

Introduction and Overview

CS524: Big Data Visualization & Analytics

Fabio Miranda

<https://fmiranda.me>

A DAY IN DATA

The exponential growth of data is undisputed, but the numbers behind this explosion – fuelled by internet of things and the use of connected devices – are hard to comprehend, particularly when looked at in the context of one day

 **500m**

tweets are sent every day

Twitter

294bn

billion emails are sent

Radicati Group

320bn

emails to be sent each day by 2021

306bn

emails to be sent each day by 2020

 **3.9bn**
people use emails

Radicati Group

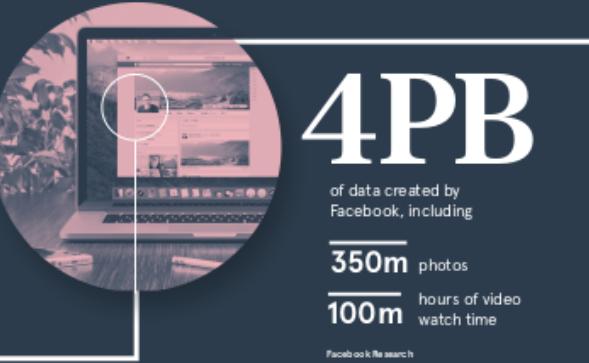
ACCUMULATED DIGITAL UNIVERSE OF DATA

4.4ZB

2010

44ZB

2020



DEMYSTIFYING DATA UNITS

From the more familiar ‘bit’ or ‘megabyte’, larger units of measurement are more frequently being used to explain the masses of data

| Unit | Value | Size |
|---------------------|--------------------|---|
| b bit | 0 or 1 | $\frac{1}{8}$ of a byte |
| B byte | 8 bits | 1 byte |
| KB kilobyte | 1,000 bytes | 1,000 bytes |
| MB megabyte | $1,000^3$ bytes | 1,000,000 bytes |
| GB gigabyte | $1,000^3$ bytes | 1,000,000,000 bytes |
| TB terabyte | $1,000^6$ bytes | 1,000,000,000,000 bytes |
| PB petabyte | $1,000^12$ bytes | 1,000,000,000,000,000,000 bytes |
| EB exabyte | $1,000^15$ bytes | 1,000,000,000,000,000,000,000 bytes |
| ZB zettabyte | $1,000^{18}$ bytes | 1,000,000,000,000,000,000,000,000 bytes |
| YB yottabyte | $1,000^{21}$ bytes | 1,000,000,000,000,000,000,000,000,000 bytes |

A lowercase ‘b’ is used as an abbreviation for bits, while an uppercase ‘B’ represents bytes.

65bn

messages sent over WhatsApp and two billion minutes of voice and video calls made

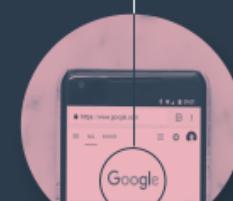
Facebook



Searches made a day → **5bn**

Searches made a day from Google → **3.5bn**

Smart Insights



463EB

of data will be created every day by 2025

idc

95m

photos and videos are shared on Instagram

Instagram Business

28PB

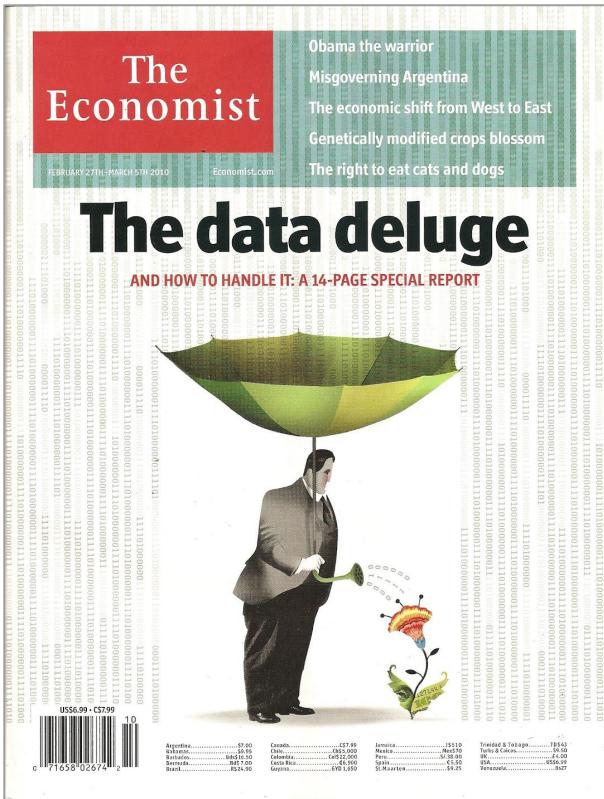
to be generated from wearable devices by 2020

Statista



Source: Raconteur

Data is everywhere



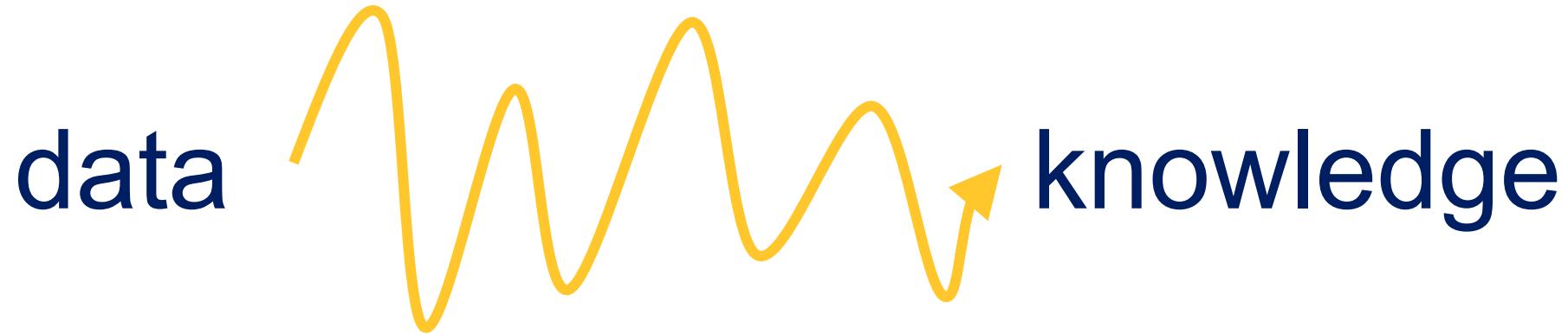
“The ability to take data—to be able to understand it, to process it, to extract value from it, to visualize it, to communicate it—that’s going to be a hugely important skill in the next decades, ... because now we really do have essentially free and ubiquitous data.”

