



UNIVERSITY *of* WASHINGTON

FAT Databases

Bill Howe

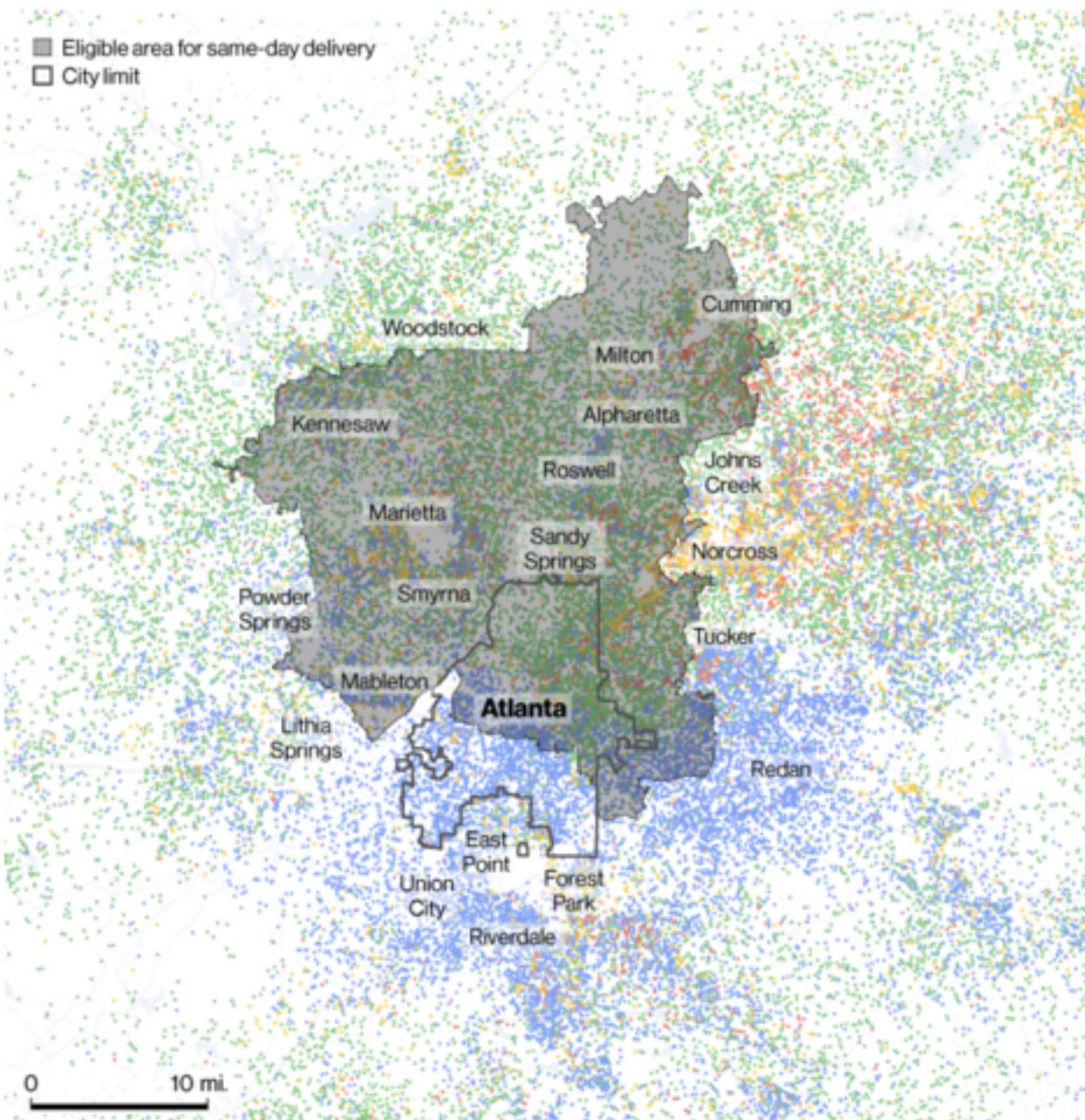
University of Washington

The next ~10 minutes...

- Some more examples of the problems
- Some topics for DB research

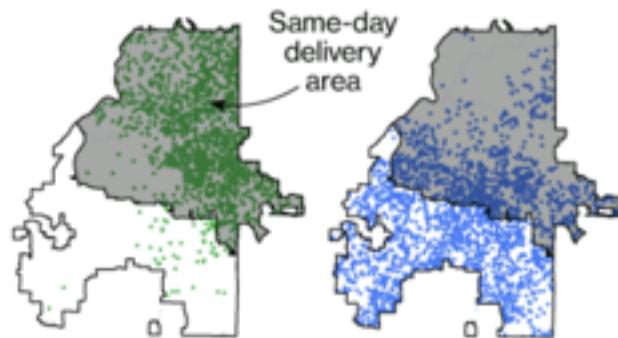
Amazon Prime Now Delivery Area: Atlanta

Bloomberg, 2016

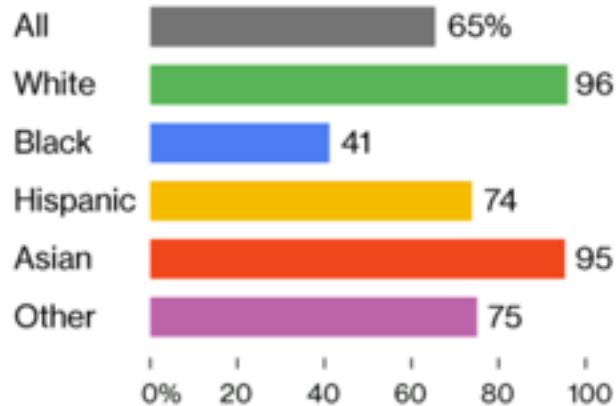


The northern half of Atlanta, home to 96% of the city's white residents, has same-day delivery. The southern half, where 90% of the residents are black, is excluded.

