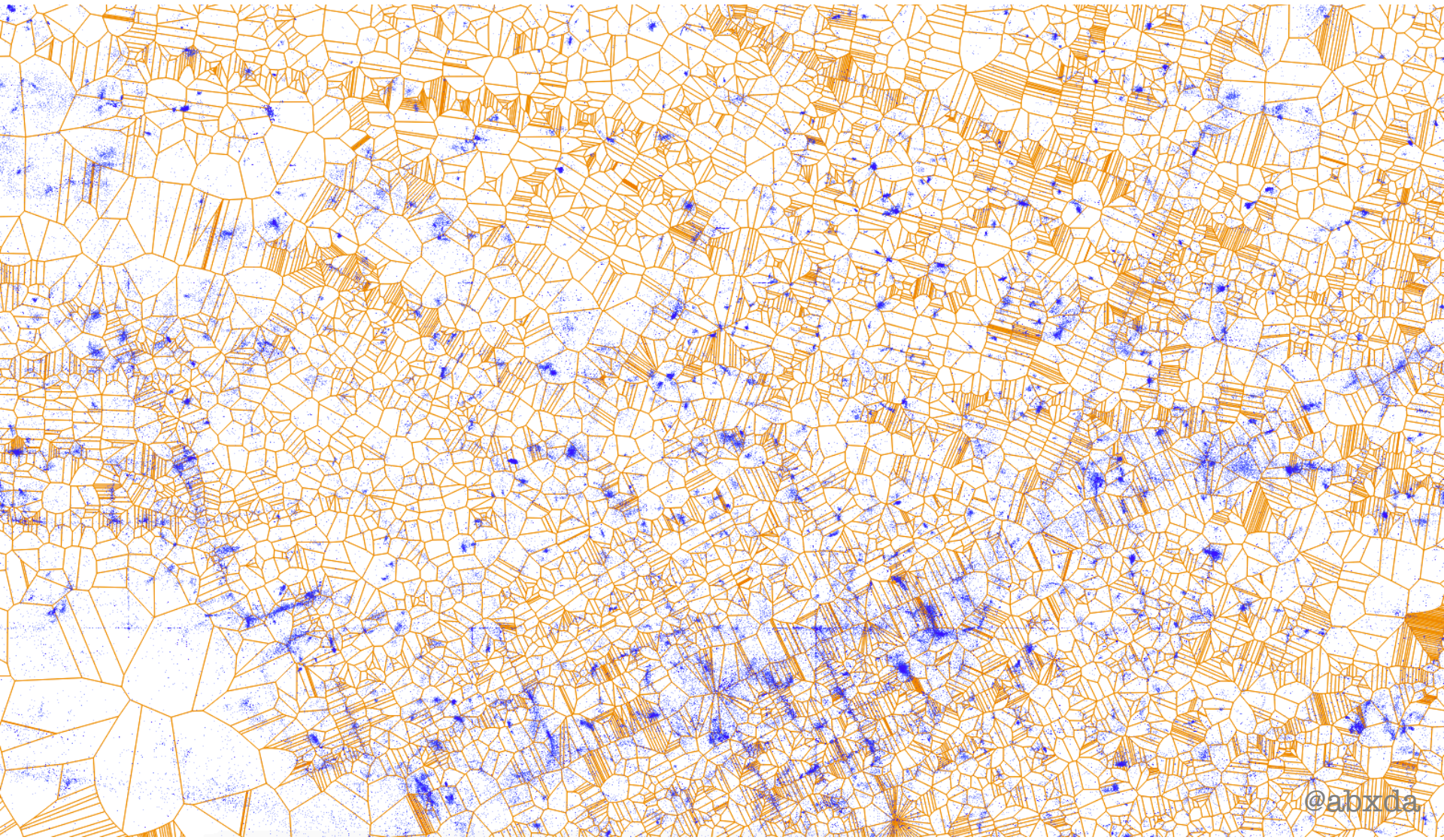
Red Nacional de Caminos y Twitter



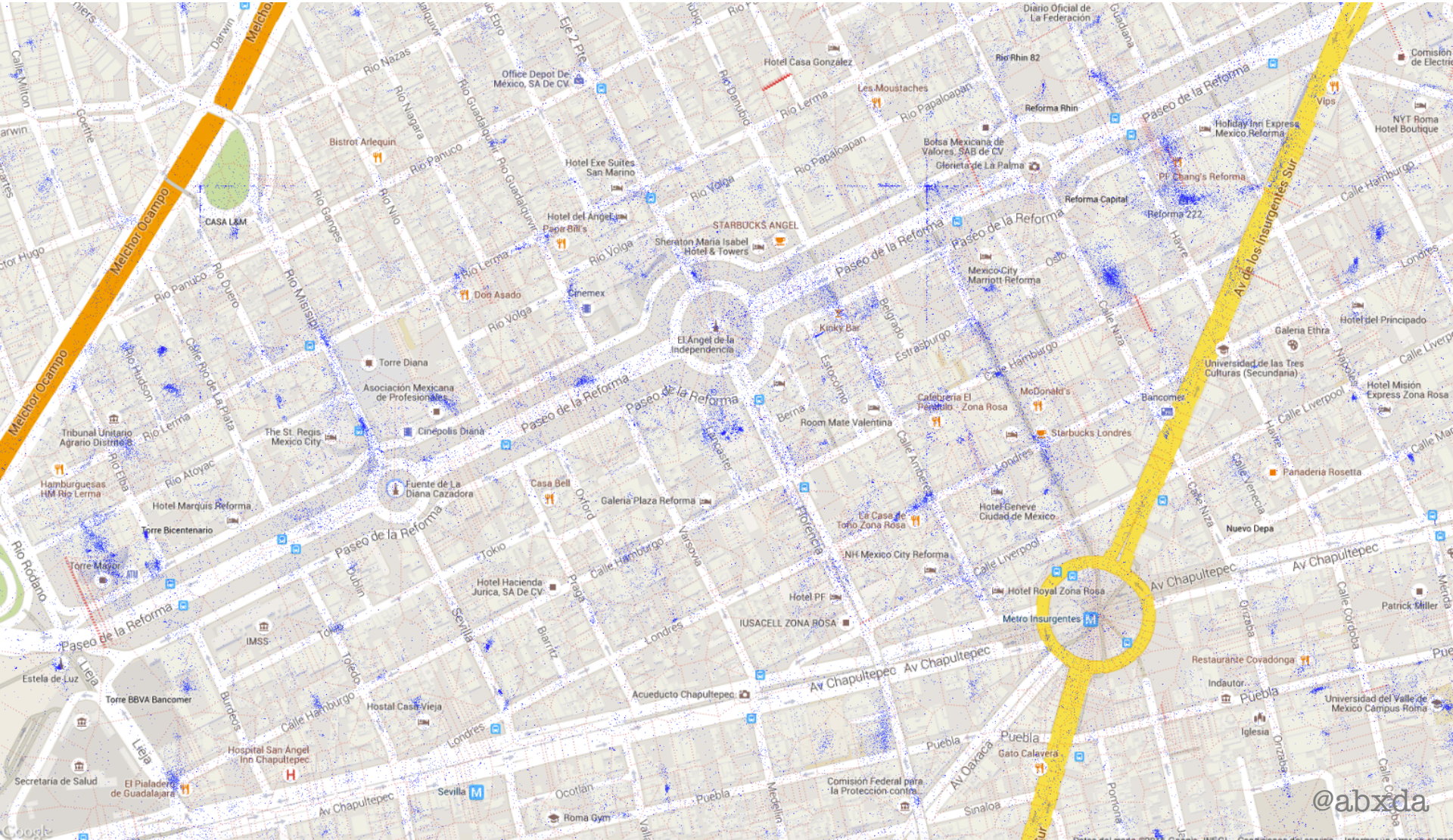