Hal Varian, Google’s Chief Economist
The McKinsey Quarterly, Jan 2009

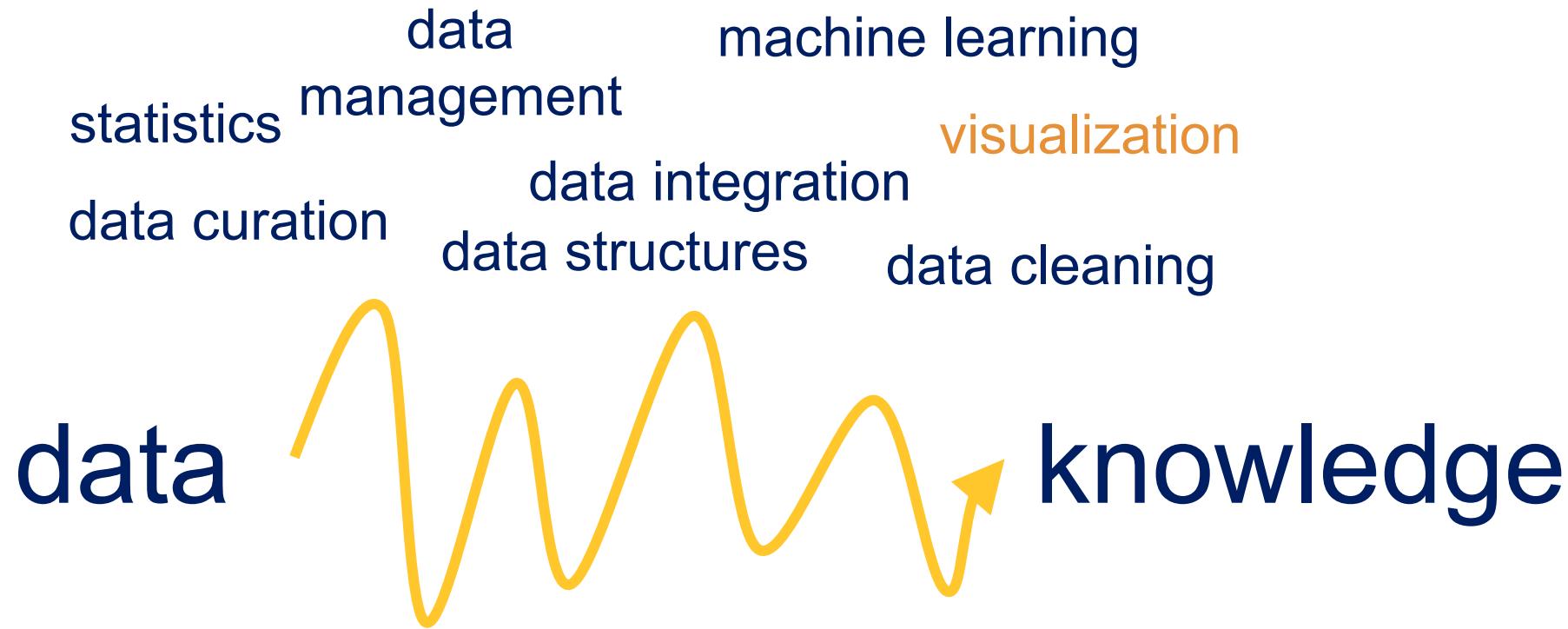
Data to knowledge

data → knowledge

Data to knowledge

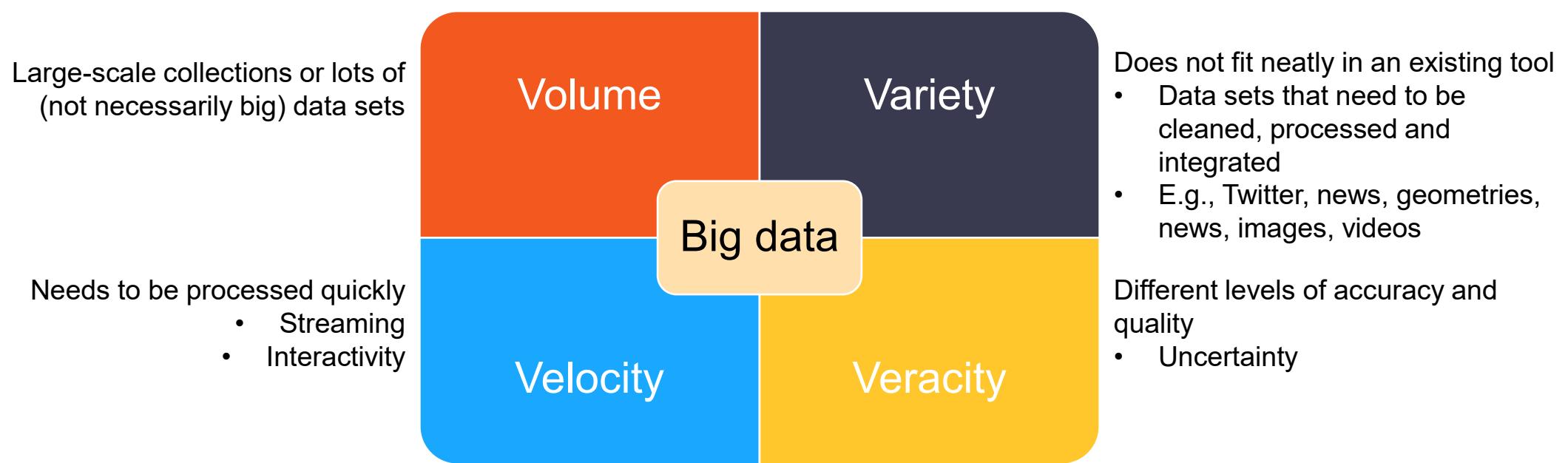


Data to knowledge

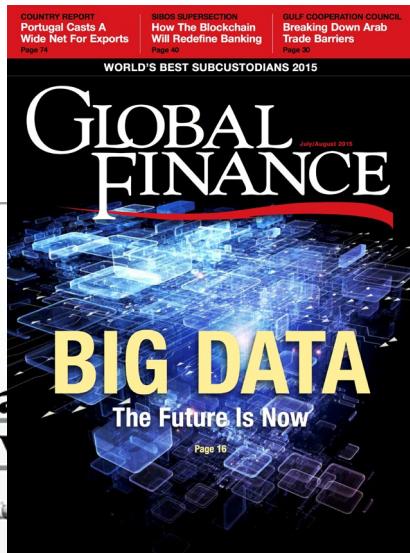


What is Big Data?

- Broad term for data so large and complex that traditional data processing applications are inadequate.

