

MOBILE SOCIAL MEDIA AND DEMOGRAPHICS OPPORTUNITIES AND RISKS

A. Chaintreau, Computer Science Department

OUTLINE

- How I came to care about demographics in big data
 - And why we need new tools for mobile social media
- Risks: Identity Reconciliation using location data
- Opportunities: Demographic study of Urban mobility
- How to scale?
- Conclusion

BIG DATA TO SOCIAL SCIENCE: WE NEED YOU!

Detailed Targeting **INCLUDE** people who match at least ONE of the following

- Behaviors > Residential profiles
 - Likely to move
- Interests > Additional Interests
 - Buying a House
 - First-time buyer
 - House Hunting

Add demograph

Narrow Audience

EXCLUDE people

- Demographics
 - African An
 - Asian Ame
 - Hispanic (I

Add demograph

Satterfield added that the “Ethnic Affinity” is not the same as race — which Facebook does not ask its members about. Facebook assigns members an “Ethnic Affinity” based on pages and posts they have liked or engaged with on Facebook.

When we asked why “Ethnic Affinity” was included in the “Demographics” category of its ad-targeting tool if it’s not a representation of demographics, Facebook responded that it plans to move “Ethnic Affinity” to another section.

Don't Miss: The Breakdown Terror in Little Saigon Dollars for Docs Surgeon Scorecard Red Cross Workers' Comp [Donate](#)

PRO PUBLICA Journalism in the Public Interest

Receive our top stories daily
Email address [SUBSCRIBE](#)

Home Investigations Data MuckReads Get Involved About Us

f t Search ProPublica

Electionland
ProPublica is covering access to the ballot and

LATEST UPDATES:
Law Enforcement Presence at Florida Election Spot After ProPublica Report of Hecklers
Nov. 3rd, 8:14 AM EDT
Twitter Takes Down False Claims Clinton Supporters Can Vote From Home, But There Are Many More



Four members of the Congressional Black Caucus wrote Facebook's Mark Zuckerberg telling him that the company should stop allowing advertisers to exclude people by race.



We're investigating algorithmic injustice and

[see all](#)

HOW I CAME TO CARE ABOUT DEMOGRAPHICS

- Personalization runs the gamut of moral hazards
 - **Commerce** (Wash. Post, Oct. 31st):
Data exchanged in “sharing economy” exacerbates discrimination
 - **Information** (Propublica, Oct. 28th):
Race induced personalization on Ads
 - **Politics** (Wash. Post, Oct. 11th):
FB, Twitter, Instagram used by law enforcement