White residents Black residents



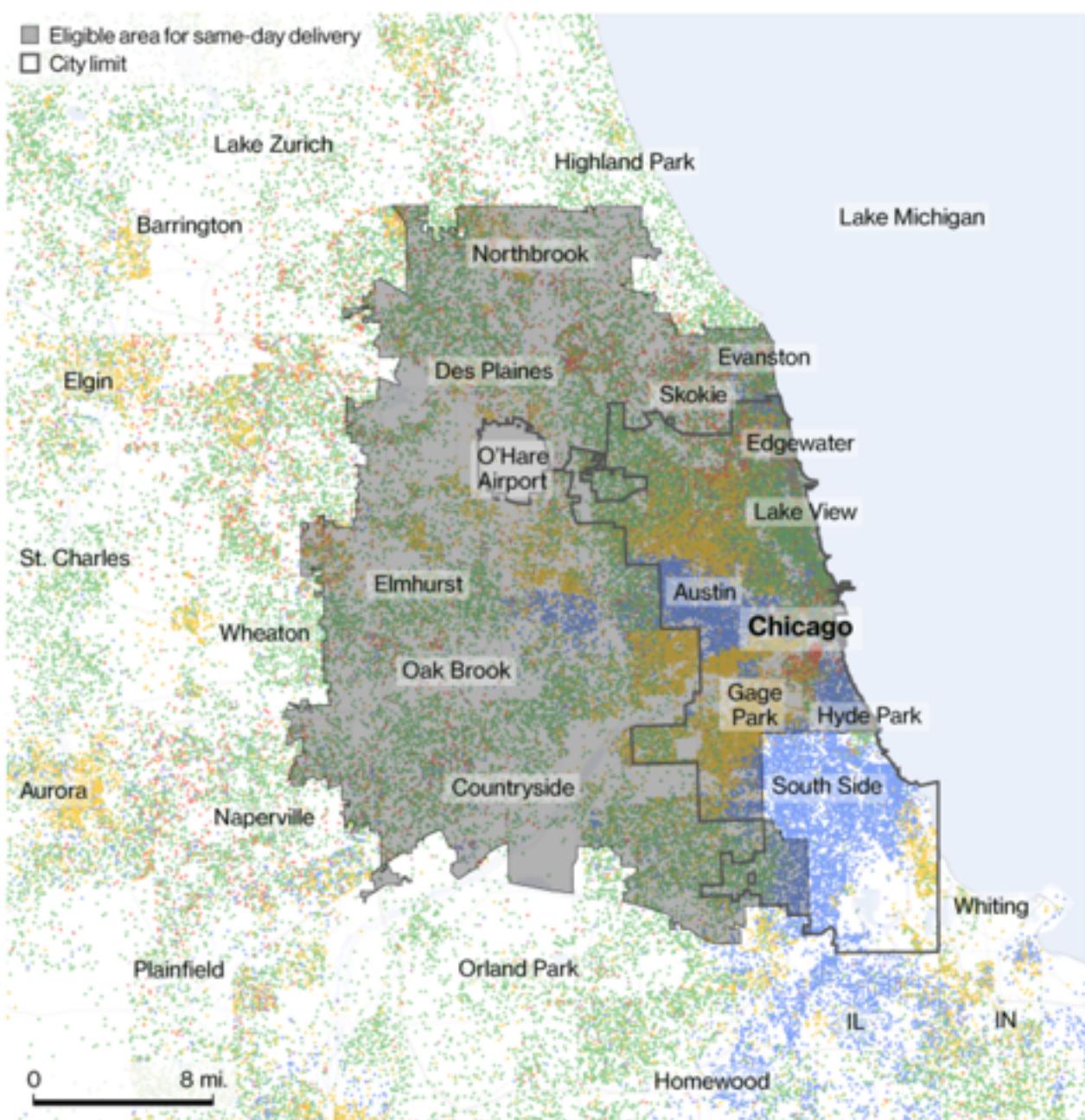
Percentage of residents living in ZIP codes with same-day delivery



Population percentages are based on American Community Survey estimates and have a 90% confidence interval.

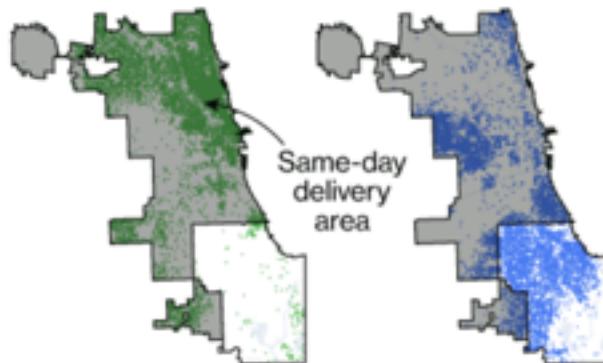
Amazon Prime Now Delivery Area: Chicago

Bloomberg, 2016

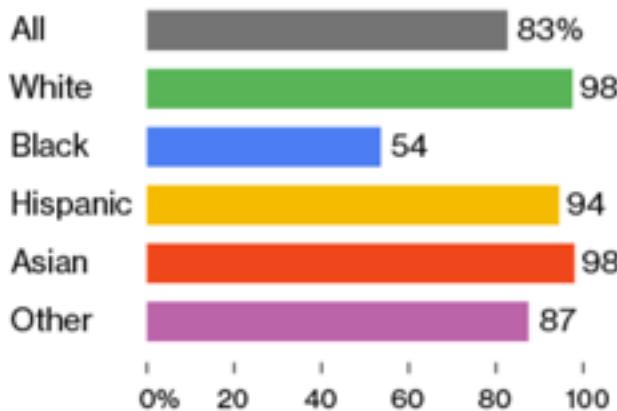


About half of Chicago's black residents live in the southern half of the city where they do not have access to Amazon's same-day delivery service.