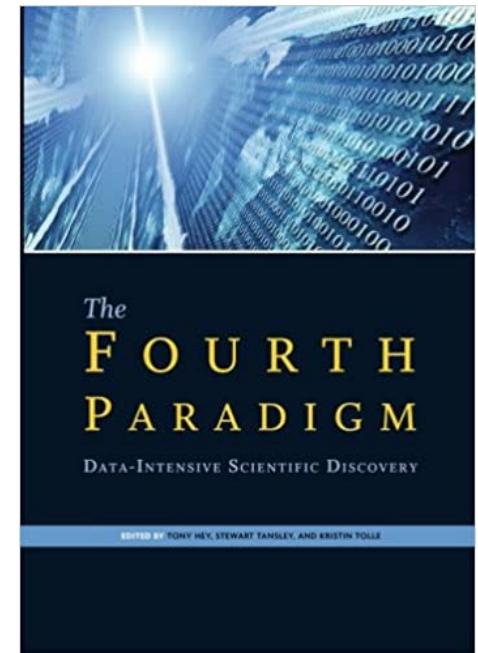


The big deal of big data

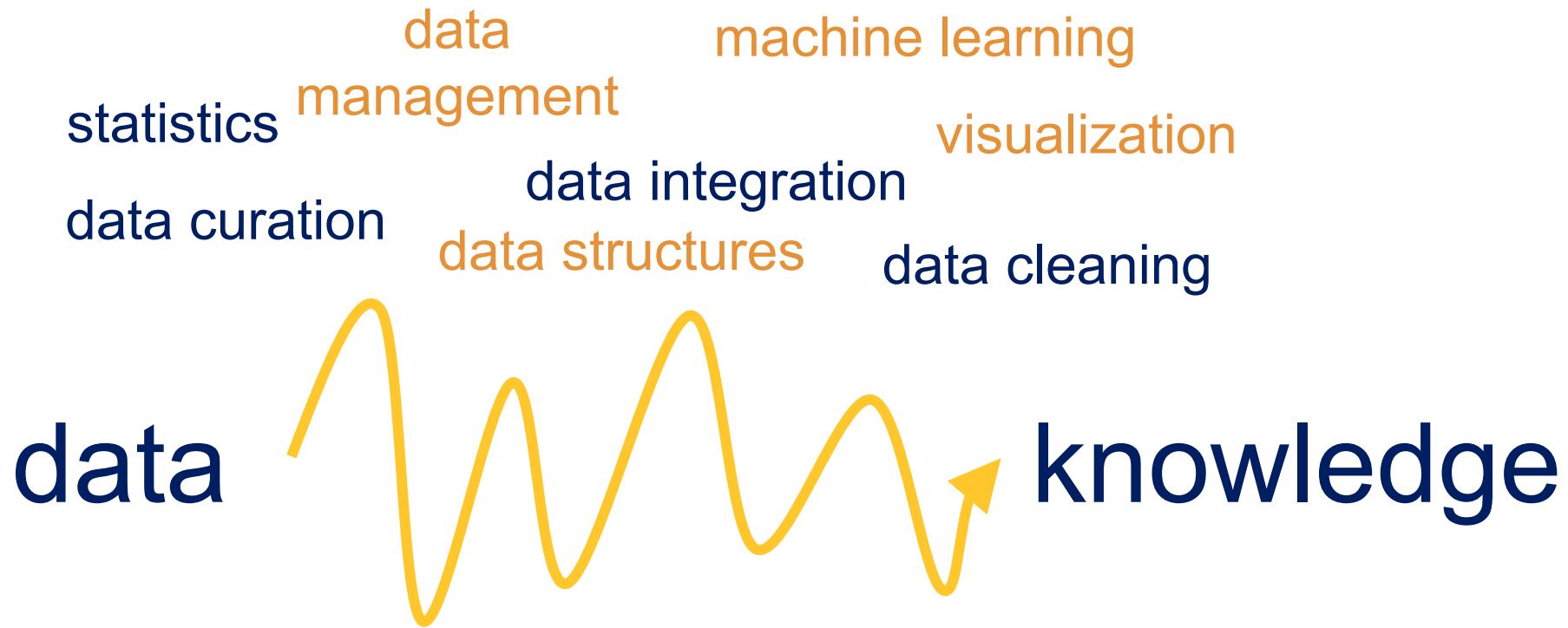


The big deal of big data

- **Science**: several domains are moving toward data-driven exploration: *increasing use of data is bringing a paradigm shift to the nature of science.*
 - Climate data, urban data, physics data, health data.
- **Industry**: companies are capitalizing on data – users are consuming and producing data.
 - Social media data, crowdsourced data, sensor data.
- **Government**: agencies use data to operate efficiently, make policies, make informed decisions.
 - Data.gov: 224,669+ datasets; NYC Open Data: 1,400 datasets.



Big data to knowledge



Big data analysis: Common practices

1. Domain experts and policy makers formulate hypotheses.
 2. Computer scientists or data scientists select data sets and slices, perform analyses, and derive plots.
 3. Domain experts examine the plots, go to step 1.
- Issues:
 - Dependency on computer scientists or data scientists distances domain experts from the data.
 - Batch-oriented analysis pipeline hampers exploration – mostly confirmatory analyses.
 - Data are complex – often multivariate spatiotemporal.
 - Analysis limited to samples.

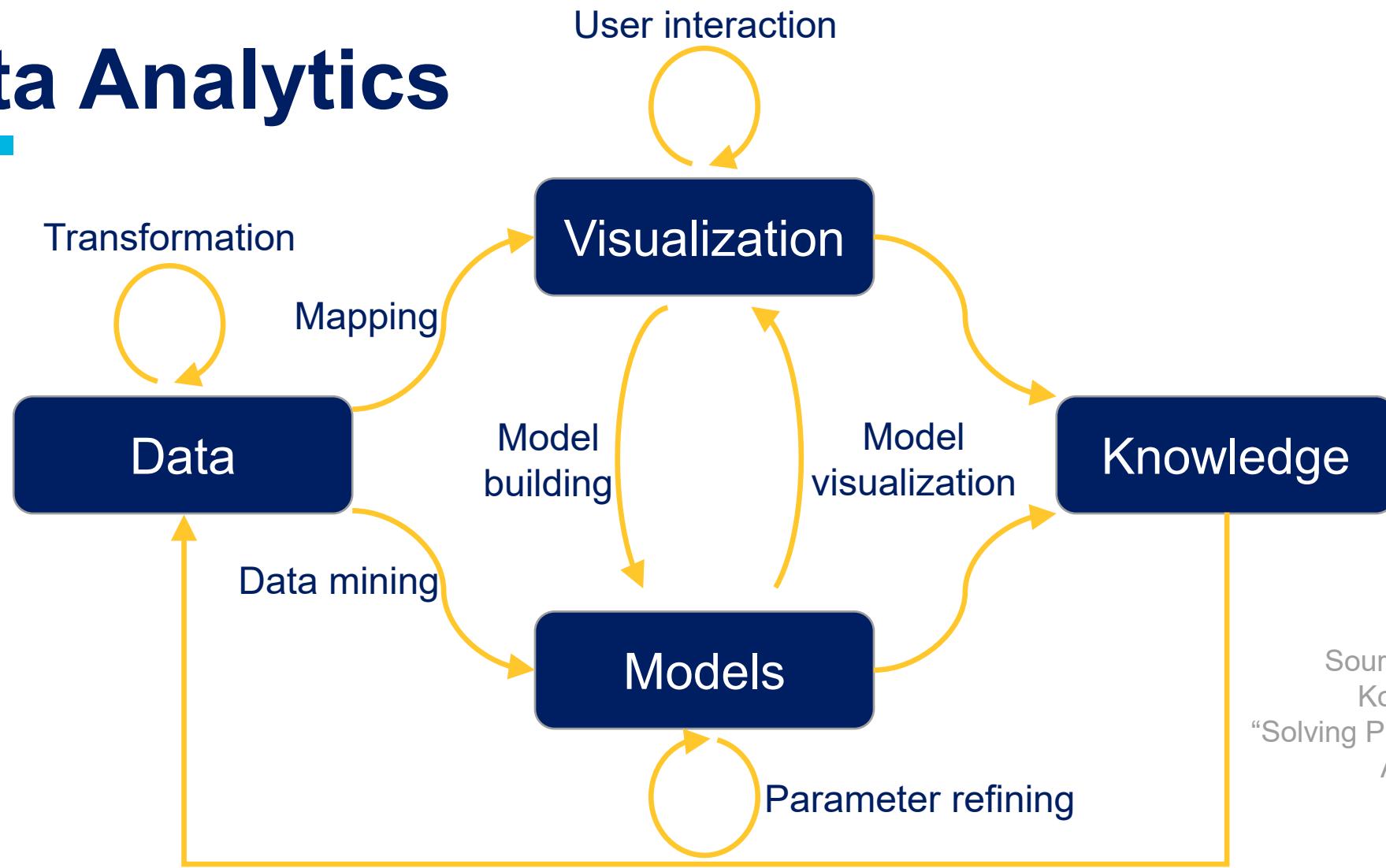
Big data analysis: What we need

- Scalable tools and techniques that help domain experts find, clean, integrate, ***interactively*** explore, and explain data.
- Guide users in the exploration process.
- Support interactive queries:
 - “*increased latency reduces the rate at which users make observations, draw generalizations, and generate hypotheses*”.
- Interdisciplinary:
 - “*as data scale and complexity increases, the novel solutions that will ultimately enable interactive, large-scale exploratory data analysis will have to come from truly interdisciplinary work*”.