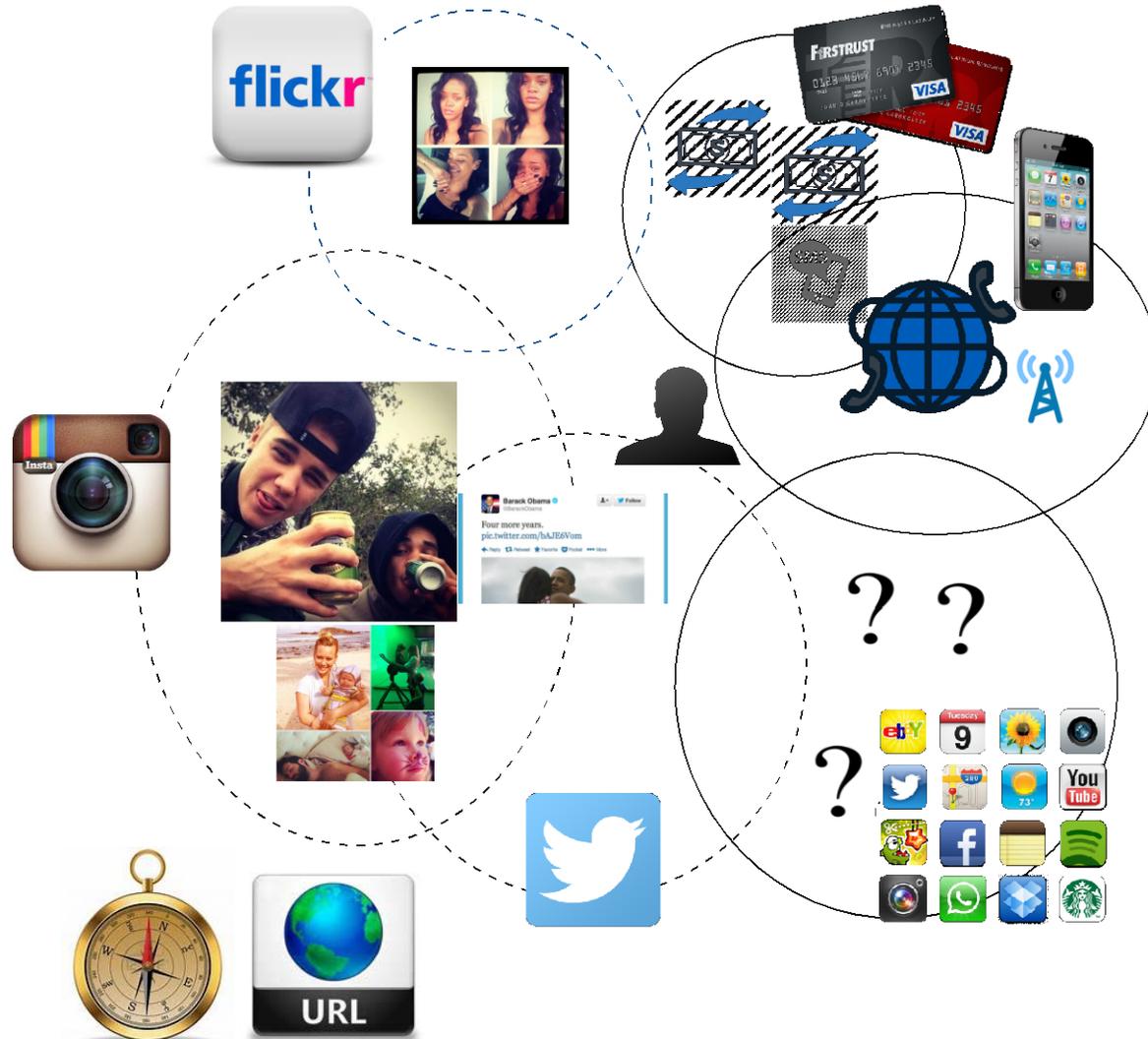
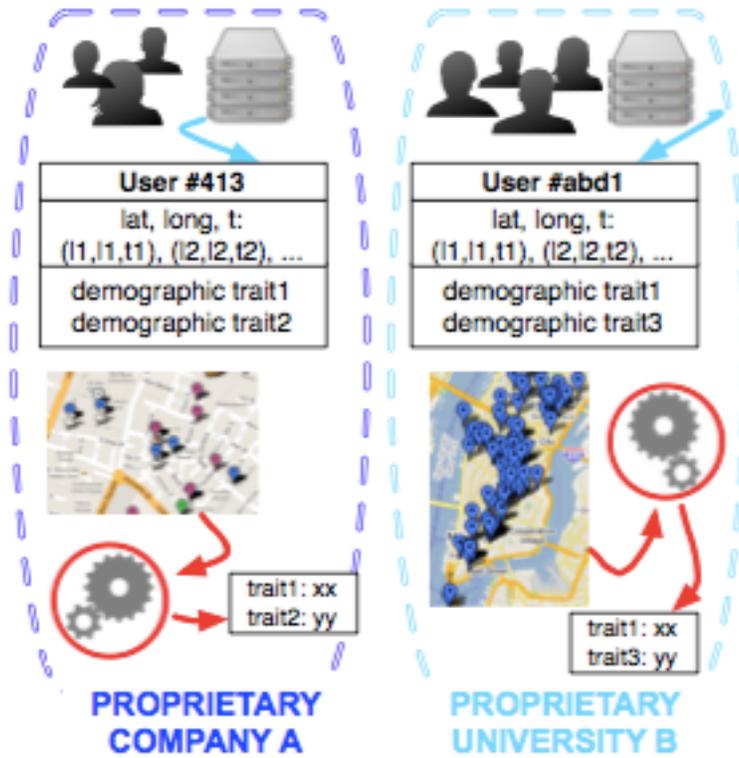
MY RESEARCH IN A NUTSHELL

- Reconciling Big Data's Disparate Impact with our values
 1. Building an infrastructure for transparency
 2. Reconciling data to empower users and social scientists
 3. When does social sharing reinforce disparate impact
- Promoting scalable reusable tools, incrementally deployable, validated with reproducible data in open access

OUTLINE

- How I came to care about demographics in big data
 - And why we need new tools for mobile social media
- Risks: Identity Reconciliation using location data
- Opportunities: Demographic study of Urban mobility
- How to scale?
- Conclusion

MOBILE SOCIAL MEDIA: BRIDGING THE CHASM



WHY WE NEED RECONCILIATION TOOLS?