White residents Black residents



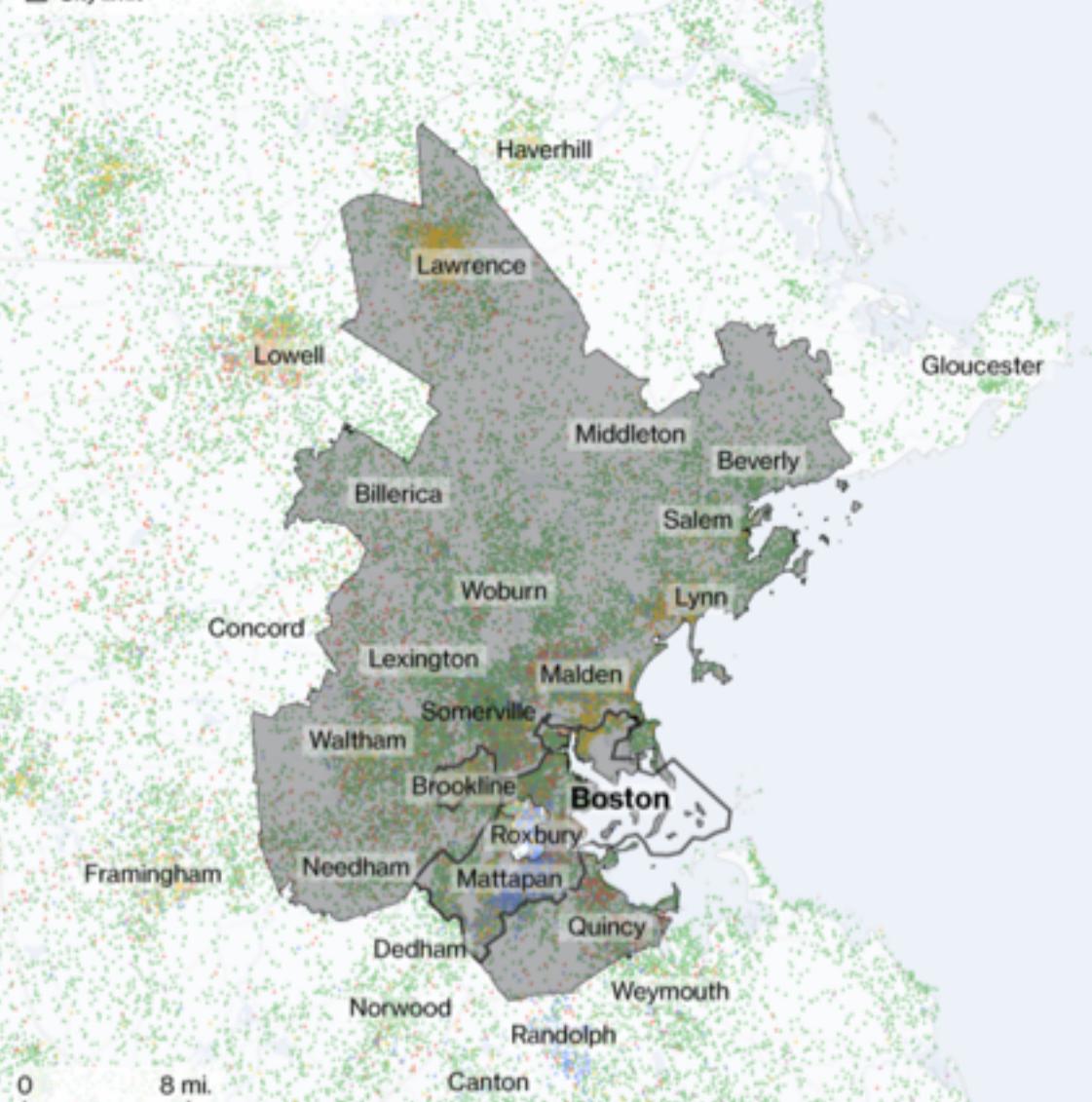
Percentage of residents living in ZIP codes with same-day delivery



Amazon Prime Now Delivery Area: Boston

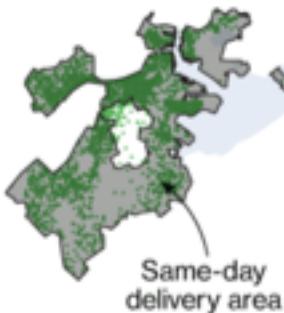
Bloomberg, 2016

- Eligible area for same-day delivery
- City limit

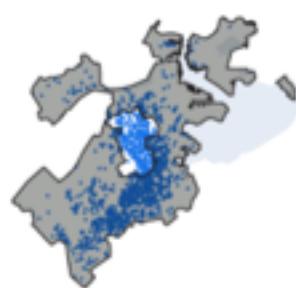


Three ZIP codes in the center of Boston, including the Roxbury neighborhood, are excluded from same-day coverage.

