

Scientific Discovery Environments: **A _{bias} View Towards “The Next Generation”**

Claudio T. Silva

Tandon School of Engineering

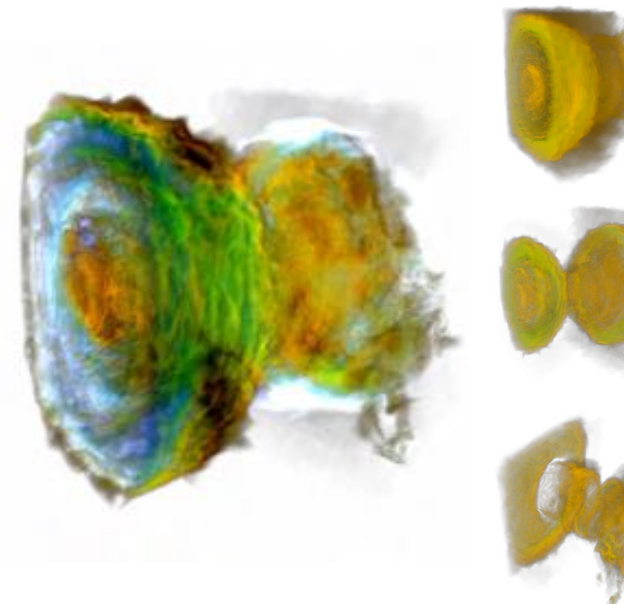
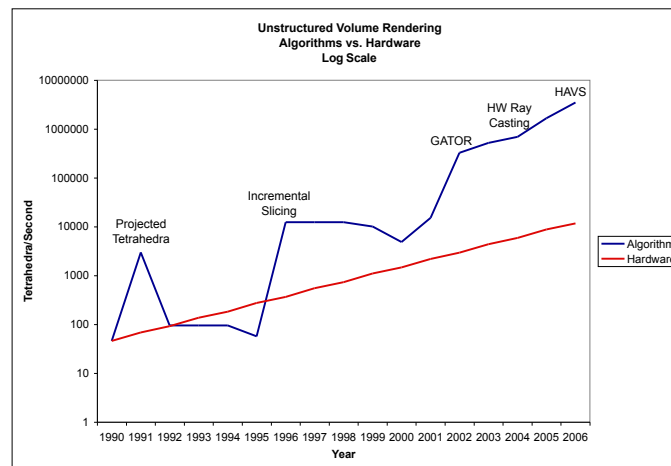
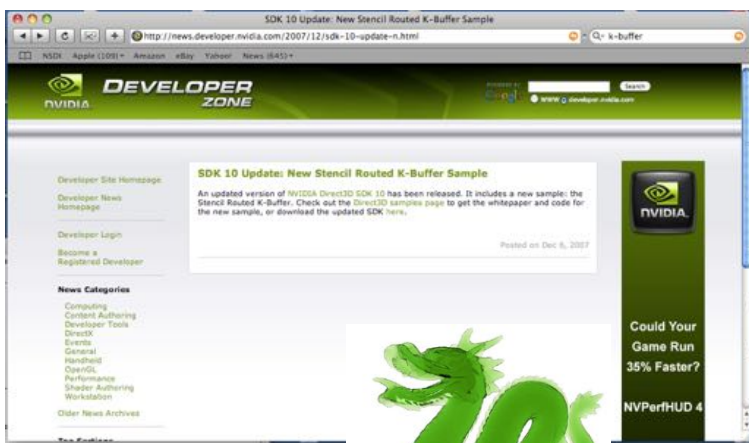
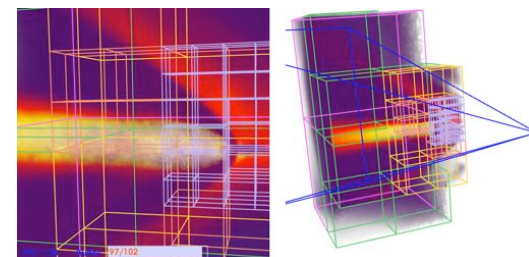
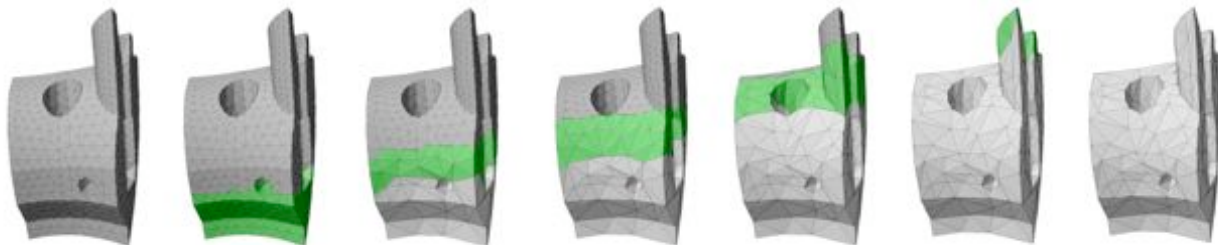
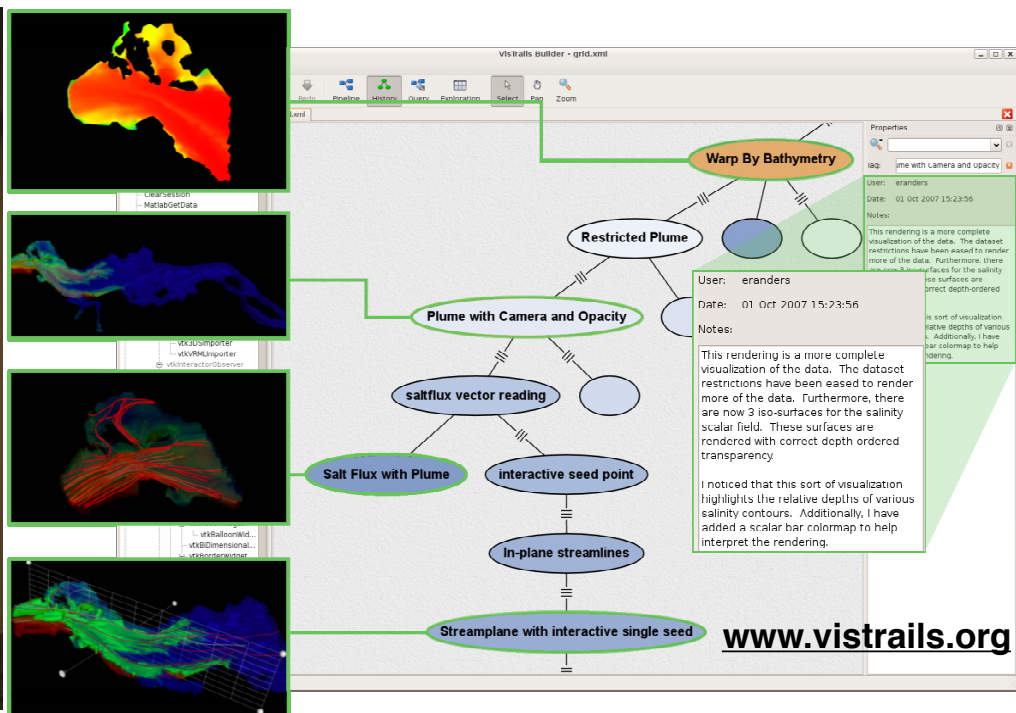
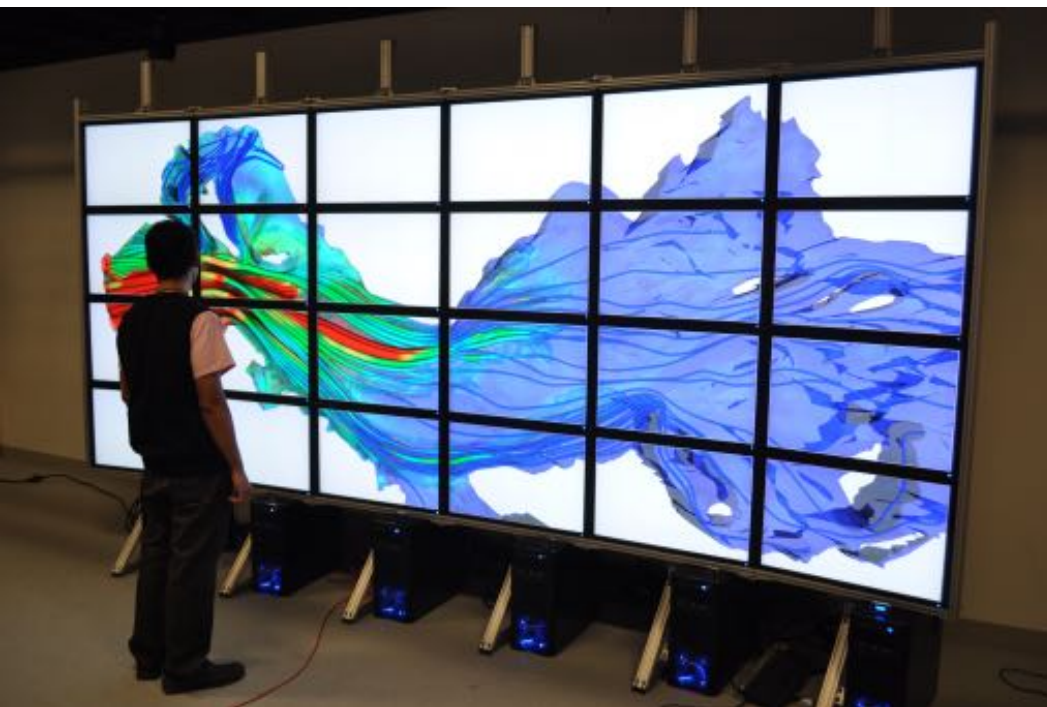
Center for Data Science