DENUE & Twitter



DENUE & Twitter

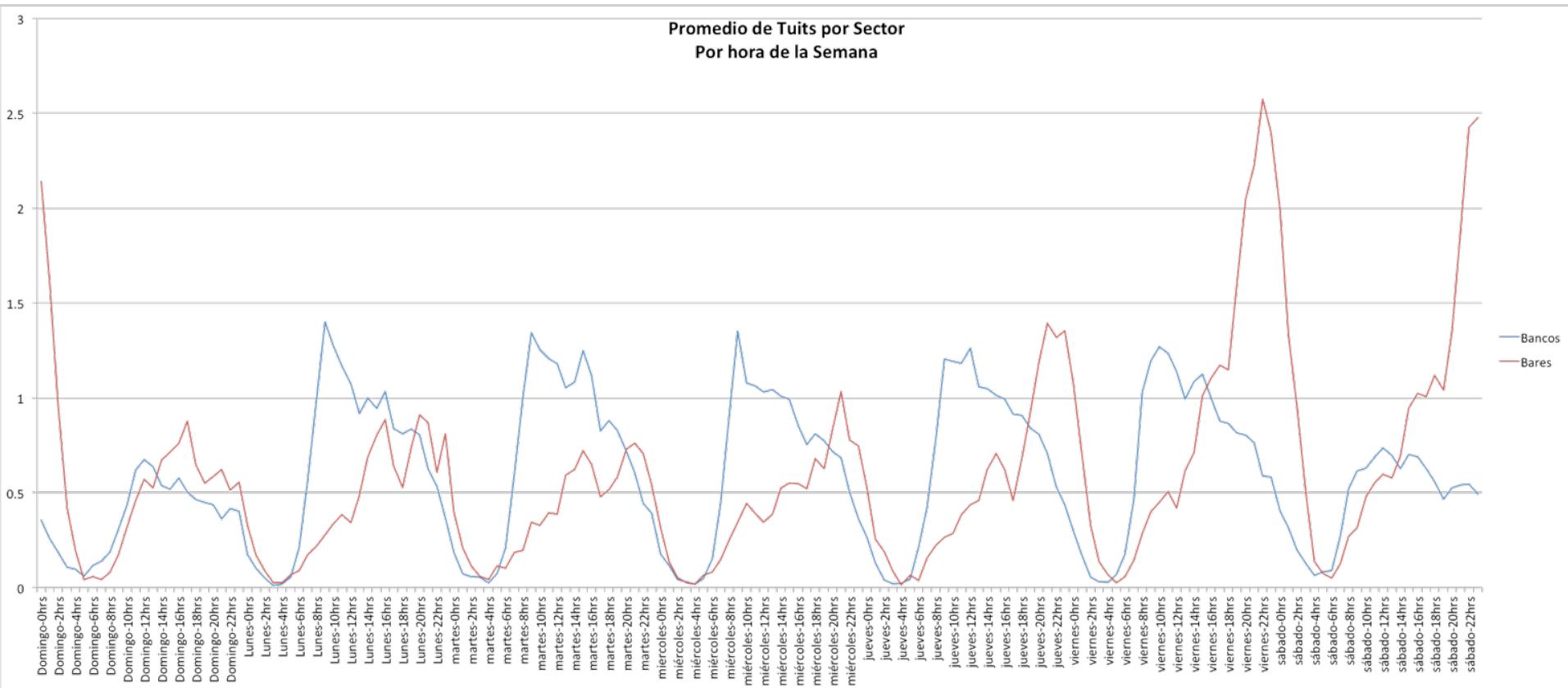


DENUE & Twitter