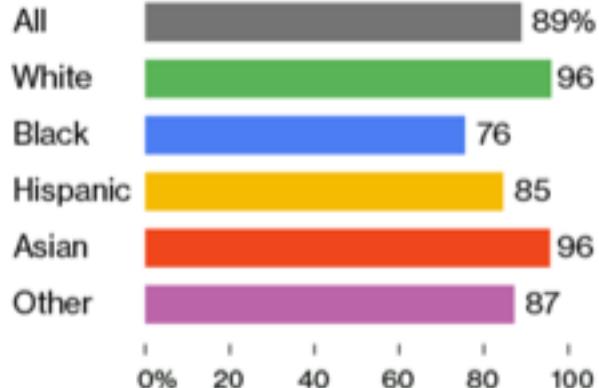
White residents



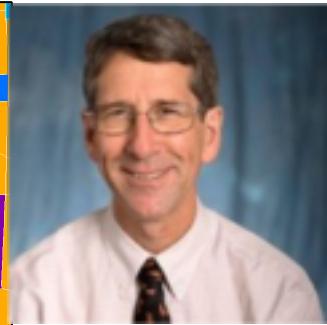
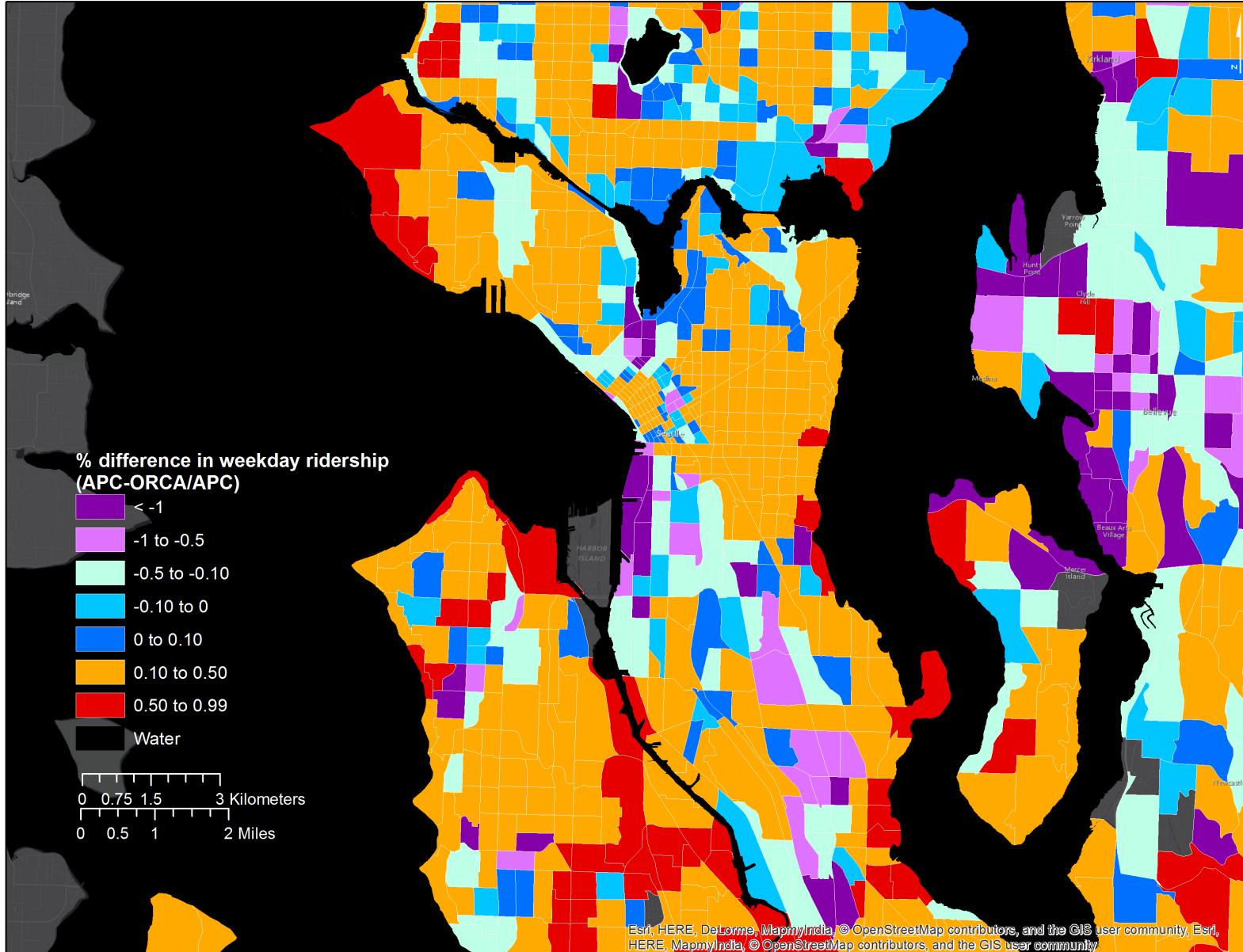
Black residents