Center for Urban Science + Progress

Courant Institute for Mathematical Sciences

New York University

Funded by Moore and Sloan Foundations,
DARPA, NSF, NASA, DOE, MLB.com, AT&T, NVIDIA, and IBM



Sports Data: MLB.com Statcast

The screenshot shows a web browser displaying a Vice Sports article. At the top, there's a banner for 'VICE Magazine's Prison Issue' with an 'OUT NOW' button. Below this is the 'MOTHERBOARD' logo and a navigation bar with links like 'Watch', 'Machines', 'Discoveries', 'Space', 'Futures', 'Gaming', and 'Earth'. The main video player shows a baseball field with overlaid statistical data and the text 'Future of the Game: Baseball's Latest Statistical Revolution'. Below the video is the article title 'Behind the Scenes of Major League Baseball's Futuristic Player Tracking System' by Jason Koehler, dated September 25, 2015. The article text begins with 'Because my beloved Orioles flamed out early, the 2015 baseball season may seem utterly forgettable. But someday even I may look back on it as the dawn of a new era: The year stats took over.' and continues to discuss the statistical revolution in baseball. A Bacardi advertisement is visible on the right side of the article.

VICE Magazine's Prison Issue OUT NOW

MOTHERBOARD Watch Machines Discoveries Space Futures Gaming Earth

Future of the Game: Baseball's Latest Statistical Revolution

VICE SPORTS

Behind the Scenes of Major League Baseball's Futuristic Player Tracking System

Written by JASON KOEHLER

September 25, 2015 // 10:20 AM EST

Because my beloved Orioles flamed out early, the 2015 baseball season may seem utterly forgettable. But someday even I may look back on it as the dawn of a new era: The year stats took over.

Baseball has always been ultra stat focused—even old timers who disavow the SABR statistical revolution detailed in *Moneyball* used bad stats like RBI and batting average (doubtlessly culled from the backs of old Topps cards) to win arguments. But even

BACARDÍ UNSTABLE SINCE 1898

A NEW SERIES UNVEILING THE ARTIST BEYOND THE STAGE



NYU

TANDON SCHOOL
OF ENGINEERING

Big Urban Data: Understanding Cities

Infrastructure



Condition, operations

Environment



Meteorology, pollution,
noise, flora, fauna

People



Relationships,
economic activities, health,
nutrition, opinions, ...

- City components interact in complex ways
- Need to look at the city *data exhaust* to understand these interactions
- Processes occur over time and space

twitter



You Tube
Broadcast Yourself

NYC OpenData



data.gov in
Open Government Data Platform India



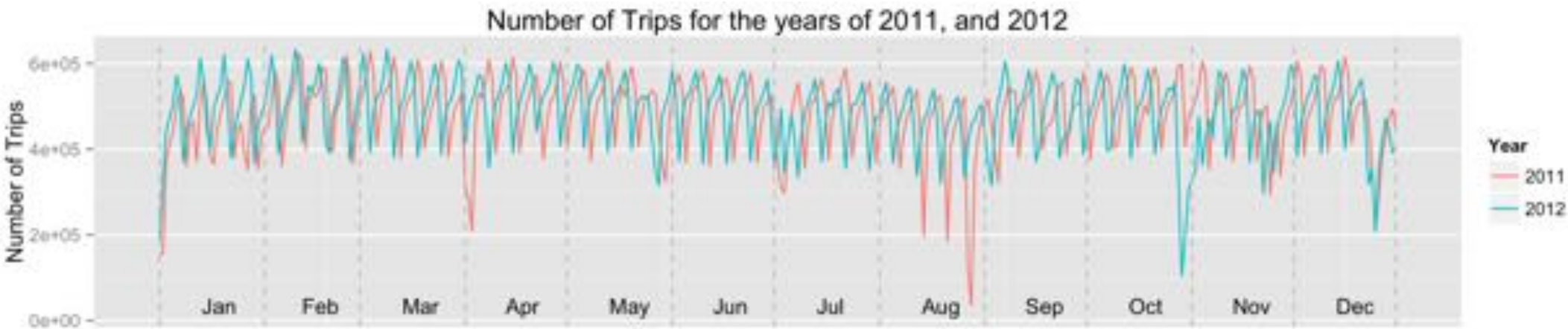
data.gouv.fr



NYU

TANDON SCHOOL
OF ENGINEERING

Exploring Urban Data: NYC Taxis



- Taxis are *sensors* that can provide unprecedented insight into city life: economic activity, human behavior, mobility patterns, ...

“How the taxi fleet activity varies during weekdays?”

“What is the average trip time from Midtown to the airports during weekdays?”

“How was activity in Midtown affected during a presidential visit?”

“How did the movement patterns change during Sandy?”

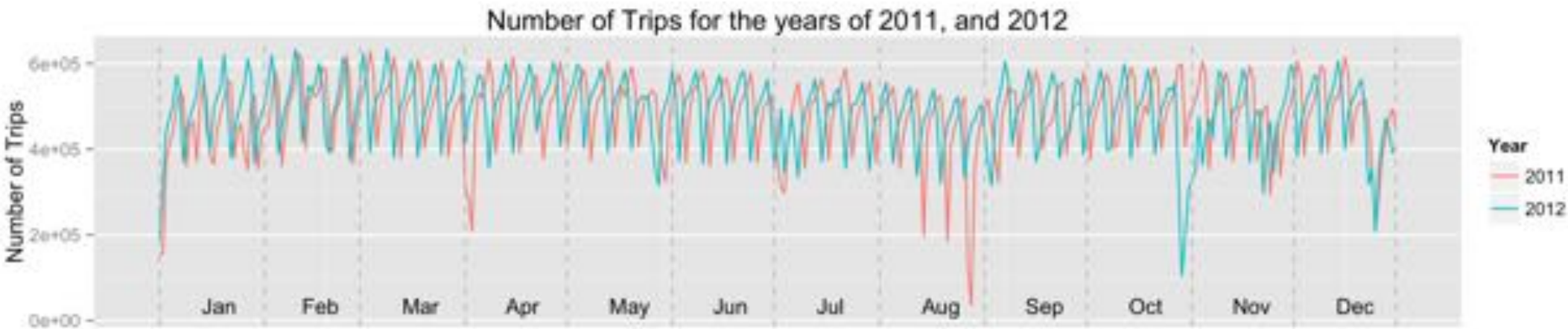
“Where are the popular night spots?”