- Privacy typically centers on anonymity in a single-domain
- Big Data has fragmented view (e.g., tweets, calls, transactions) with overlaps in various domains
- “Can two profiles I maintained be linked?”, “Can my behavior in be used to discriminate?”

Anonymization of Location Data Does Not Work: A Large-Scale Measurement Study

Hui Zang
Sprint
1 Adrian Ct
Burlingame, CA 94010, USA

Jean Bolot*
Technicolor
735 Emerson St
Palo Alto, CA 94301, USA



Unique in the Crowd: The privacy bounds of human mobility

SUI
/
APPLIED
COMPUTATI

the abstract numerical magnitude, independently of non-numerical cues. Our results indicate that a disposition to map numerical magnitudes onto a left-to-right-oriented MNL exists independently of cultural factors and can be observed in animals with very little nonsymbolic numerical experience, supporting a nativistic foundation of such orientation. Spatial mapping of numbers from left to right may be a universal cognitive strategy available soon after birth. Experience and, in humans, culture and education (e.g., reading habits and formal mathematics education) may modulate or even be modulated by this innate number sense.

IDENTITY AND PRIVACY

Unique in the shopping mall: On the reidentifiability of credit card metadata

Yves-Alexandre de Montjoye,^{1*} Laura Radaelli,² Vivek Kumar Singh,^{1,2} Alex “Sandy” Pentland¹

Large-scale data sets of human behavior have the potential to fundamentally transform the way we fight diseases, design cities, or perform research. Metadata, however, contain sensitive information. Understanding the privacy of these data sets is key to their broad use and, ultimately, their impact. We study 3 months of credit card records for 1.1 million people and show that four spatiotemporal points are enough to uniquely reidentify 90% of individuals. We show that knowing the price of a transaction increases the risk of reidentification by 22%, on average. Finally, we show that even data sets that provide coarse information at any or all of the dimensions provide little anonymity and that women are more reidentifiable than men in credit card metadata.

Large-scale data sets of human behavior have the potential to fundamentally transform the way we fight diseases, design cities, or perform research. Ubiquitous technologies create personal metadata on a vast scale.

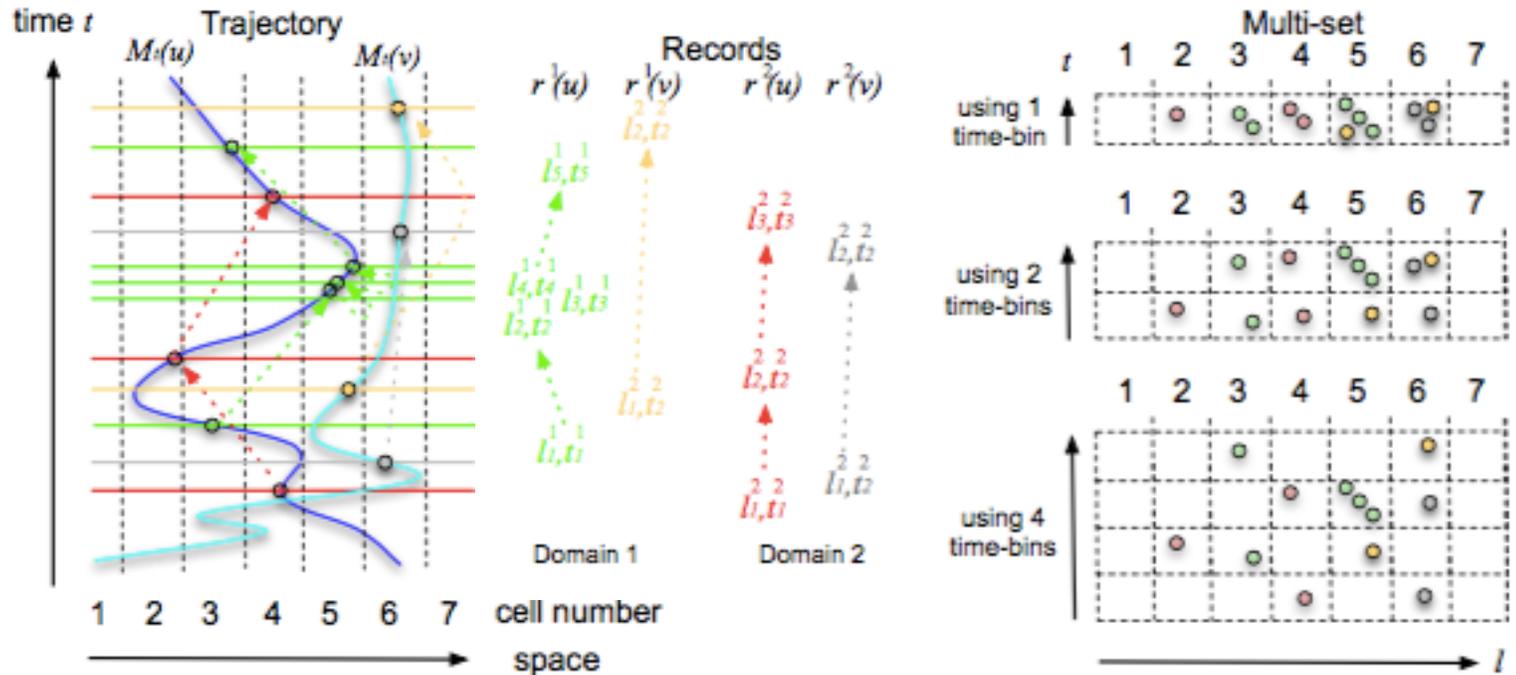
scale behavioral data sets to the invention of the microscope (1). New fields such as computational social science (2–4) rely on metadata to address crucial questions such as fighting malaria, studying the spread of information, or monitoring the