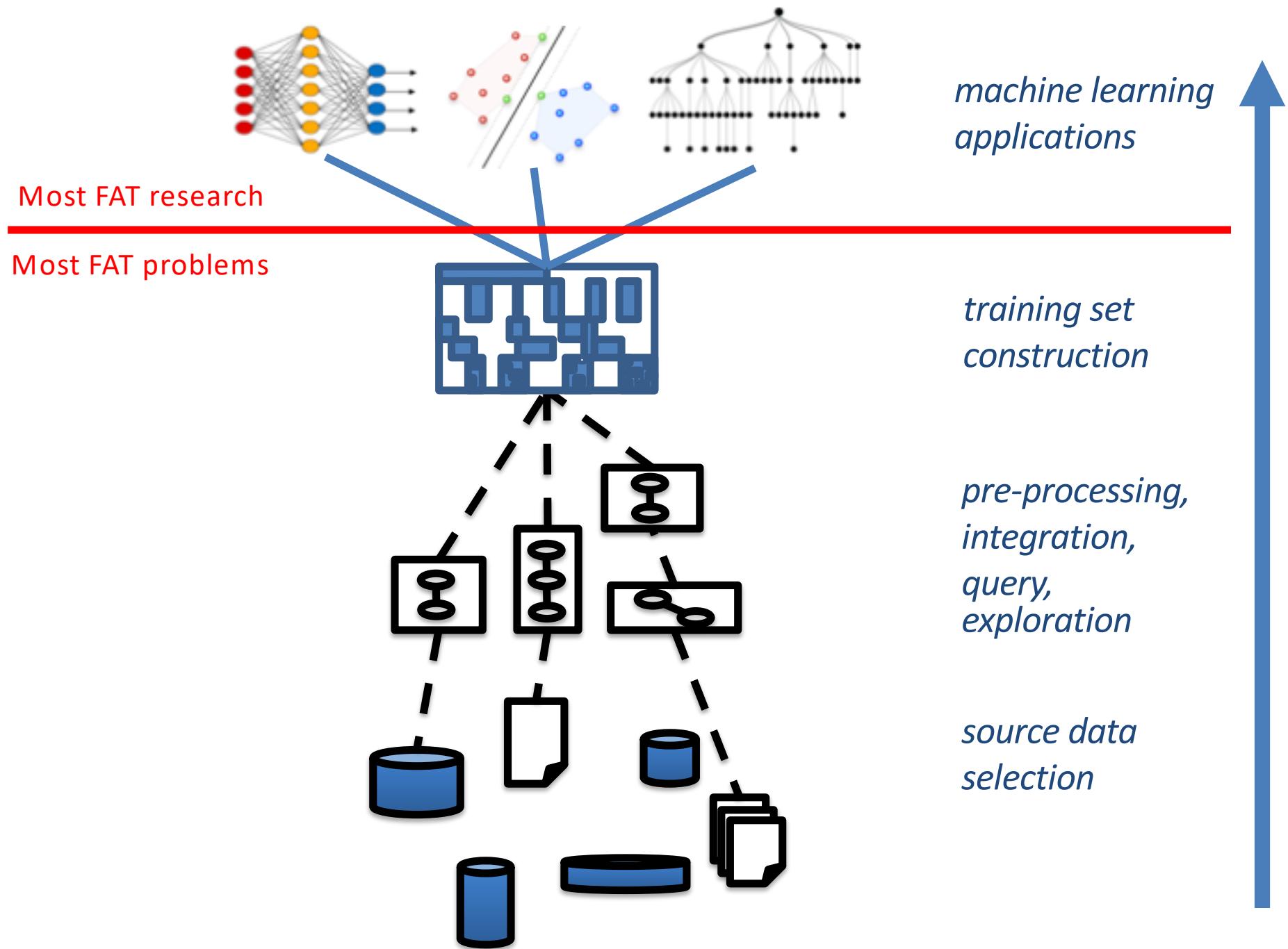
Percentage of residents living in ZIP codes with same-day delivery



Data source selection: Bias in transportation measurements



Mark
Hallenbeck
TRAC



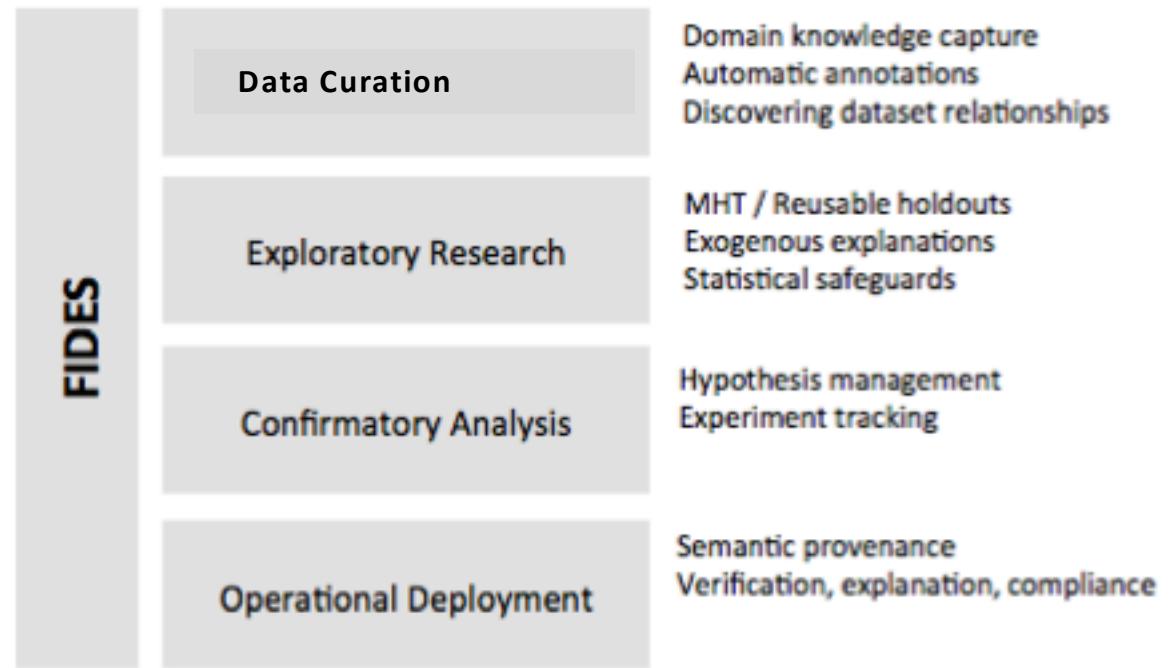
Databases are becoming “training set management systems”

- Claim: Query results are increasingly being used to train models
- So the exact answer to the query result is not that important. It's important that the trained model gets the right answer on unseen data
- We need declarative specification (i.e., SQL) and management of **high-quality training sets**
- What's a high-quality training set?
 - It's a bad training set if the resulting classifier doesn't work, or overfits
 - It's a bad training set if it deviates too far from the specification (fidelity)
 - It's a bad training set if it's too small (significance), and it *might* be a bad training set if it's too big (scale)
 - It's a bad training set if it leaks private information (privacy-preserving)
 - It's a bad training set if you can't tell where it came from (provenance)
 - It's a bad training set if it reinforces discrimination (bias-correcting)

So what do we do about it?