NYU

TANDON SCHOOL
OF ENGINEERING

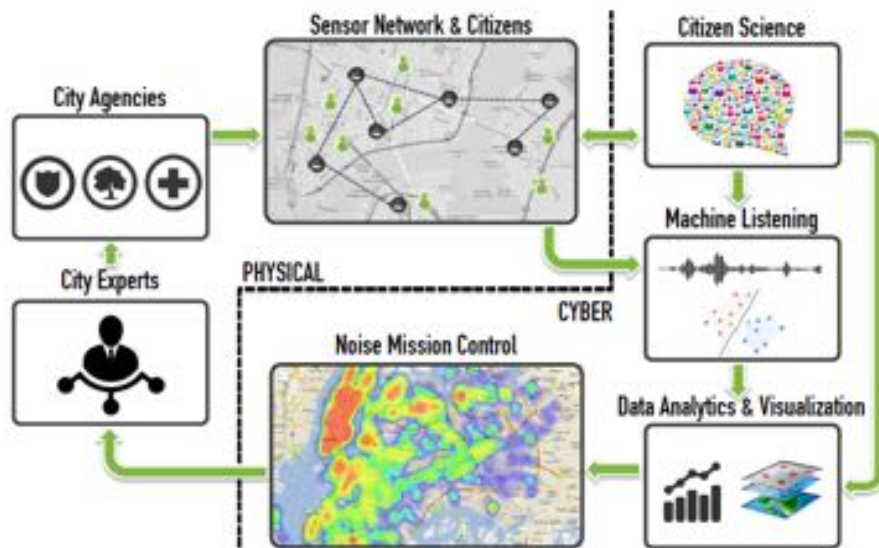
Exploring Urban Data: NYC Taxis



NYU

TANDON SCHOOL
OF ENGINEERING

Sounds of New York City



NYU

TANDON SCHOOL
OF ENGINEERING



NYU

TANDON SCHOOL
OF ENGINEERING

Major Trends

AI

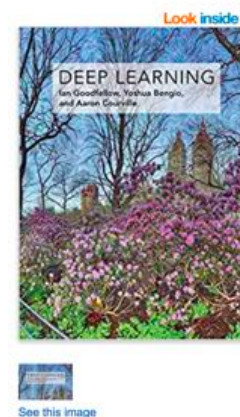
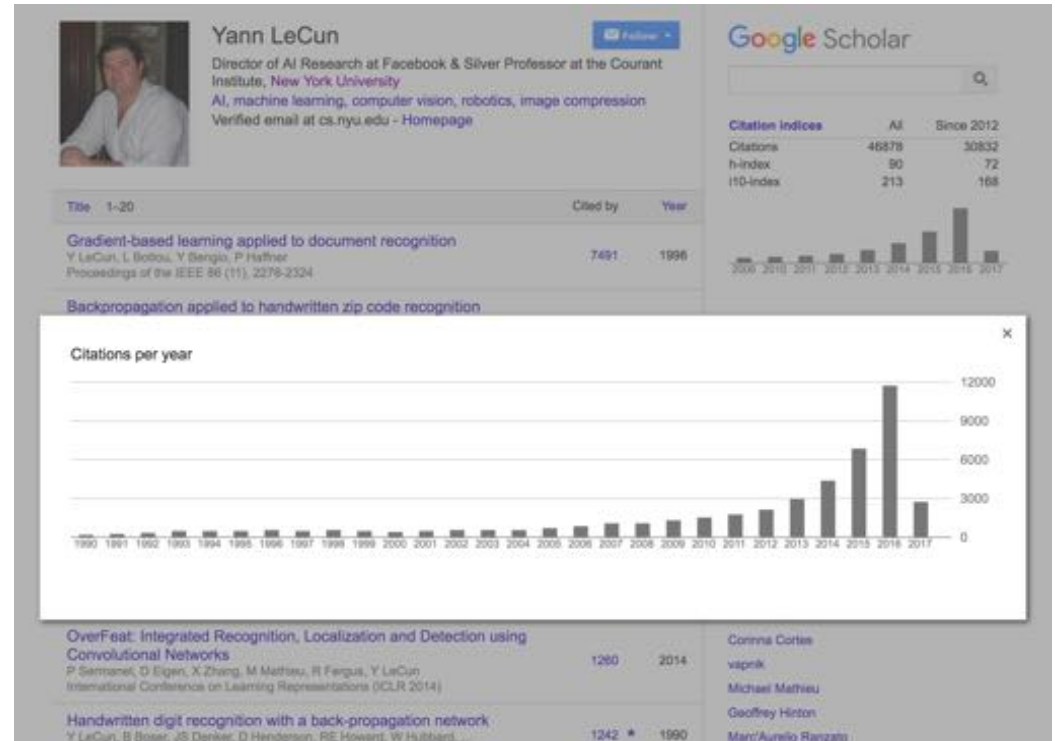
Deep learning
Machine learning
Natural language
Automatic analysis
Anomaly detection

Data management

Linking analysis and data
Provenance
Collaborative support
Lots of simultaneous datasets
Lots of visualizations
Lots of images

GUI

Desktop
Touch-enabled interfaces
Large Displays
VR/AR
General versus problem specific
System support/architecture
Cloud environments
Interactive support/programming
Progressive Visualization
Client-server support
Parallelism



Deep Learning (Adaptive Computation and Machine Learning series) Hardcover – November 18, 2016
by Ian Goodfellow (Author), Yoshua Bengio (Author), Aaron Courville (Author)
★★★★★ 57 customer reviews
#1 Best Seller in Artificial Intelligence & Semantics

See all 2 formats and editions

Kindle \$60.80	Hardcover \$71.50 Prime
Read with Our Free App	16 Used from \$79.82 24 New from \$67.51

Note: Available at a lower price from other sellers, potentially without free Prime shipping.

"Written by three experts in the field, *Deep Learning* is the only comprehensive book on the subject." – **Elon Musk**, cochair of OpenAI, cofounder and CEO of Tesla and SpaceX

Deep learning is a form of machine learning that enables computers to learn from experience and understand the world in terms of a hierarchy of concepts. Because

Report incorrect product information.