[Liu and Heer, IEEE TVCG 2014]

[Chang, Fekete, Freire and Scheidegger, Dagstuhl Reports 2017]

Data Analytics



Source: Kleim and Kohlhammer,
“Solving Problems with Visual
Analytics”

Big data challenges

- “*Although modern database management systems (DBMS’s) allow users to perform complex scientific analyses over large datasets, DBMS’s are not designed to respond to queries at interactive speeds.*”

[Battle, Chang, Stonebraker SIGMOD 2016]

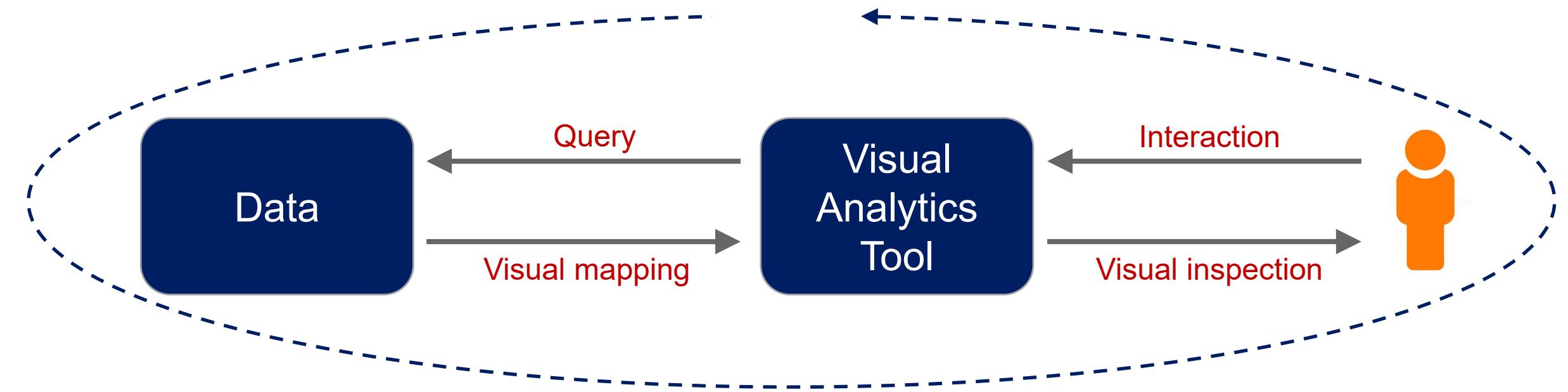


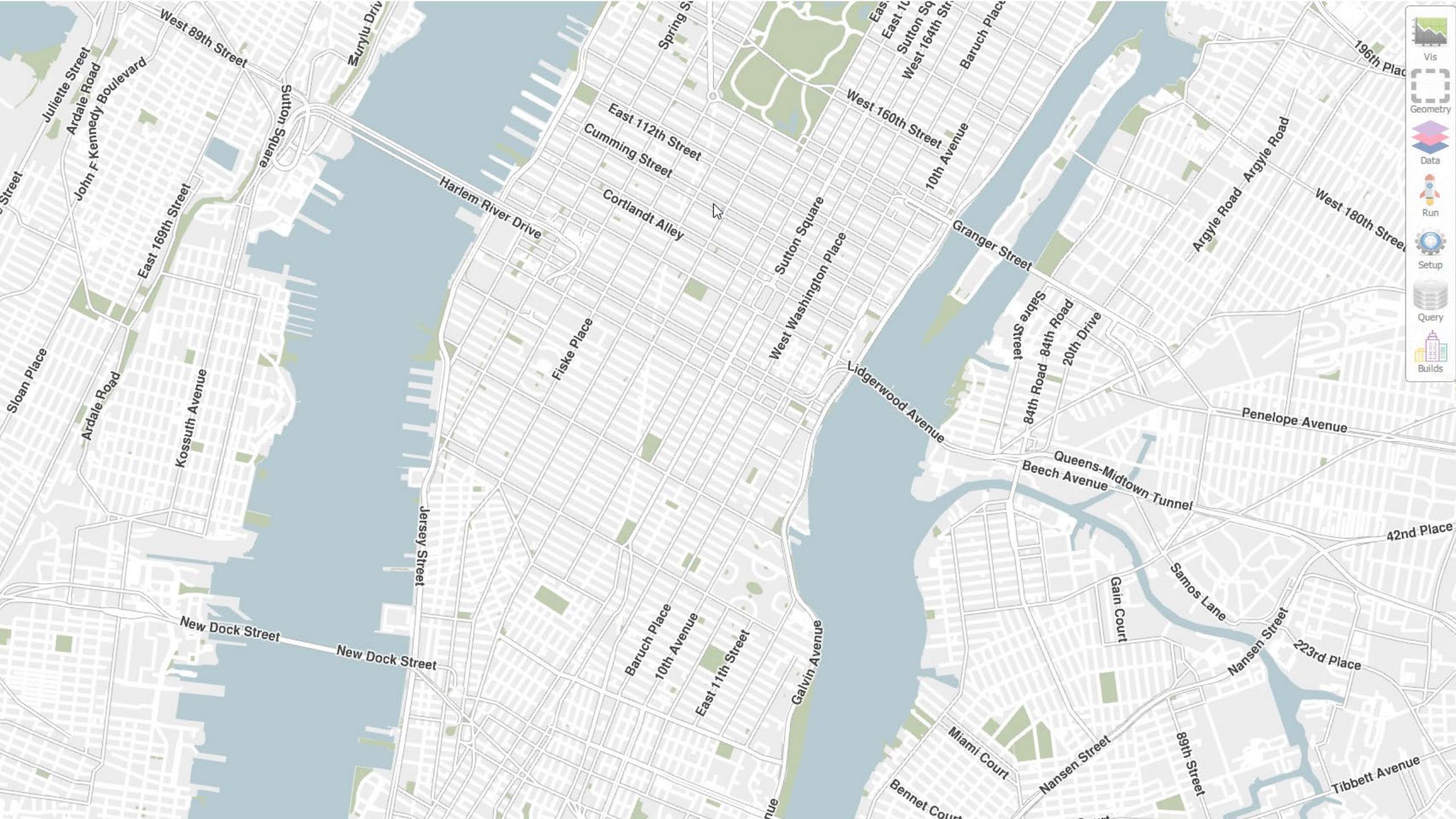
Big data challenges

- Data are vast and produced at unprecedented rates.
 - Sources are broad, varied, and unreliable.
- Computational processes are required to extract insight.
 - Hard to assemble and require expertise in a wide range of topics.
- Exploratory task are inherently iterative as one tests and formulates hypotheses:
 - “*An analysis has 30 different steps. It is tempting to just do this then that and then this. You have no idea in which ways you are wrong and what data is wrong*”.

[Kandel et al., VAST 2012]