@abxda

Horarios de Tuiteo cerca de algún sector



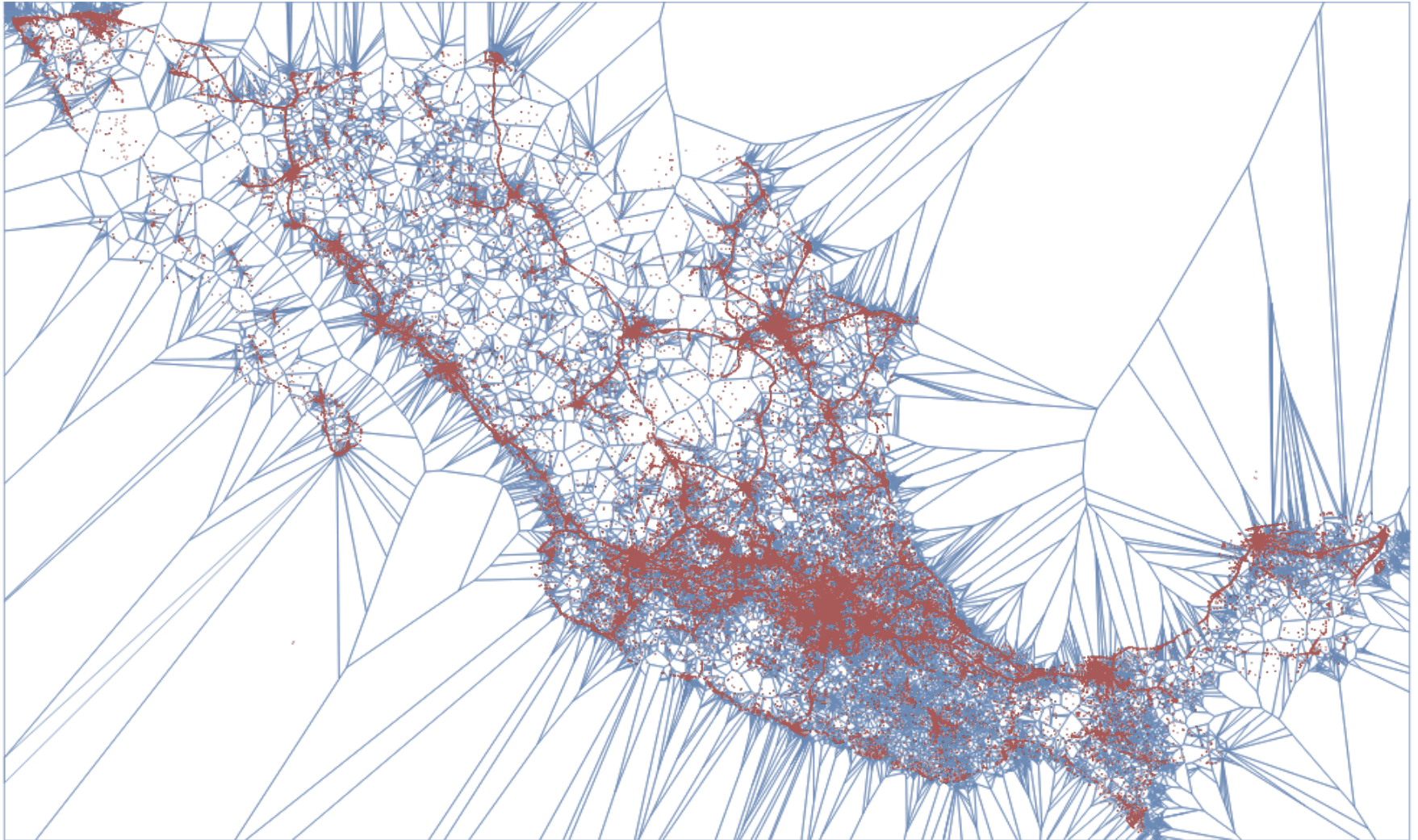
PRIMER EJERCICIO NACIONAL

DENUE - TWITTER

4.9 M de Polígonos de Voronoi (DENUE)



Big Spatial Join (4.9 M DENUE +60 M Tweets)