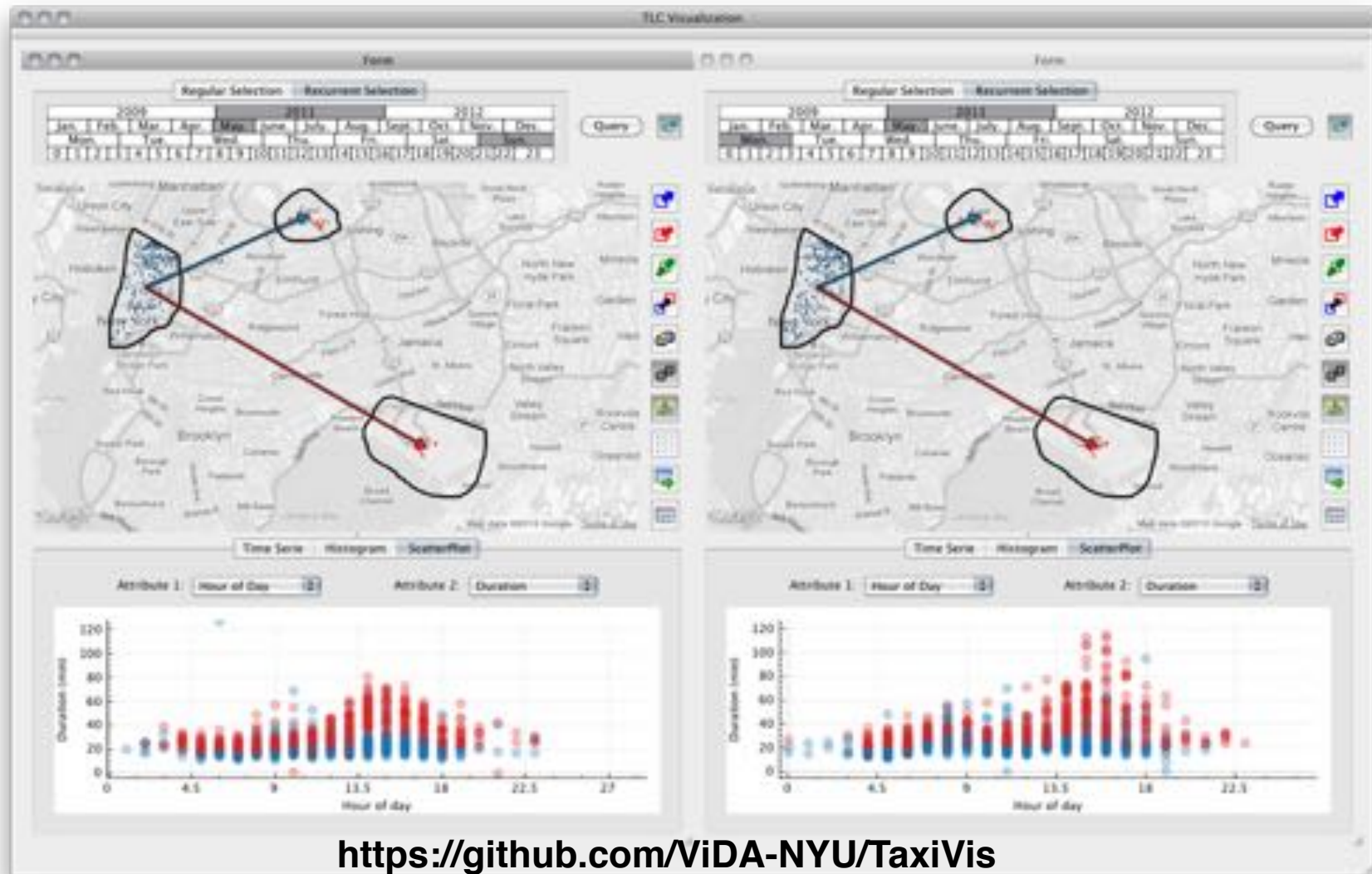


NYU

TANDON SCHOOL
OF ENGINEERING

Graphical User Interfaces

General versus problem specific



Ferreira et al, 2013



NYU

TANDON SCHOOL
OF ENGINEERING



VISUALIZATION
IMAGING AND
DATA ANALYSIS
CENTER

Provenance

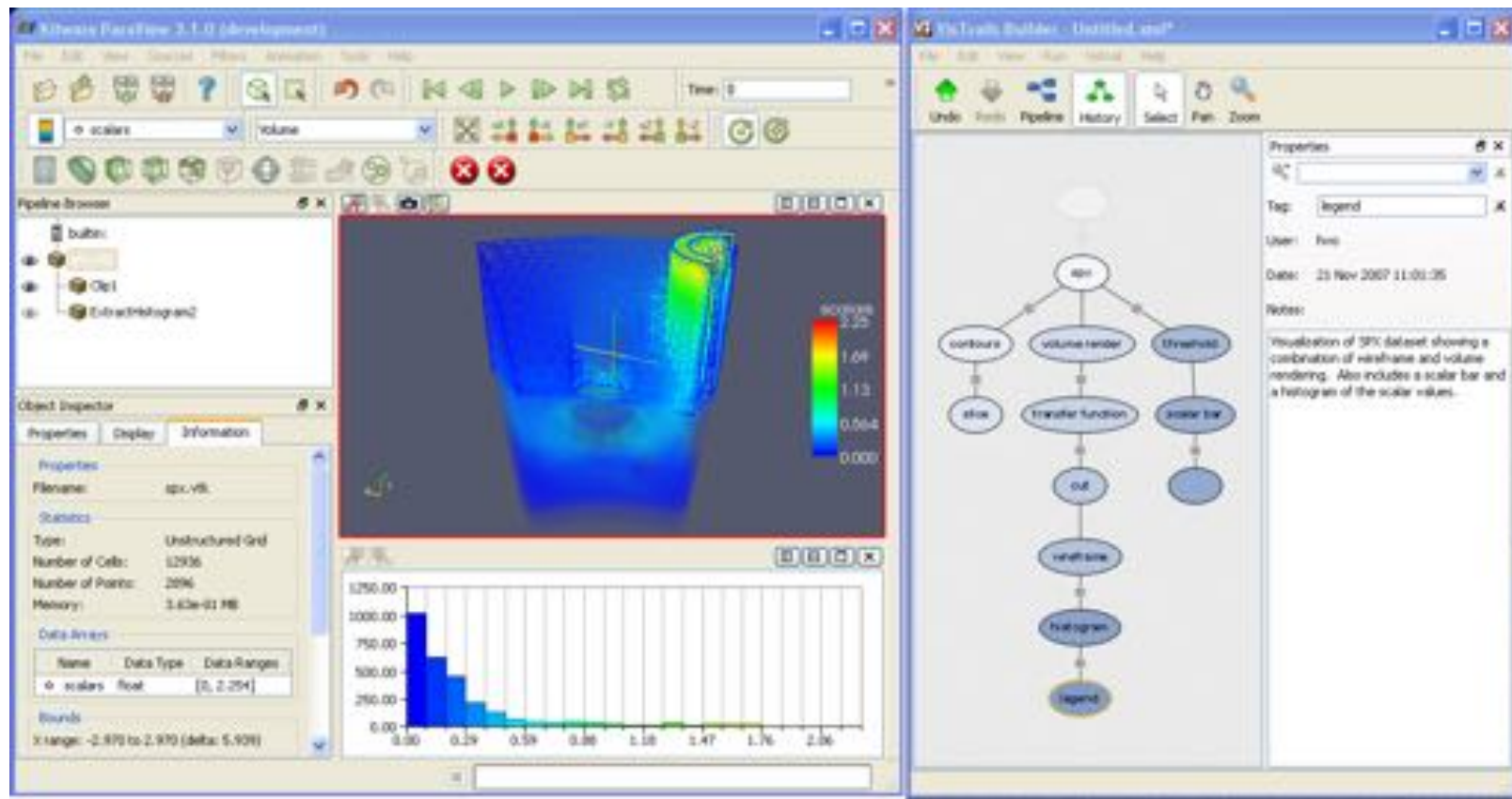


Fig. 2. A screenshot of ParaView (left) with the provenance captured by VisTrails and displayed as a version tree in a separate window (right). This preliminary prototype taps into ParaView undo/redo mechanism to capture the exploration process.

Callahan et al, 2008



NYU

TANDON SCHOOL
OF ENGINEERING

Provenance

VisTrails Plugin for ParaView

Data management / UrbanGIS

Urbane: A 3D Framework to Support Data Driven Decision Making in Urban Development

Nivan Ferreira*
New York University
Marcos Lage†
Universidade Federal Fluminense
Harish Doraiswamy‡
New York University
Huy Vo§
New York University
Luc Wilson¶
Kohn Pedersen Fox Associates PC
Heidi Werner||
Muchan Park §
Claudio Silva ||
New York University



Figure 1: Urbane provides architects, developers, and planners with a new, data and analysis rich way of reading the city with the goal of improving decision making in urban development. Users can explore properties of neighborhoods and buildings using the data exploration view to identify underdeveloped sites for potential development. Then, using the visual interface together with the map view, they can simulate the impact of such development. For example, the views of the freedom tower (highlighted in green) from the buildings highlighted in red would be adversely impacted (positively impacted buildings are highlighted in blue) if the new constructions (colored yellow) are built. The supplemental video shows the different features and visualizations supported by Urbane.