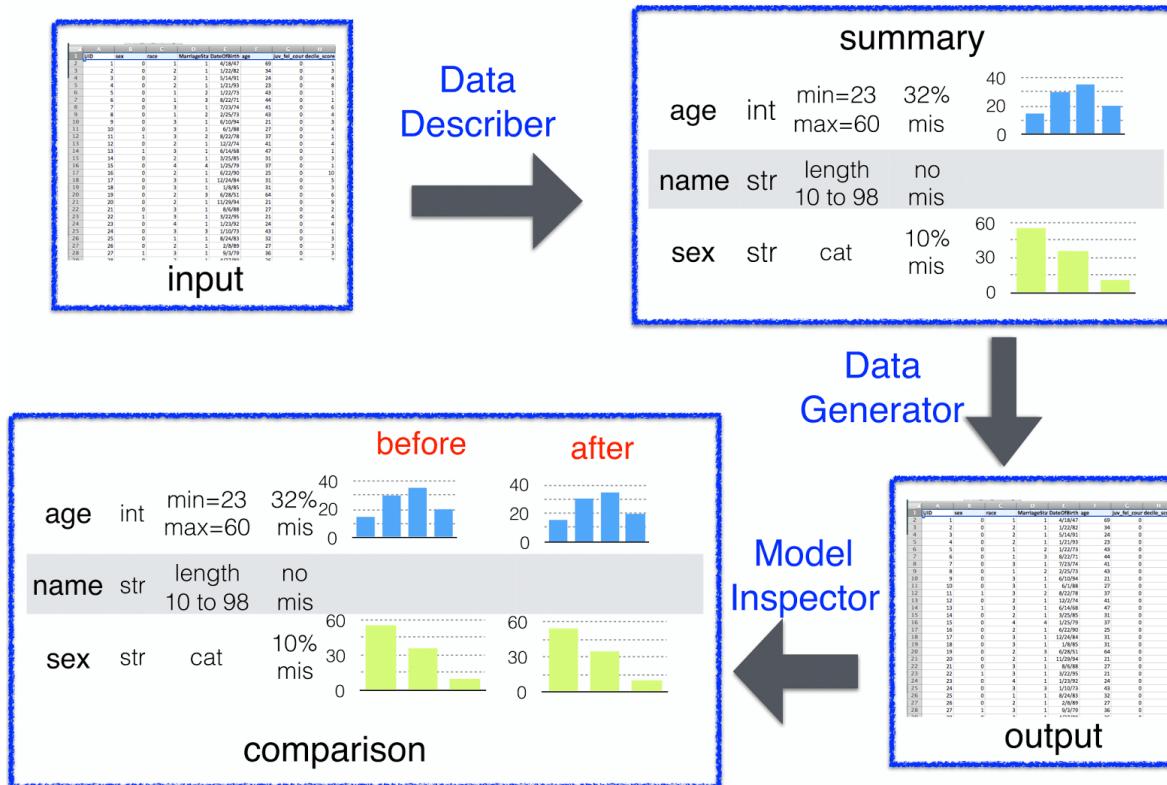
Fides: Responsible Data Management

Fairness
Accountability
Transparency
Privacy
Reproducibility



joint with Stoyanovich [US], Abiteboul [FR], Miklau [US], Sahuguet [US], Weikum [DE]

Data Synthesizer: Privacy-preserving synthetic data



With Stoyanovich (Drexel), Gee (Chicago), Ping (Drexel), Herman (UW)

A Nutritional Label for Rankings

[Yang, Stoyanovich, Asudeh, Howe, Jagadish, Miklau SIGMOD 2018 demo]

Recipe			
Top 10:			
Attribute	Maximum	Median	Minimum
PubCount	18.3	9.6	6.2
Faculty	122	52.5	45
GRE	800.0	796.3	771.9
Overall:			
Attribute	Maximum	Median	Minimum
PubCount	18.3	2.9	1.4
Faculty	122	32.0	14
GRE	800.0	790.0	757.8

Stability	
Score	
950	Stability
900	
850	
800	
750	
0	Rank Position
5	
10	
15	
20	
25	
30	
35	
40	
45	
50	

Slope at top-10: -6.91, over-all: -1.61.
A ranking is unstable when the absolute value of the slope of the line that is fit to the score distribution falls below 0.25.

Ranking Facts

← Recipe	
Attribute	Weight
PubCount	1.0
Faculty	1.0
GRE	1.0

Ingredients	
Attribute	Importance
PubCount	1.0
CSRankingAllArea	0.24
Faculty	0.12

Importance of an attribute in a ranking is quantified by the correlation coefficient between attribute values and items scores, computed by a linear regression model. Importance is high if the absolute value of the correlation coefficient is over 0.75, medium if this value falls between 0.25 and 0.75, and low otherwise.