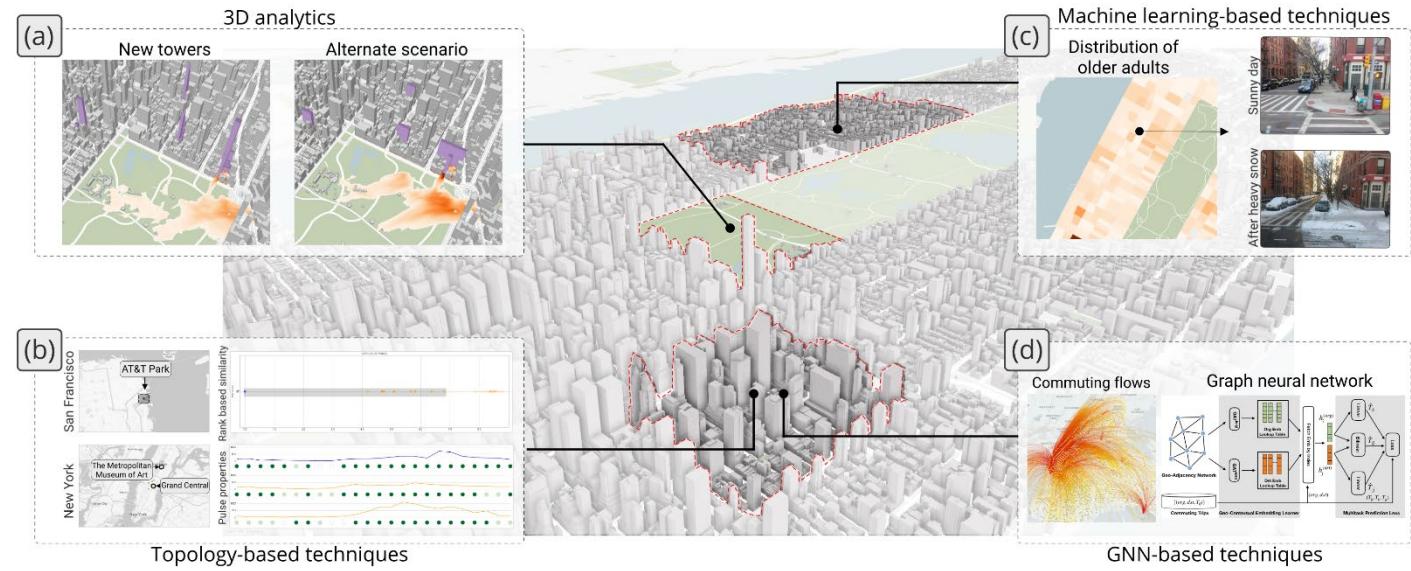
Interactive visual analysis



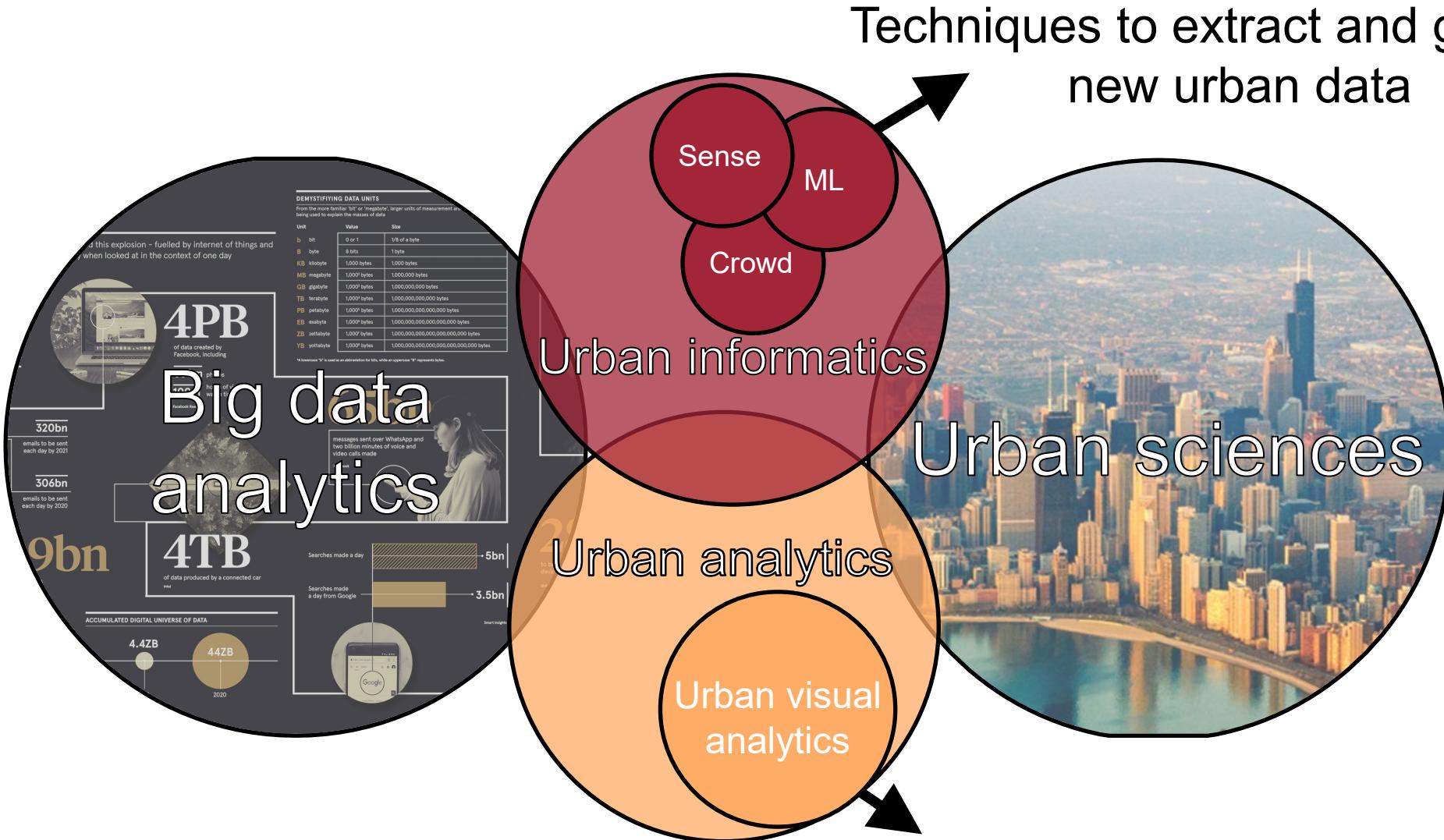


Research

Assistant Professor, CS (UIC)
PhD, CS 2018 (NYU)
MSc, CS 2012 (PUC-Rio)
BSc, CS 2009 (UFMG)



- Methods and techniques that follow a **human-centered approach to data science**, fostering the involvement of domain experts in the analysis process of big data.
- **Interactive tools and frameworks** that combine visualization, data management, human-computer interaction, and machine learning to support data-driven decision making by domain experts.

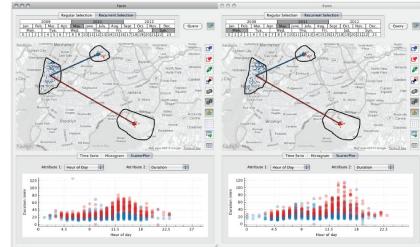


Techniques to extract and generate new urban data



Visual analytics tools and systems to analyze large urban data

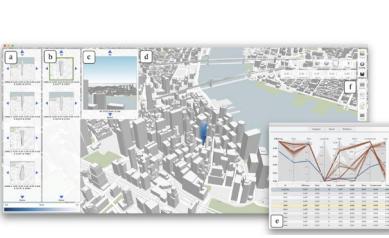




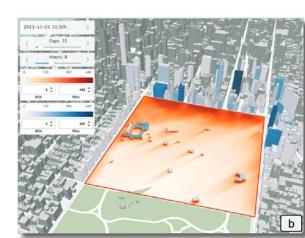
TaxiVis
(Ferreira et al., 2013)



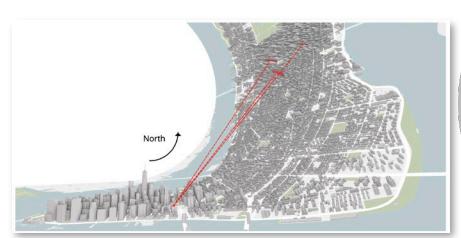
Urbane
(Ferreira et al., 2015)



Catalogue
(Doraiswamy et al., 2015)



Shadow Profiler
(Miranda et al., 2019)



UrbanRama
(Chen et al., 2020)



UTK
(Moreira et al., 2023)

2014

2016

2018

2020

2022

2023

2013

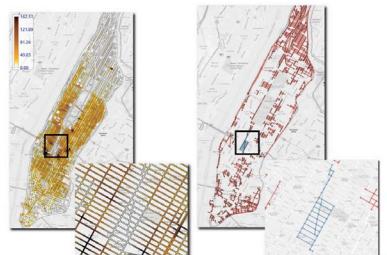
2015

2017

2019

2021

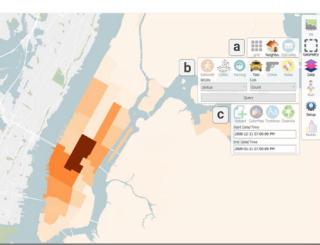
Taxi Patterns
(Doraiswamy et al., 2016)



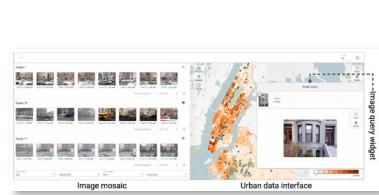
Urban Pulse
(Miranda et al., 2016)



Raster-Join
(Doraiswamy et al., 2018)