SpatialSpark

Large-Scale Spatial Join Query Processing in Cloud

Simin You

Dept. of Computer Science
CUNY Graduate Center
New York, NY, USA
syou@gc.cuny.edu

Jianting Zhang

Department of Computer Science
The City College of New York
New York, NY, USA
jzhang@cs.ccnyc.cuny.edu

Le Gruenwald

Dept. of Computer Science
The University of Oklahoma
Norman, OK, USA
ggruenwald@ou.edu

Abstract— The rapidly increasing amount of location data available in many applications has made it desirable to process their large-scale spatial queries in Cloud for performance and scalability. We report our designs and implementations of two prototype systems that are ready for Cloud deployments: SpatialSpark based on Apache Spark and ISP-MC based on Cloudera Impala. Both systems support indexed spatial joins based on point-in-polygon test and point-to-polyline distance computation. Experiments on the pickup locations of ~170 million taxi trips in New York City and ~10 million global species occurrences records have demonstrated both efficiency and scalability using Amazon EC2 clusters.

Existing works on processing large-scale spatial join query processing mainly fall into two categories: a) improving single-node efficiency by exploiting the massive data-parallel computing power of hardware accelerators, such as Graphics Processing Units (GPUs) that are capable of general computing, and b) exploiting scalability provided by Hadoop/MapReduce based systems. Here MapReduce is referred to both as a computing model and as a component in Hadoop (together with the Hadoop Distributed File System—HDFS). Our previous work on GPU-based spatial joins [2,3] have demonstrated that it is quite possible to achieve orders of magnitude of performance improvements by re-designing and