Diversity at top-10

DeptSizeBin = Regional Code =

Large Small NE W MW SA SC

Diversity overall

DeptSizeBin = Regional Code =

Large Small NE W MW SA SC

Stability

Top-K	Stability
Top-10	Stable
Overall	Stable

Fairness

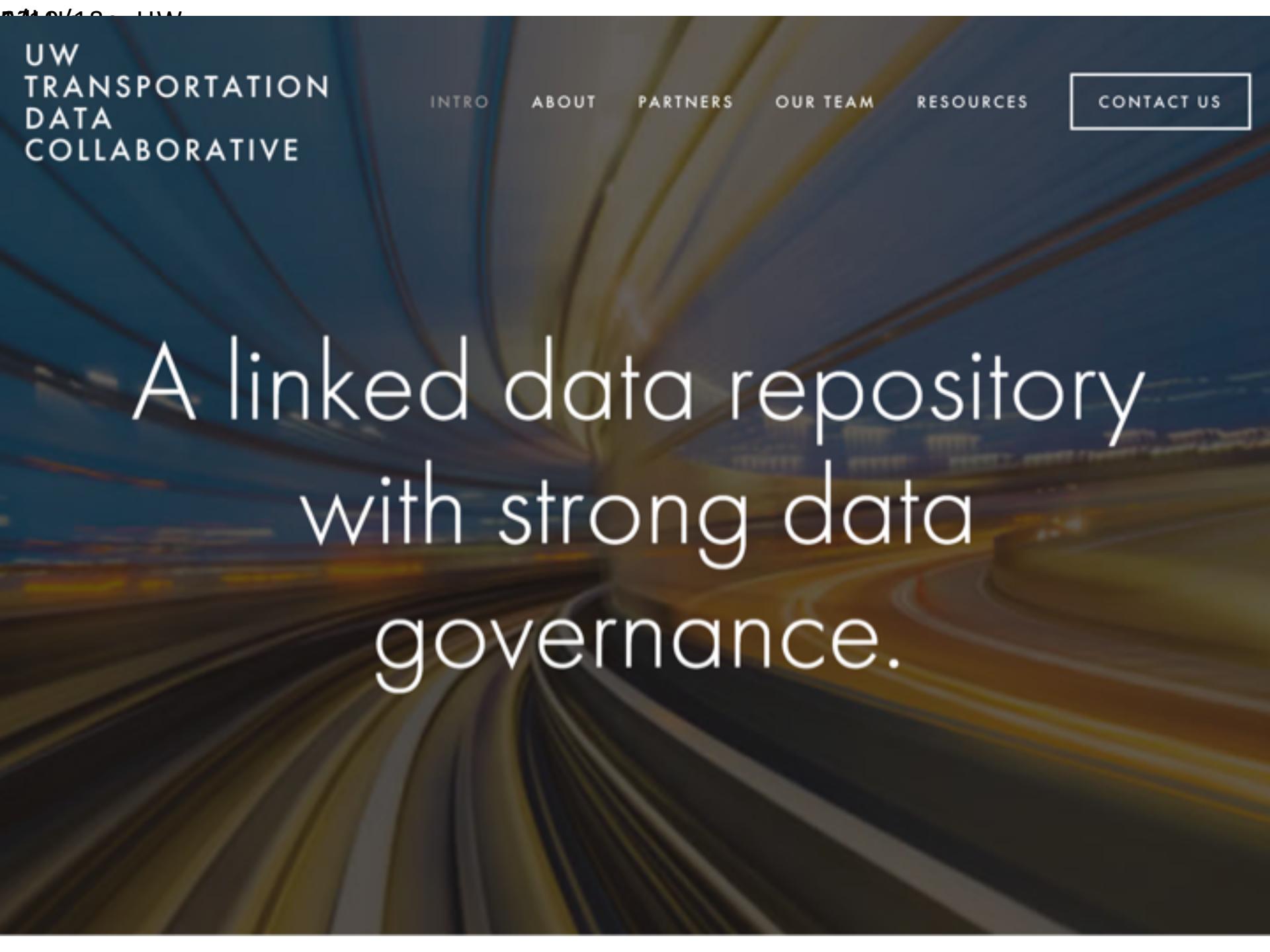
DeptSizeBin	FA'IR	Pairwise	Proportion
Large	Fair	Fair	Fair
Small	Unfair	Unfair	Unfair

FA'IR and difference in proportions (Proportion) are measured with respect to 26 highest-scoring items (the top-K). The top-K contains 100 items or one half of the input, whichever is smaller.

← Ingredients			
Top 10:			
Attribute	Maximum	Median	Minimum
PubCount	18.3	9.6	6.2
CSRankingAllArea	13	6.5	1
Faculty	122	52.5	45
Overall:			
Attribute	Maximum	Median	Minimum
PubCount	18.3	2.9	1.4
CSRankingAllArea	48	26.0	1
Faculty	122	32.0	14

← Fairness	
FA'IR	Pairwise
p-value	adjusted p-value
DeptSizeBin	p-value

A ranking is considered unfair when the p-value of the corresponding statistical test falls below 0.05.



A linked data repository
with strong data
governance.

THE BIKE SHARE WAR IS SHARING UP SEATTLE LIKE NOWHERE ELSE

Residents are divided over whether the city's dockless bike share program is revolutionizing transit—or creating an unwieldy, dangerous mess.

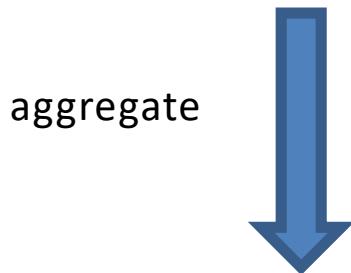
BY [MARK HARRIS](#)



Example: Bike Share data

- Companies need to release trip data to comply with Seattle permits and civic transparency
- But there are concerns about privacy, misuse, and competitive advantage
- Setup:

Trip(time, userid, company, orig, dest, gender, helmet)



Domain info:

company in {Lime, Spin, Ofo}

origin, dest one of 94 neighborhoods in Seattle}

gender in {M, F, other, null}

Helmetuser in {true, false}

OD(company, origin, dest, gender, helmet, count)