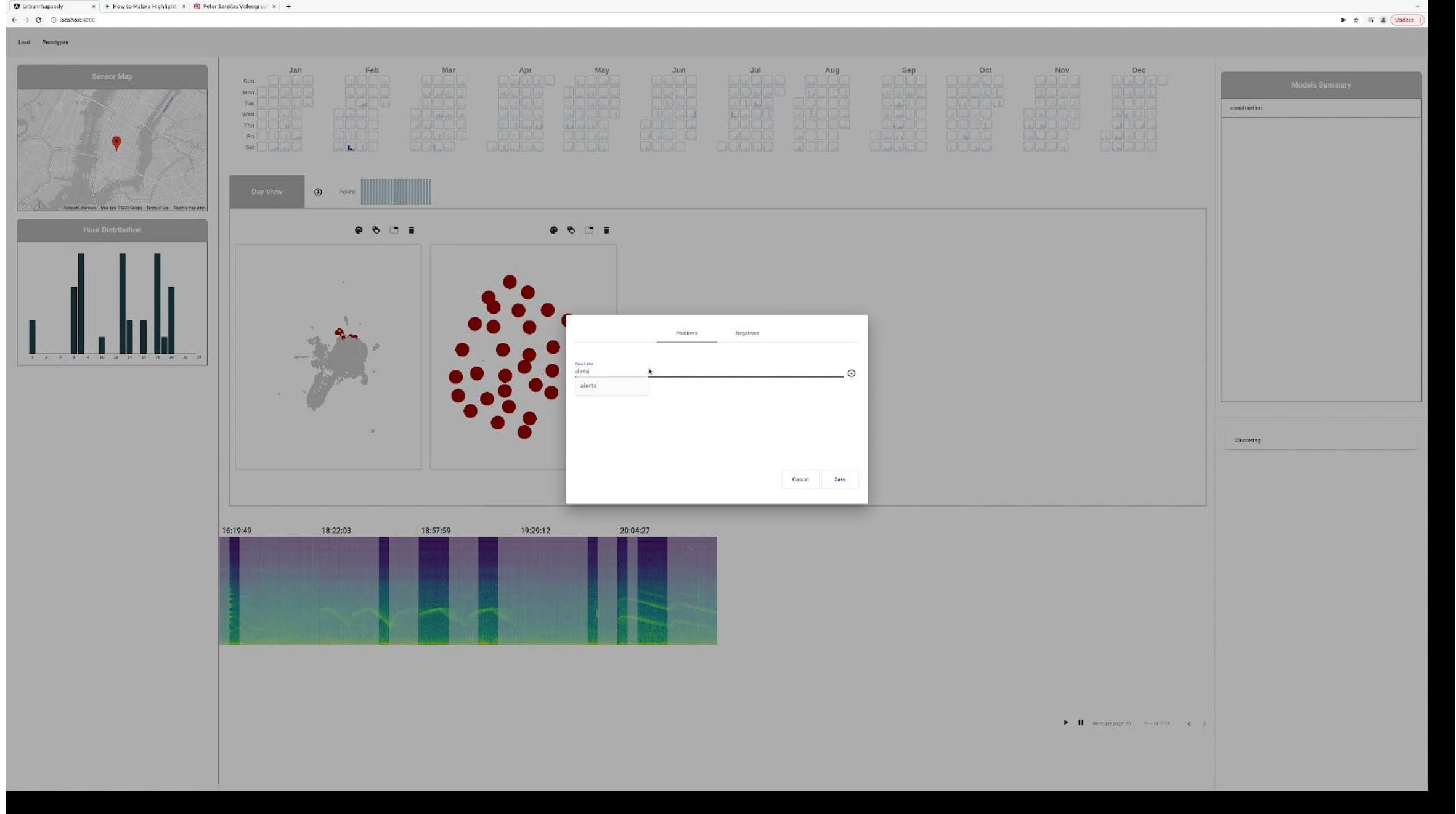


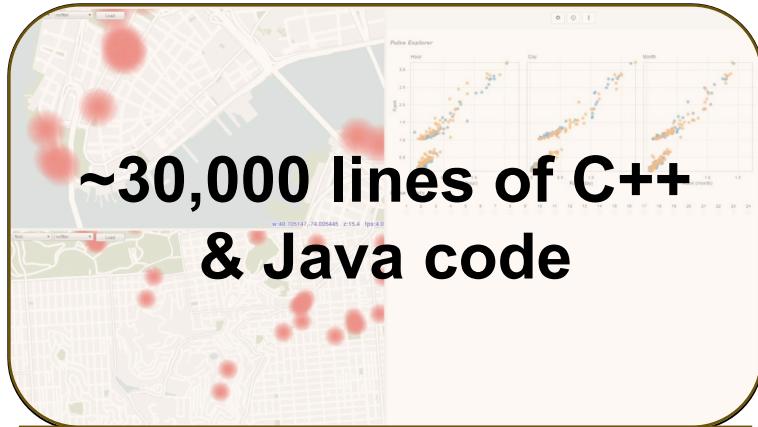
Urban Mosaic
(Miranda et al., 2020)



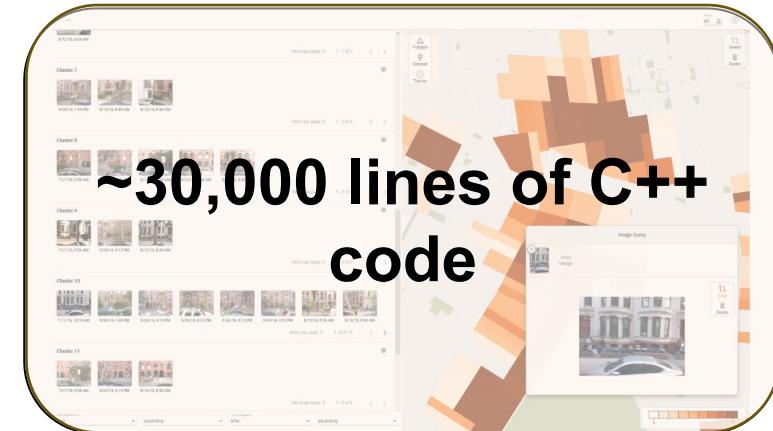
Urban Rhapsody
(Rulff et al., 2022)



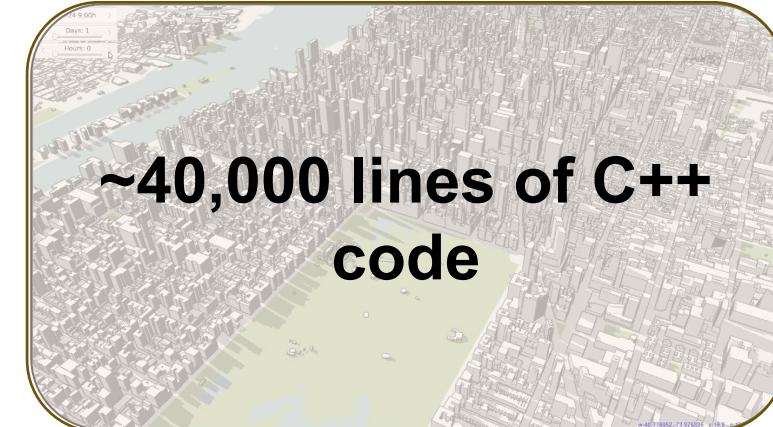




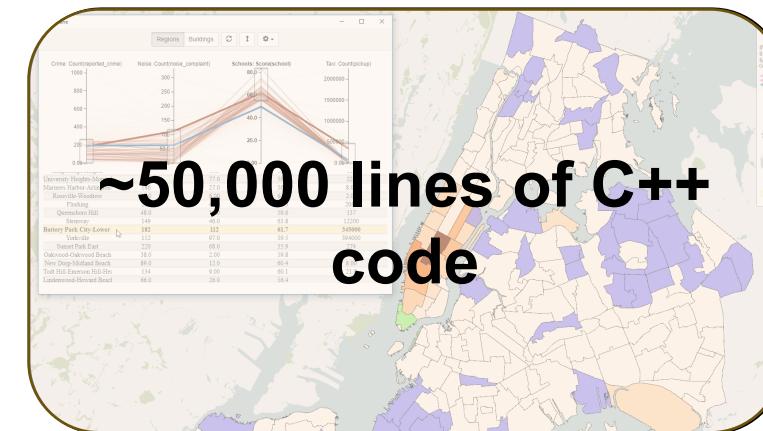
Urban Pulse:
Large-scale data mining of social media data