ABSTRACT

Architects working with developers and city planners typically rely on experience, precedent and data analyzed in isolation when making decisions that impact the character of a city. These decisions are critical in enabling vibrant, sustainable environments but must also negotiate a range of complex political and social forces. This requires those shaping the built environment to balance maximizing the value of a new development with its impact on the character of a neighborhood. As a result architects are focused on two issues throughout the decision making process: a) what defines the character of a neighborhood? and b) how will a new development change its neighborhood? In the first, character can be influenced by a variety of factors and understanding the interplay between diverse data sets is crucial; including safety, transportation access, school quality and access to entertainment. In the second, the impact of a new development is measured, for example, by how it impacts the view from the buildings that surround it. In this paper, we work in collaboration with architects to design Urbane, a 3-dimensional

multi-resolution framework that enables a data-driven approach for decision making in the design of new urban development. This is accomplished by integrating multiple data layers and impact analysis techniques facilitating architects to explore and assess the effect of these attributes on the character and value of a neighborhood. Several of these data layers, as well as impact analysis, involve working in 3-dimensions and operating in real time. Efficient computation and visualization is accomplished through the use of techniques from computer graphics. We demonstrate the effectiveness of Urbane through a case study of development in Manhattan depicting how a data-driven understanding of the value and impact of speculative buildings can benefit the design-development process between architects, planners and developers.

Keywords: Urban data analysis; GIS; impact analysis; visual analytics; architecture; city development

1 INTRODUCTION

Why do two neighborhoods feel similar? Or different? Why does a new building change the quality of a neighborhood and another doesn't? While the experience of a city is inherently subjective, the characteristics that shape the quality of it are not. These characteristics can be difficult to obtain, measure or analyze by those shaping the future of a city. Architects working with developers and city planners typically rely on experience, precedent and data analyzed in isolation when making decisions that impact the character of a city. These decisions, while being critical in enabling vibrant and sustainable environments, must also negotiate a range of complex political and social forces. This requires those shaping the built environment to balance maximizing the value of new development

*e-mail:nivan.ferreira@nyu.edu

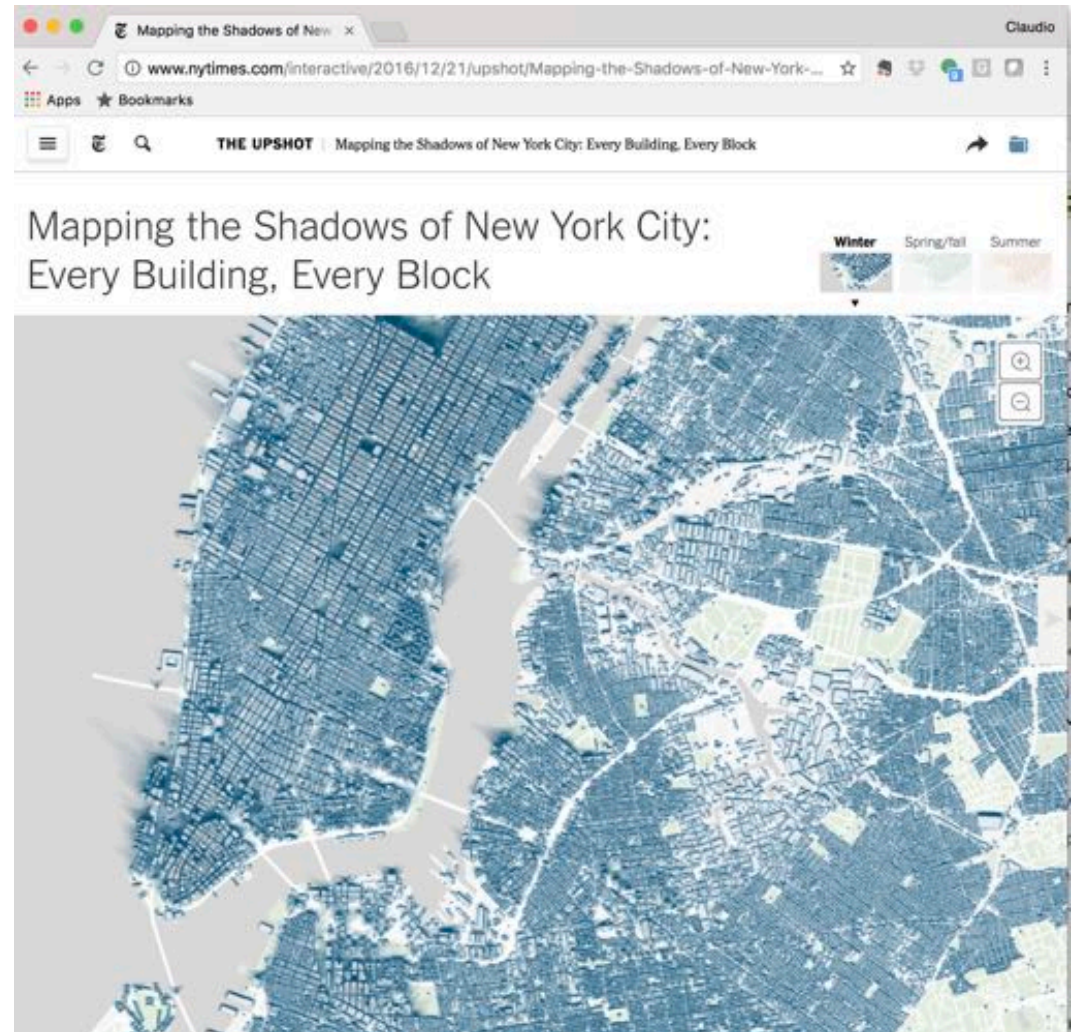
†e-mail:mlage@ic.uff.br

‡e-mail:harishd@nyu.edu

§e-mail:huy.vo@nyu.edu

¶e-mail:{twilson,hwerner,mpark}@kpf.com

||e-mail:cssilva@nyu.edu



Ferreira et al, 2015



NYU

TANDON SCHOOL
OF ENGINEERING



VISUALIZATION
IMAGING AND
DATA ANALYSIS
CENTER

Data management / UrbanGIS

URBANE:

A 3D Framework to Support Data Driven Decision Making in Urban Developments

IEEE VAST 2015
Submission ID: 268

Harish Doraiswamy¹, Nivan Ferreira¹, Huy Vo¹, Claudio Silva¹, Marcos
Lage², Muchan Park³, Heidi Werner³, Luc Wilson³

New York University¹, Universidade Federal Fluminense²,
Kohn Pederson Fox Associates PC³



NYU

TANDON SCHOOL
OF ENGINEERING

Data management / Image Collections

ARIES: Enabling Visual Exploration and Organization of Art Image Collections

Lhaylla Crissaff, Louisa Ruby, Samantha Deutch, Luke DuBois, Jean-Daniel Fekete, *Senior Member, IEEE*, Juliana Freire, *Member, IEEE*, Cláudio T. Silva, *Fellow, IEEE*



Fig. 2. The ARIES interface includes a toolbar (a) and four views: image menu (b), metadata (c), lightbox canvas (d) and group menu (e). Image menu, metadata and group menu are retractable, enlarging the lightbox canvas. Works of art on the lightbox canvas are displayed in relative size.



NYU

TANDON SCHOOL
OF ENGINEERING

Data management / Image Collections



NYU

TANDON SCHOOL
OF ENGINEERING

Data Management / Cloud

Live Demo: Statcast

The screenshot shows the MLB.com Tracking page for a game between the Colorado Rockies and Milwaukee Brewers on April 3, 2017. The page is divided into several sections:

- Scoreboard:** A table showing the scores of various MLB teams. The Rockies are listed as the home team with a score of 1, and the Brewers as the away team with a score of 4.
- Innings:** A list of the top 9 innings. The top 4 innings are highlighted, showing plays such as "Pitching Change: Tommy Milone replaces Junior Guerra, batting 9th." and "At Bat 25: Carlos Gonzalez called out on strikes."
- Plays:** A section titled "Top 4" showing the top 4 plays of the game. The top 4 plays are highlighted, showing plays such as "At Bat 25: Carlos Gonzalez called out on strikes." and "At Bat 26: Nolan Arenado singles on a fly ball to right fielder Domingo Santana."
- StatCast Metrics:** A section titled "StatCast Metrics" showing the bottom 3 of the game. It includes a "Watch Simulation" button and a "View Summary" button. Below this is a "New - Probabilities" section showing a 10% probability.
- Pitch Details:** A table showing the details of the pitch. The table has columns for Release Speed, Perceived Speed, Spin Rate, Type, Extension, Horizontal Break, and Vertical Break. The values are: Release Speed: 90.32 mph, Perceived Speed: 90.35 mph, Spin Rate: 1524.04 rpm, Type: SL, Extension: 5.33 ft, Horizontal Break: 4.71 in, Vertical Break: -10.54 in.
- Hit Details:** A table showing the details of the hit. The table has columns for Speed, Distance, Proj. Distance, Confidence, Angle, Height, Hang Time, Proj. Hang Time, Direction, and Generated Speed. The values are: Speed: 38.07 mph, Distance: 5.03 ft, Proj. Distance: 5.00 ft, Confidence: High, Angle: -46.71 deg, Height: 3.55 ft, Hang Time: 0.49 sec, Proj. Hang Time: 0.49 sec, Direction: 55.27 deg, Generated Speed: -52.45 mph.
- Calculated Metrics:** A section titled "Calculated Metrics" showing the calculated metrics for the hit.