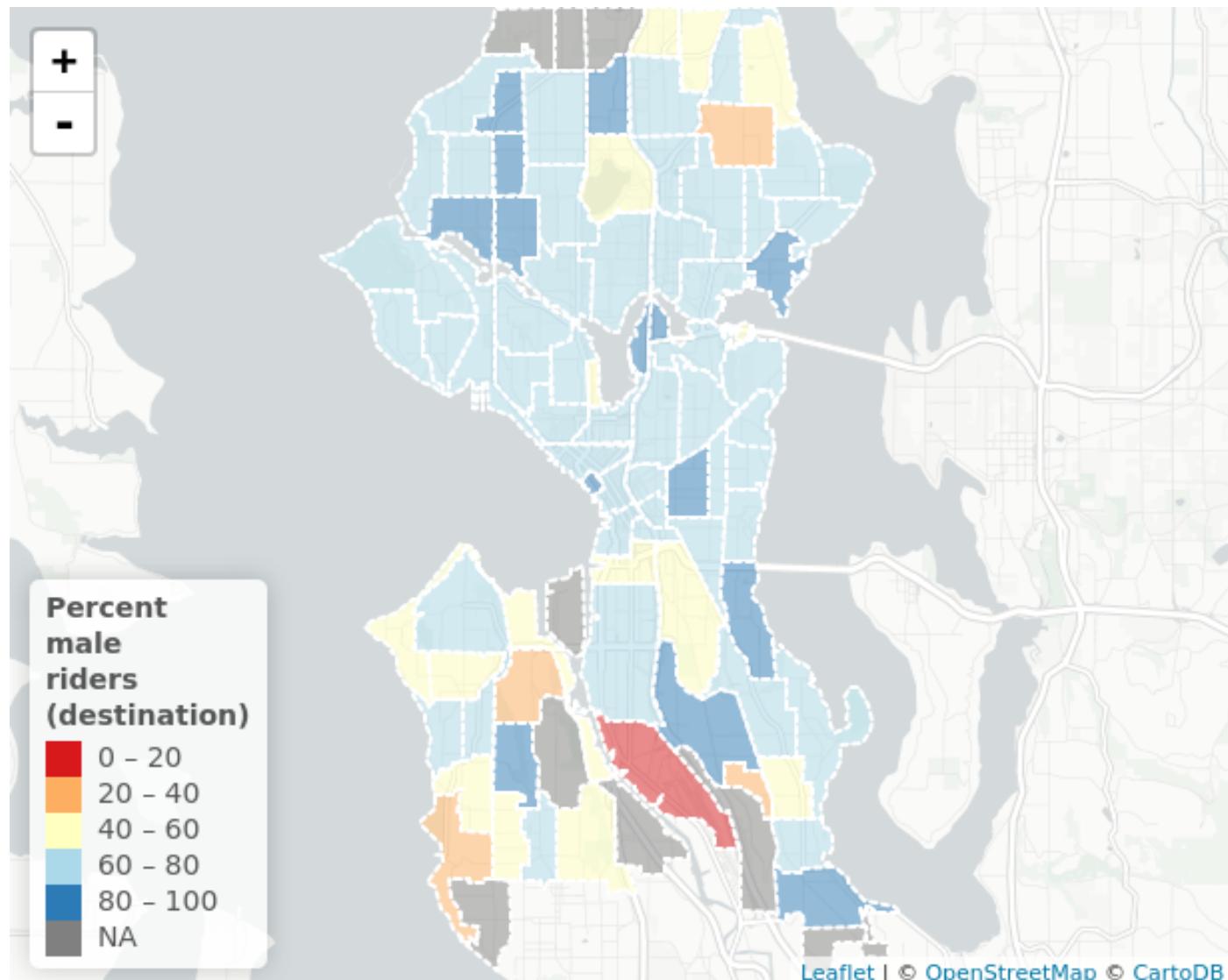


We can release the joint distribution of company, origin, dest, gender, helmet

To release
plots like this:

For privacy we can add
noise to the counts.

But we also want to
remove bias...



Bias-Corrected Data Sharing by Breaking Causal Relationships



Babak Salimi

Luke Rodriguez

Ex: We don't want race to influence hiring, so set the mutual information between, say, race and GPA to zero before releasing the dataset.

Different strategies:

- You could remove and insert tuples
- You could directly edit the GPA
- You could change the “weight” of each tuple

Bias-Corrected Data Sharing by Breaking Causal Relationships



Babak Salimi



Luke Rodriguez

Back to bike share: Hiding competitive advantage

Company doesn't mind releasing data, but doesn't want to reveal that they have a marketing campaign targeting women.

Set the mutual information between company (X) and gender (Y) to zero, conditioned on the other attributes (Z).

Compute a new joint distribution of trips, asserting independence between X and Y conditioned on Z

$$P_{R'}(\mathbf{A}) = P_{R'}(X|\mathbf{Z})P_{R'}(Y|\mathbf{Z})P_{R'}(\mathbf{U}|XYZ)$$

FAT* 2019 Call for Papers

Important Dates and Links

Submission site	TBD
Abstract Deadline	pre-registration at 11:59PM August 16, 2018 AoE
Full paper submission	11:59PM August 23, 2018 AoE
Notification Date	October 12, 2018
Conference Date	late January/early February 2019

FAT* is an international and interdisciplinary peer-reviewed conference that seeks to publish and present work examining the **fairness, accountability, and transparency of algorithmic systems**.

Topics of Interest

The FAT* conference solicits work from a wide variety of disciplines, including computer science, statistics, the humanities, and law. FAT* welcomes submissions that touch on any of the following topics (broadly construed):

- Fairness
 - Techniques and models for fairness-aware data mining, information retrieval, recommendation, etc.
 - Formalizations of fairness, bias, discrimination; trade-offs and relationships between them
 - Defining, measuring and mitigating biases in data sets; improving data collection processes; combining different sources of information
 - Translation of legal, social, and philosophical models of fairness into mathematical objectives
 - Qualitative, quantitative, and experimental studies on perceptions of algorithmic bias and unfairness
 - Design interventions to mitigate biases in systems, or discourage biased behavior from users
 - Measurement and data collection regarding potential unfairness in systems
 - Understanding how tools from causal inference can help us to better reason about fairness and the interplay between prediction and intervention
 - Analyses of the impact of algorithmic experimentation and exploration