Urban Mosaic:
Interactive exploration of large imagery data



Shadow Profiler:
City-scale assessment of sunlight access



Urbane:
Interactive exploration of large data

```

1 #include "UrbaneMapView.hpp"
2
3 #include <QApplication>
4 #include "../MapView/BuildingRenderingLayer.hpp"
5
6 #include "../MassingGeneration/massinggeneration.h"
7 #include "../Util/ColorMapDivergent.hpp"
8 #include "UrbaneManager.hpp"
9
10 #include <QElapsedTimer>
11 #include <QThread>
12 #include <QDir>
13
14 #include <vector>
15
16 UrbaneMapView::UrbaneMapView(const QString &filename, const QRectF &vp, QWidget *parent)
17 | : MapView(filename, vp, parent), graphLayer(NULL)
18 {
19     initialized = false;
20     skyExposureData = false;
21     this->centerIndex = GridIndex(1024, 1024);
22     this->currentLayer = NULL;
23     this->lotUpdate = true;
24 }
25
26 UrbaneMapView::~UrbaneMapView() {}
27
28 void UrbaneMapView::initializeGL() {
29     if(!initialized) {
30         MapView::initializeGL();
31         this->buildingScore.initComputeShader();
32         this->skyScore.initComputeShader();
33     }
34     initialized = true;
35 }
36
37 void UrbaneMapView::paintGL()
38 {
39     this->showOsd(false);
40
41     // Lot data initialization in manager
42     // TODO Don't know of a better place to do this
43     if(lotUpdate &amp; this->parcelLayer->isDataReady()) {
44         updateLotDataDB();
45         lotUpdate = false;
46     }
47
48     UrbaneManager *manager = UrbaneManager::getInstance();
49     QPair<RenderingOperation, UIOperation> state = manager->getState();
50
51     RenderingOperation operation = state.first;
52     UIOperation what = state.second;
53     switch(operation) {
54     case RenderingOperation::UpdateVis:
55     {
56         bool updateFunction = false;

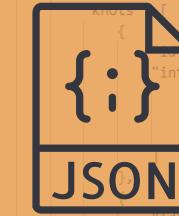
```



Abstracts low-level functionalities



Easy access to data analytics



Self-contained & sharable JSON file

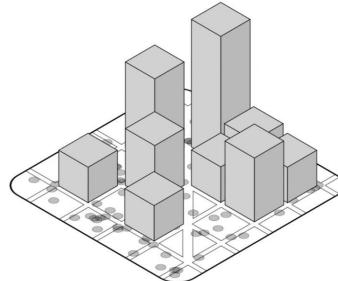


Lower the barrier for the construction of urban tools & systems

```

1 {
2     "components": [
3         "map": {
4             "position": [-13961, -115, 149, -56, -42],
5             "right": [946, 6354370117188, -423, 6624084472656, 497, 6396560824219],
6             "at": [1138, 2, 390, 1, 62, 63, 344, 13, 11, 56, 1],
7             "up": [0.018835, 0.26612, 0.15, 0.85, 0.1, 0.342, 0.01],
8             "interactions": ["NONE", "NONE", "NONE", "NONE"]
9         },
10        "plots": [
11            {
12                "id": "pureparks",
13                "integration_scheme": {
14                    "in": {"bin": true, "field": "shadowToBuildings_abstract"}, "out": {"name": "parks", "level": "OBJECTS"}
15                }
16            },
17            {
18                "id": "purewater",
19                "integration_scheme": {
20                    "in": {"bin": true, "field": "shadowToBuildings_abstract"}, "out": {"name": "water", "level": "OBJECTS"}
21                }
22            },
23            {
24                "id": "pureroads",
25                "integration_scheme": {
26                    "in": {"bin": true, "field": "shadowToBuildings_abstract"}, "out": {"name": "roads", "level": "OBJECTS"}
27                }
28            }
29        ],
30        "knots": [
31            {
32                "id": "pureparks",
33                "integration_scheme": {
34                    "in": {"bin": true, "field": "shadowToBuildings_abstract"}, "out": {"name": "parks", "level": "OBJECTS"}
35                }
36            },
37            {
38                "id": "purewater",
39                "integration_scheme": {
40                    "in": {"bin": true, "field": "shadowToBuildings_abstract"}, "out": {"name": "water", "level": "OBJECTS"}
41                }
42            },
43            {
44                "id": "pureroads",
45                "integration_scheme": {
46                    "in": {"bin": true, "field": "shadowToBuildings_abstract"}, "out": {"name": "roads", "level": "OBJECTS"}
47                }
48            }
49        ],
50        "interactions": [
51            {
52                "id": "shadowToBuildings",
53                "integration_scheme": {
54                    "in": {"bin": true, "field": "shadowToBuildings_abstract"}, "out": {"name": "buildings", "level": "OBJECTS"}
55                }
56            }
57        ]
58    ]
59 }

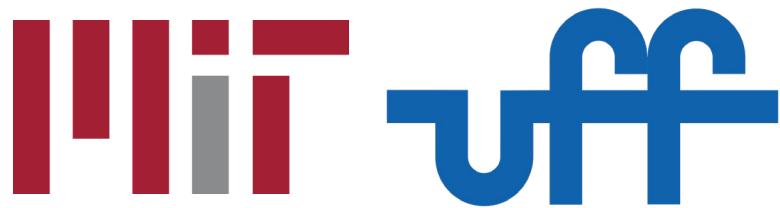
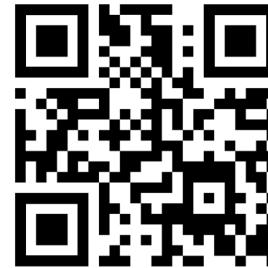
```



UrbanTK

.....

Code & tutorials:
urbantk.org



Gustavo Moreira, University of Illinois Chicago
Maryam Hosseini, Massachusetts Institute of Technology
Md Nafiul Alam Nipu, University of Illinois Chicago
Marcos Lage, Universidade Federal Fluminense
Nivan Ferreira, Universidade Federal de Pernambuco
Fabio Miranda, University of Illinois Chicago

The Urban Toolkit
A Grammar-based Framework for Urban Visual Analytics

Getting Started GitHub Tutorials

While cities around the world are looking for smart ways to channel new advances in data collection, management, and analysis to address their day-to-day problems, the complex nature of urban issues and the overwhelming amount of available structured and unstructured data have posed significant challenges in translating these efforts into actionable insights. In the past few years, urban visual analytics tools have significantly helped tackle these challenges. With this in mind, we present the Urban Toolkit, a flexible and extensible visualization framework that enables the easy authoring of web-based visualizations through a new high-level grammar specifically built with common urban use cases in mind.

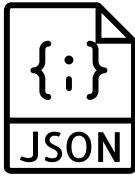
The toolkit is described in the [paper](#):
The Urban Toolkit: A Grammar-based Framework for Urban Visual Analytics
Gustavo Moreira, Maryam Hosseini, Md Nafiul Alam Nipu, Marcos Lage, Nivan Ferreira and Fabio Miranda
IEEE Transactions on Visualization and Computer Graphics (Accepted at IEEE VIS 2023, to appear)

[Moreira, VIS 2023]



Transportation experts

- What-if scenarios
- Model inspection
- ...



Weather experts

- What-if scenarios
- Model inspection
- Data wrangling



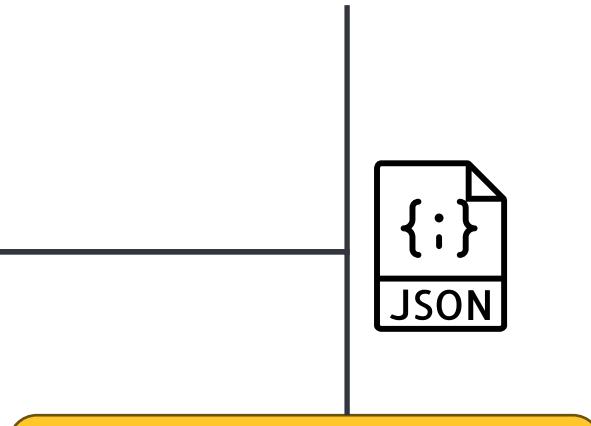
The Urban Toolkit

- Open urban data
- Modeling & simulation results
- Crowdsourced data
- Urban sensing data



Communities

- Engagement
- ...