SpatialSpark: Open Source

← → ↻ GitHub, Inc. [US] <https://github.com/syoummer/SpatialSpark> B 🔧 🔍 ⭐


GitHub Explore Features Enterprise Pricing

 [syoummer](#) / **SpatialSpark** 👁 Watch 9

Big Spatial Data Processing using Spark <http://simin.me/projects/spatialspark/>

🕒 29 commits 🌿 1 branch 🏷 1 release 👤 3 contributors

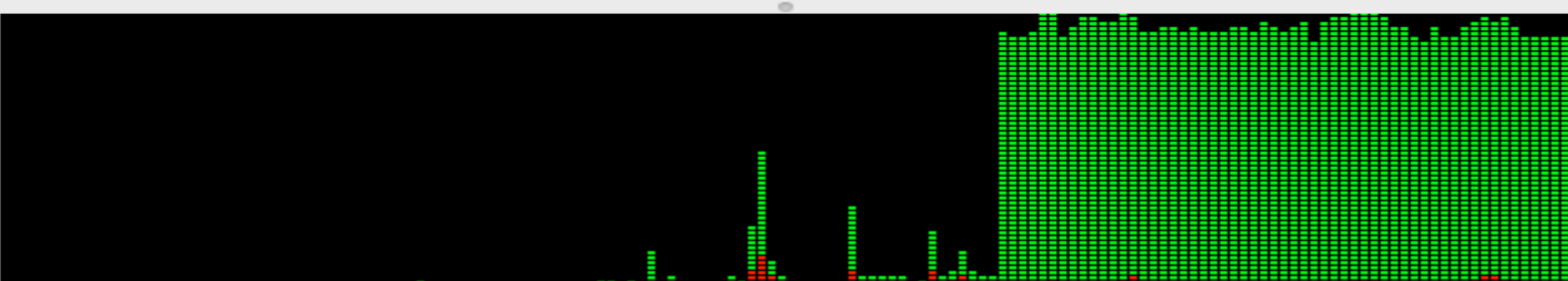
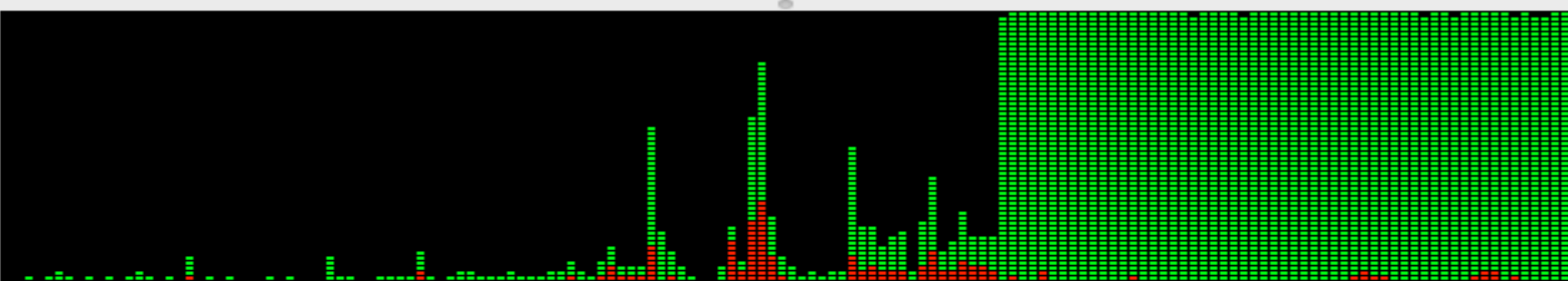
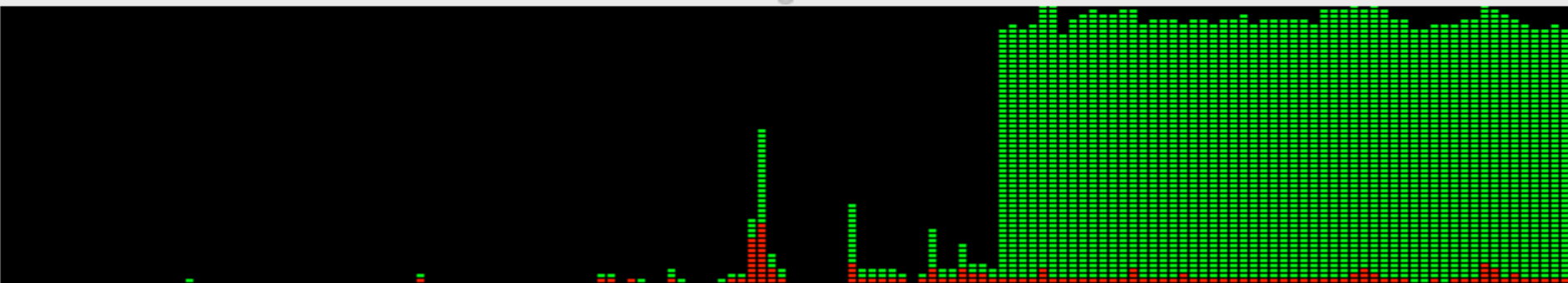
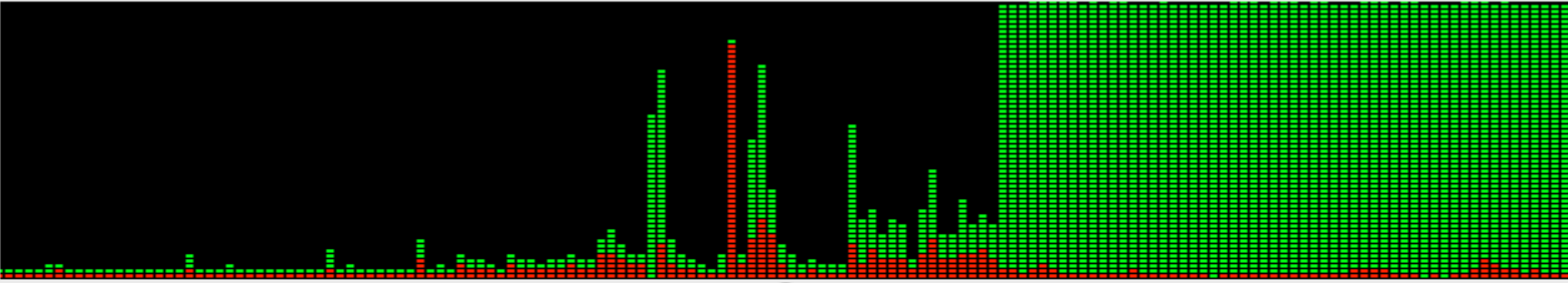
 Branch: **master** ▾ **SpatialSpark** / + ☰