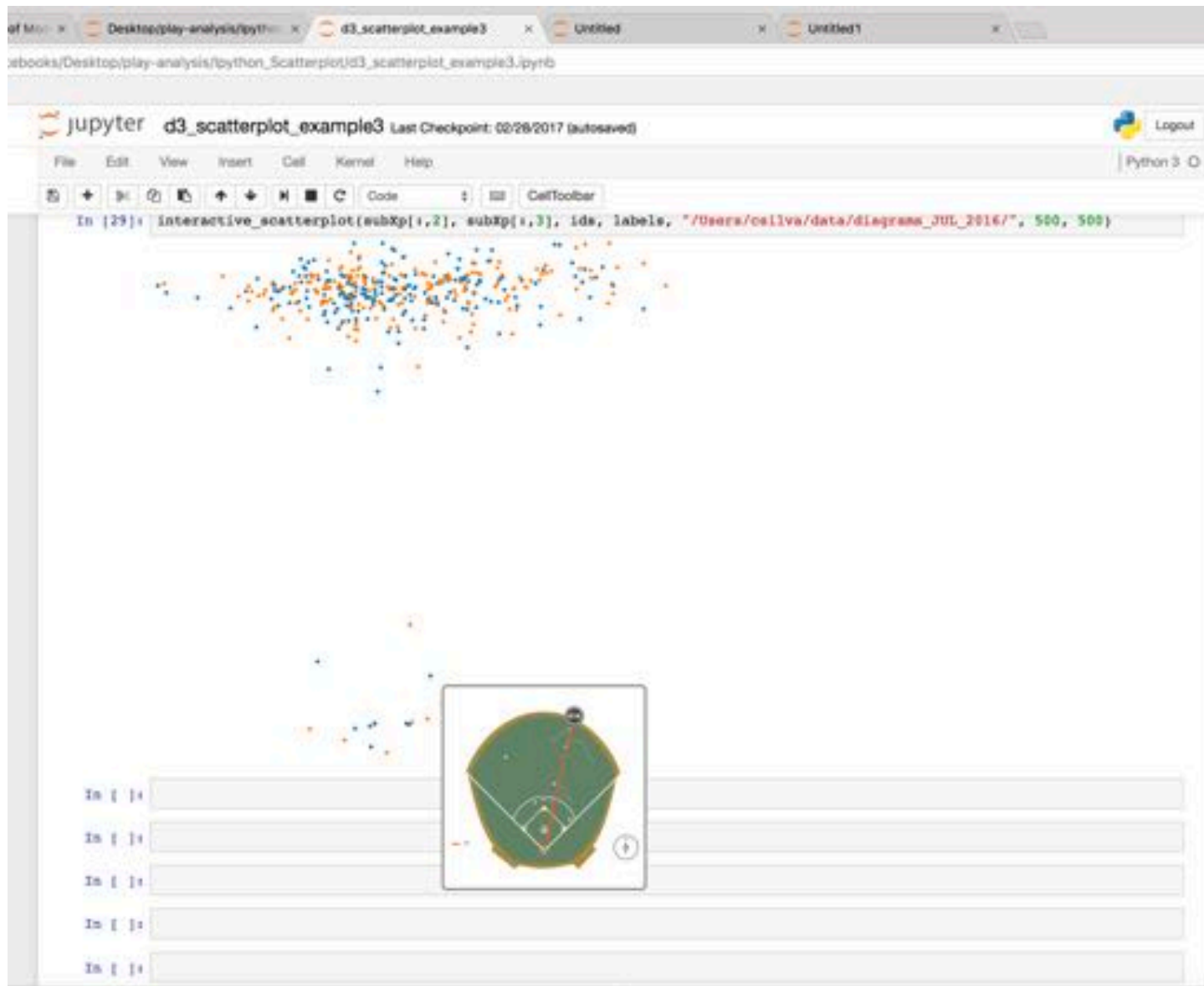
NYU

TANDON SCHOOL
OF ENGINEERING



VISUALIZATION
IMAGING AND
DATA ANALYSIS
CENTER

System support / Interactive Programming



NYU

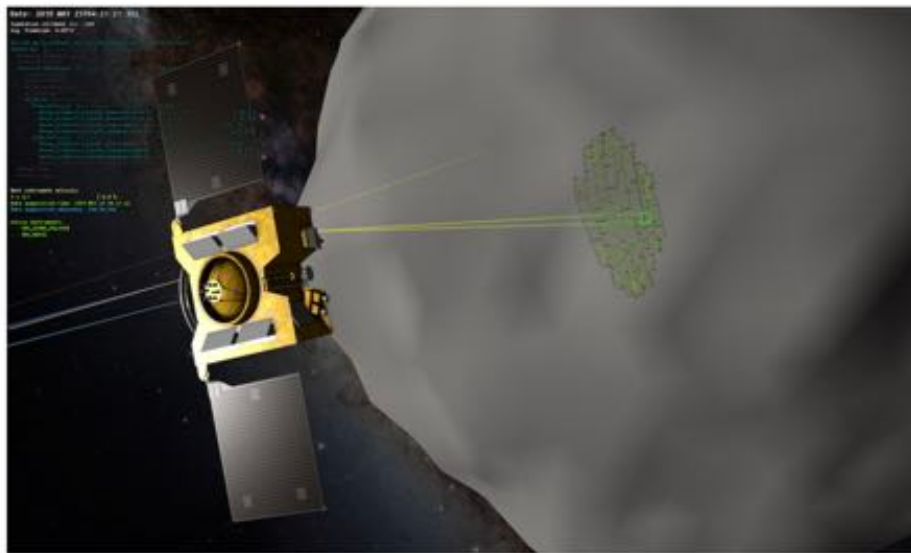
TANDON SCHOOL
OF ENGINEERING

OpenSpace is open source interactive data visualization software designed to visualize the entire known universe and portray our ongoing efforts to investigate the cosmos.

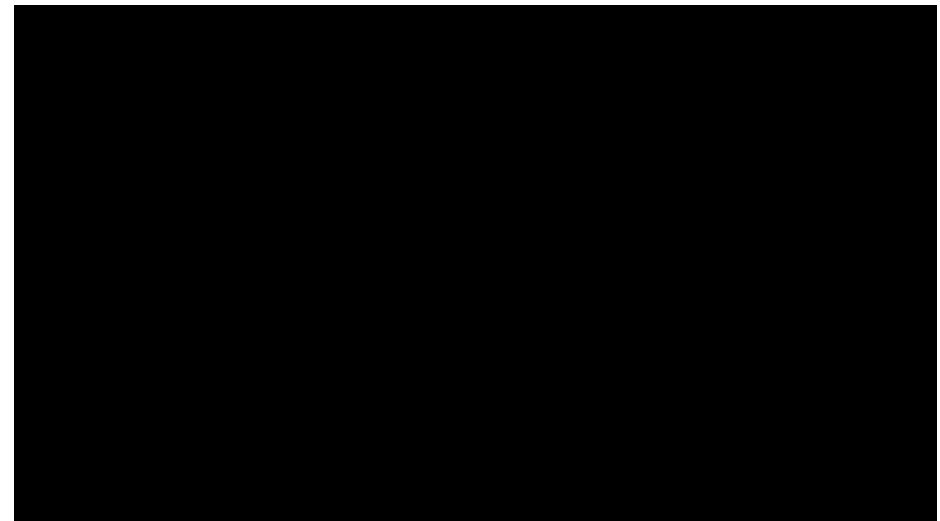
OpenSpace brings the latest techniques from data visualization research to the general public. OpenSpace supports interactive presentation of dynamic data from observations, simulations, and space mission planning and operations. OpenSpace works on multiple operating systems, with an extensible architecture powering high resolution tiled displays and planetarium domes, and makes use of the latest graphic card technologies for rapid data throughput. In addition, OpenSpace enables simultaneous connections across the globe, creating opportunity for shared experiences among audiences worldwide.