Policy makers

- What-if scenarios
- Engagement
- ...



CS524: Big Data Visualization & Analytics

- Intersection between visualization, data management and data mining, covering the necessary topics to build visual analytics tools to handle big data.
- Broad definition of big data: any dataset with size (or complexity) that goes beyond the ability of standard tools and techniques to **interactively** manage and process it.
- At the end of the course, you will be able to:
 - Design and implement visual analytics systems capable of handling large data (combining visualization and data mining techniques, data structures and algorithms)

CS524: Big Data Visualization & Analytics

- Necessary techniques to build visual analytics tools to handle big data:
 1. Building blocks: current technologies and libraries to build visual analytics systems for big data.
 2. Visualization: visual implications of handling big data.
 3. Data management: techniques to handle big data.
 4. Analytics: data mining and technical frameworks to extract patterns or features that can drive visual analytics systems.

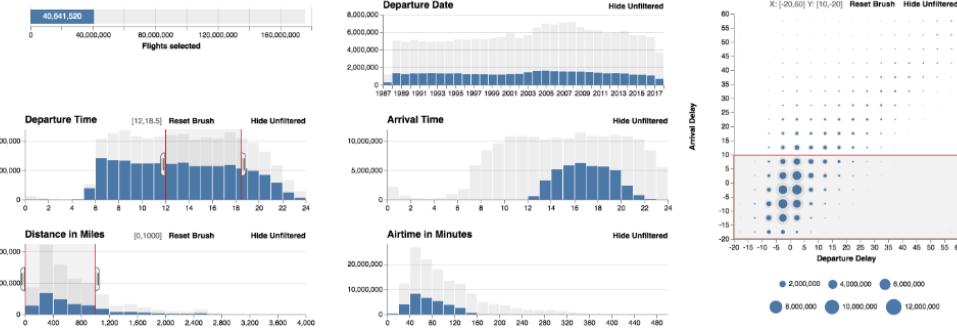
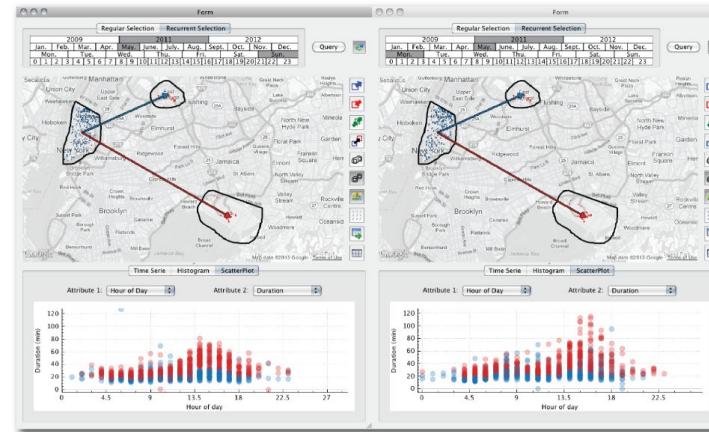
Building blocks:

- Necessary tools and frameworks for the development of big data analytic systems:
 - Angular (Javascript)
 - Boost (C++)
 - Qt (C++)
 - Flask (Python)



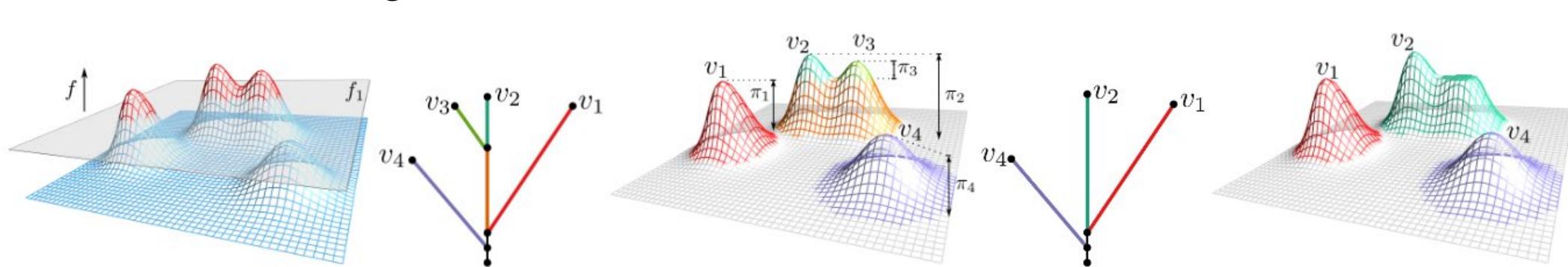
Visualization:

- Visual analytics systems
 - Interactivity requirements.
 - Components of the system.
- Progressive visualization
 - Progressively build visualizations, maintaining interactivity.
- Uncertainty visualization
 - Defining uncertainty.
 - Displaying uncertainty.
- Visualization for machine learning



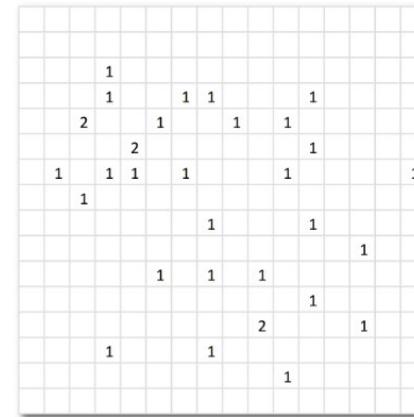
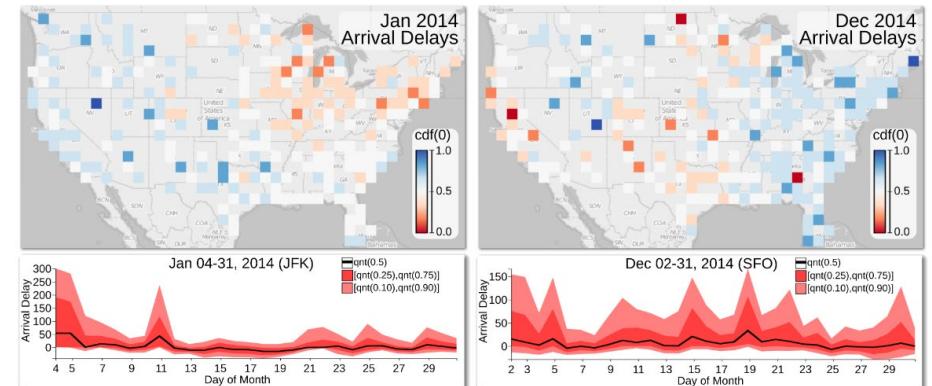
Analytics:

- Data mining and analytics techniques to extract features and patterns from big data:
 - Computational topology
 - Wavelet
 - Techniques for streaming data
- Machine learning for visualization

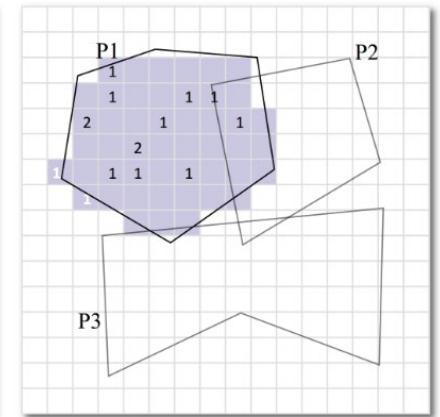


Data management:

- Approximate queries, hashing, learned indices
- Spatial structures
 - Nanocubes
- Spatial queries
 - GPU-based indices
- MapReduce
 - Hadoop, Spark



(a)



(b)

CS524: Big Data Visualization & Analytics