 kgs Spark tests refactor	Latest commit ea573d1 on 8 Oct	
📁 data	init commit	8 months ago
📁 project	Publishing code using sonatype-sbt	a month ago
📁 src	Spark tests refactor	a month ago
📄 citation	Some code cleanup	2 months ago

Running Code into Local Apache Spark

```
[MacBook-Pro-de-Abel:~ abxda$ bin/spark-submit --driver-memory 1800M  
--executor-memory 1800M --jars /Users/abxda/Development/BigData/Spark/GeoSpark-0.1-GeoSpark.jar --class spatialspark.main.SpatialJoinApp  
/Users/abxda/Development/BigData/Spark/SpatialSpark-master/target/sca  
la-2.10/spatial-spark_2.10-1.1-SNAPSHOT.jar --left /Users/abxda/Devel  
opment/BigData/Spark/tuits_32_ent_mun.tsv --geom_left 2 --right /User  
s/abxda/Development/BigData/Spark/VORONOI_MINI_ZAC.csv --geom_right 0  
--output /Users/abxda/Development/BigData/Spark/salidaZacatecas.tsv ]  
--predicate within --parallel_part true --method stp --conf 32:32:0.1  
--partition 100
```

Historial de la CPU



DENUE - Twitter

Actividad Económica	Total de Tuits
Comercio al por menor en tiendas de abarrotes, ultramarinos y misceláneas	7,741,777
Salones y clínicas de belleza y peluquerías	2,378,443
Comercio al por menor de artículos de papelería	1,470,637
Restaurantes con servicio de preparación de alimentos a la carta o de comida corrida	1,445,646
Comercio al por menor en minisupers	1,246,794
Restaurantes con servicio de preparación de antojitos	1,088,333
Restaurantes con servicio de preparación de tacos y tortas	1,080,876
Reparación mecánica en general de automóviles y camiones	983,615
Asociaciones y organizaciones religiosas	967,601
Cafeterías, fuentes de sodas, neverías, refresquerías y similares	901,126
Comercio al por menor de ropa, excepto de bebé y lencería	691,736
Consultorios dentales del sector privado	664,654
Panificación tradicional	612,145
Lavanderías y tintorerías	592,689
Servicios de acceso a computadoras	564,462
Elaboración de tortillas de maíz y molienda de nixtamal	556,633
Fabricación de productos de herrería	515,880
Restaurantes con servicio de preparacióm de pizzas, hamburguesas, hot dogs y pollos r	509,583
Hoteles con otros servicios integrados	488,394
Restaurantes que preparan otro tipo de alimentos para llevar	476,261
Alquiler sin intermediación de salones para fiestas y convenciones	446,707
Banca múltiple	445,290
Comercio al por menor de cerveza	408,057
Comercio al por menor en ferreterías y tlapalerías	405,750
Escuelas de educación preescolar del sector privado	401,189
Administración pública en general	385,663
Centros de acondicionamiento físico del sector privado	366,161
Hojalatería y pintura de automóviles y camiones	357,171
Reparación y mantenimiento de otros artículos para el hogar y personales	353,976
Servicios de preparación de otros alimentos para consumo inmediato	351,555
Comercio al por menor de dulces y materias primas para repostería	338,091
Escuelas del sector privado que combinan diversos niveles de educación	329,189
Servicios de contabilidad y auditoría	327,615