Osiris-REx Launch Event at AMNH

September 9, 2016 — Kayla Nussbaum — News



Today, NASA



NYU

TANDON SCHOOL
OF ENGINEERING



VISUALIZATION
IMAGING AND
DATA ANALYSIS
CENTER

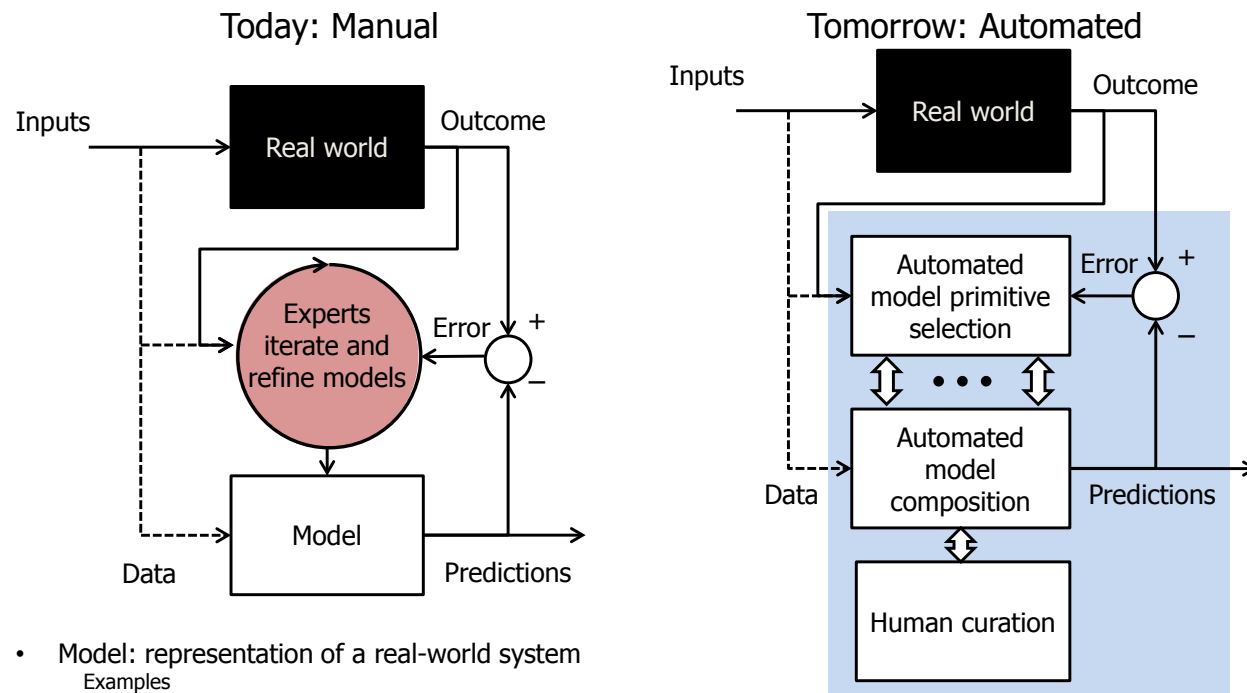
AI

DARPA D3M, Mr Wade Shen

<http://www.darpa.mil/program/data-driven-discovery-of-models>



D³M: Data-driven discovery of models



- Model: representation of a real-world system

Examples

- Inferring locations of images
- Prediction of election outcomes
- Estimation model for disease outbreaks
- Manual process: 10-1000s of person-years
- Teams of experts required to develop the model

- Automatically select problem-specific model primitives
 - Extend the library of modeling primitives
- Automatically compose complex models from primitives
- Facilitate user interaction with composed models

Approved for Public Release, Distribution Unlimited

3



NYU

TANDON SCHOOL
OF ENGINEERING

AI / Natural Language

FlowSense: A Natural Language Interface for Visual Data Exploration with Data Flow

Bowen Yu and Cláudio T. Silva *Fellow, IEEE*

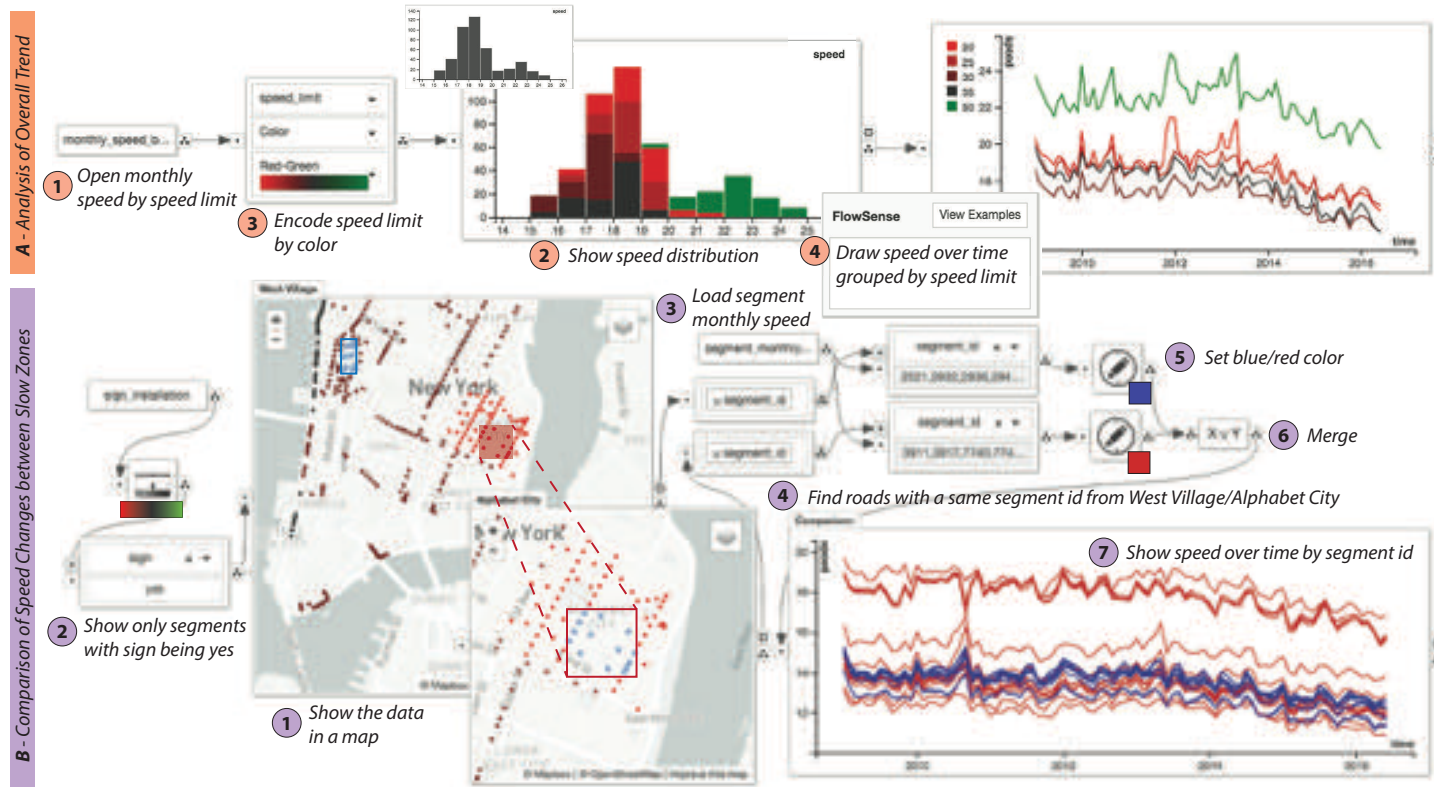


Fig. 1. Applying FlowSense to study the speed reduction in New York City. The important steps in the analysis and their natural language queries are shown in order. (A) FlowSense visualizes the overall speed reduction trend for streets of different speed limit. Step 4 shows the FlowSense dialog for typing queries. (B) A comparative study on the street speed changes between the West Village slow zone (blue) and the Alphabet City slow zone (red).



NYU

TANDON SCHOOL
OF ENGINEERING

AI / Natural Language

FlowSense: A Natural Language Interface for Visual Data Exploration with Data Flow

Bowen Yu, Claudio Silva
New York University

Submission ID: 267



NYU

TANDON SCHOOL
OF ENGINEERING

AI / Features Over 1000s of Datasets

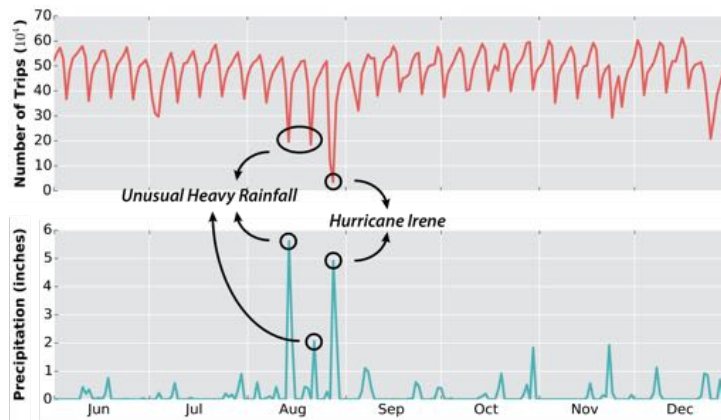


Fig. 2. Variation of the number of taxi trips in NYC over time (top) and its relationship with precipitation (bottom).

Chirigati et al, 2016

Chan et al, 2017

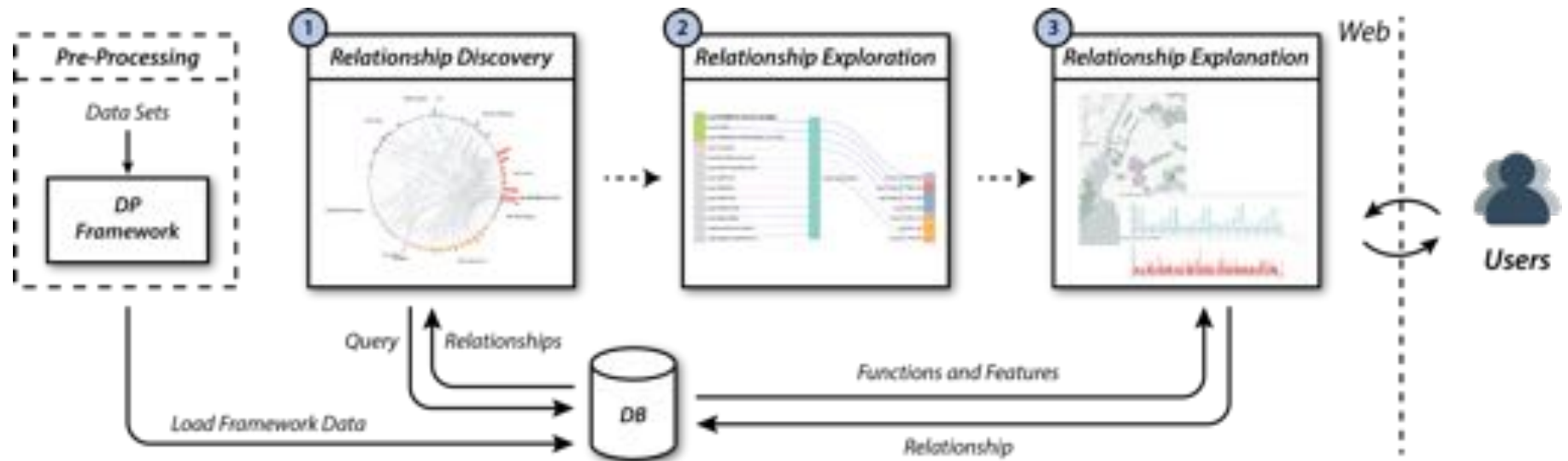


Fig. 4. DPER overview: the three components of the system correspond to the different stages of the data analysis pipeline. First, users query for interesting relationships (1); then, they can browse and filter the query results based on the relationships properties (2); finally, given a relationship, users can further inspect the data behind it to assess its validity (3).



NYU

TANDON SCHOOL
OF ENGINEERING

AI / Features Over 1000s of Datasets

DPer: A Deeper Dive into
Polygamous Relationships in Urban Data

Online Submission ID: 264



N

Major trends - Summary

Artificial Intelligence
Deep learning

Data management

User Interfaces

System architecture
Cloud environments
Interactive support/programming

FOCUS ON THE USER!

—

*DO NOT BE AFRAID TO GO OFF THE BEATEN PATH,
DEVELOP NEW SYSTEMS,
IDEAS ARE MORE IMPORTANT THAN CODE,
AND TECHNOLOGY MOVES VERY FAST*

—

*TRY NEW IDEAS!
RELEASE AS OPEN SOURCE!*



NYU

TANDON SCHOOL
OF ENGINEERING

Thank you!

csilva@nyu.edu

Urban Data Analysis @ NYU VIDA

Urbane: A 3D Framework to Support Data Driven Decision Making in Urban Development

A GPU-Based Index to Support Interactive Spatio-Temporal Queries over Historical Data

Harish Doraiswamy* Huy T. Vo† Cláudio Silva* Juliana Freire*
*New York University †The City College of the City University of New York
{harishd, huy.vo, csilva, juliana.freire}@nyu.edu

Topology-based Catalogue Exploration Framework for Identifying View-Enhanced Tower Designs

A Scalable Approach for Data-Driven Taxi Ride-Sharing Simulation

Masayo Ota^{1,2}, Huy Vo^{1,3}, Cláudio Silva^{1,2}, and Juliana Freire^{1,2}
¹Center for Urban Science and Progress, New York University
²Department of Computer Science and Engineering, New York University
³Department of Computer Science, the City College of the City University of New York
{masayo.ota, huy.vo, csilva, juliana.freire}@nyu.edu

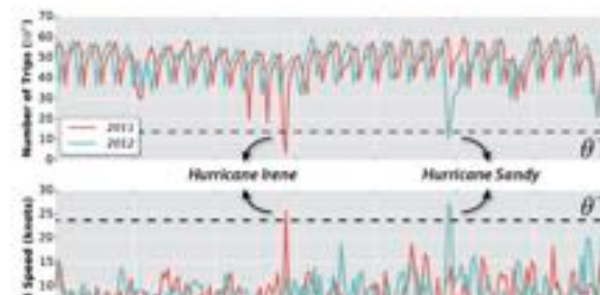
Data Polygamy: The Many-Many Relationships among Urban Spatio-Temporal Data Sets

Fernando Chirigati* Harish Doraiswamy* Theodoros Damoulas*† Juliana Freire*

* New York University * University of Warwick † Alan Turing Institute
{fchirigati, harishd, juliana.freire}@nyu.edu damoulas@warwick.ac.uk

ABSTRACT

The increasing ability to collect data from urban environments, coupled with a push towards openness by governments, has resulted in the availability of numerous spatio-temporal data sets covering diverse aspects of a city. Discovering relationships between these data sets can produce new insights by enabling domain experts to not only test but also generate hypotheses. However, discovering these relationships is difficult. First, a relationship between two data sets may occur only at certain locations and/or time periods. Second, the sheer number and size of the data sets,



Collaborations with social scientists, architects, and city agencies