OUTLINE

- How I came to care about demographics in big data
 - And why we need new tools for mobile social media
- Risks: Identity Reconciliation using location data
- Opportunities: Demographic study of Urban mobility
- How to scale?
- Conclusion

IDENTITY RECONCILIATION USING LOCATION

- Using (space, time)
 - No username, friends, tastes, etc.
- \neq indiv. anonymity
 - Non-overlapping set
 - cross domain
- Difficult & lack data



HOW TO COMBINE POSITIVE AND NEGATIVE SIGNALS

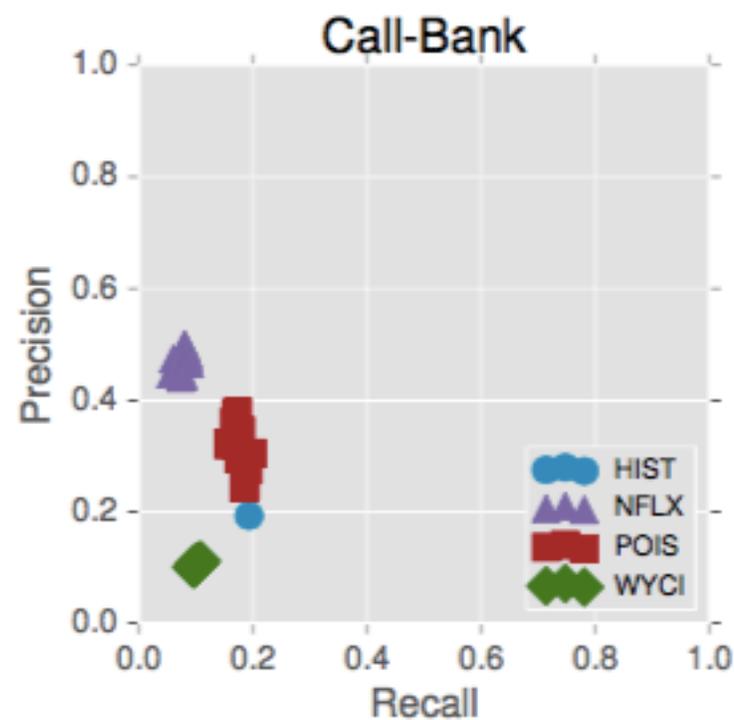
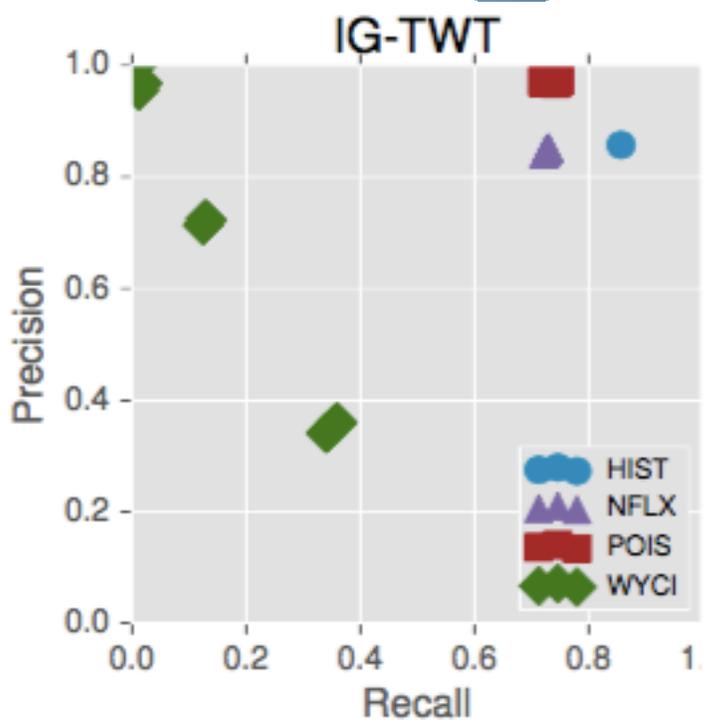
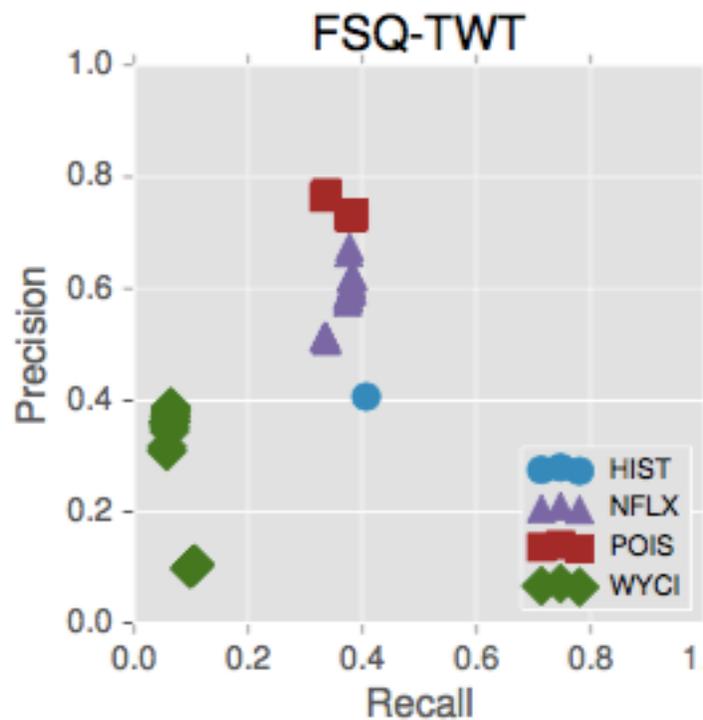
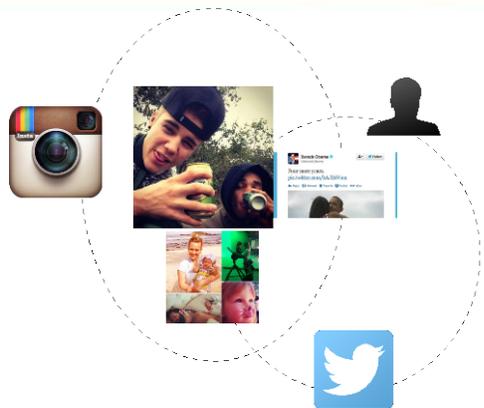
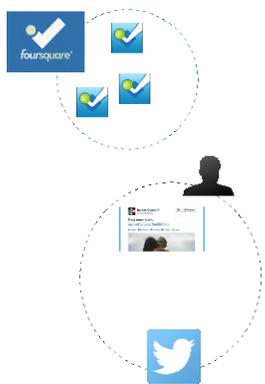
- 2 Poisson Processes sampled
 - From the same original pattern
 - 3 parameters: λ , p_1 , p_2
- Our analysis:
 - Computer a max matching.
 - Thm: In expectation, there is only one score maximizer, which is the correct identity reconciliation
 - Parameters? can be estimated from imperfect matching

Algorithm 1: Our reconciliation algorithm

Require: $\forall u \in U : r^1(u), \forall v \in V : r^2(v), \{\lambda_{\ell,t}\}$
for $(u, v) \in (U \times V)$ **do**
 $w(u, v) = \sum_{t \in T} \sum_{\ell \in L} \ln \phi_{\ell,t}(a_1(u, \ell, t), a_2(v, \ell, t))$
end for
Let $E = \{w(u, v) : (u, v) \in (U \times V)\}$
Compute the maximum weighted matching on the bipartite graph $B(U, V, E)$
return the function that maps matched vertices.

$$\frac{e^{-\lambda(1-p_1-p_2)}}{(\lambda(1-p_1))^{a_1}(\lambda(1-p_2))^{a_2}} \sum_{k \geq \max(a_1, a_2)} \frac{(\lambda(1-p_1)(1-p_2))^k k!}{(k-a_1)!(k-a_2)!}$$

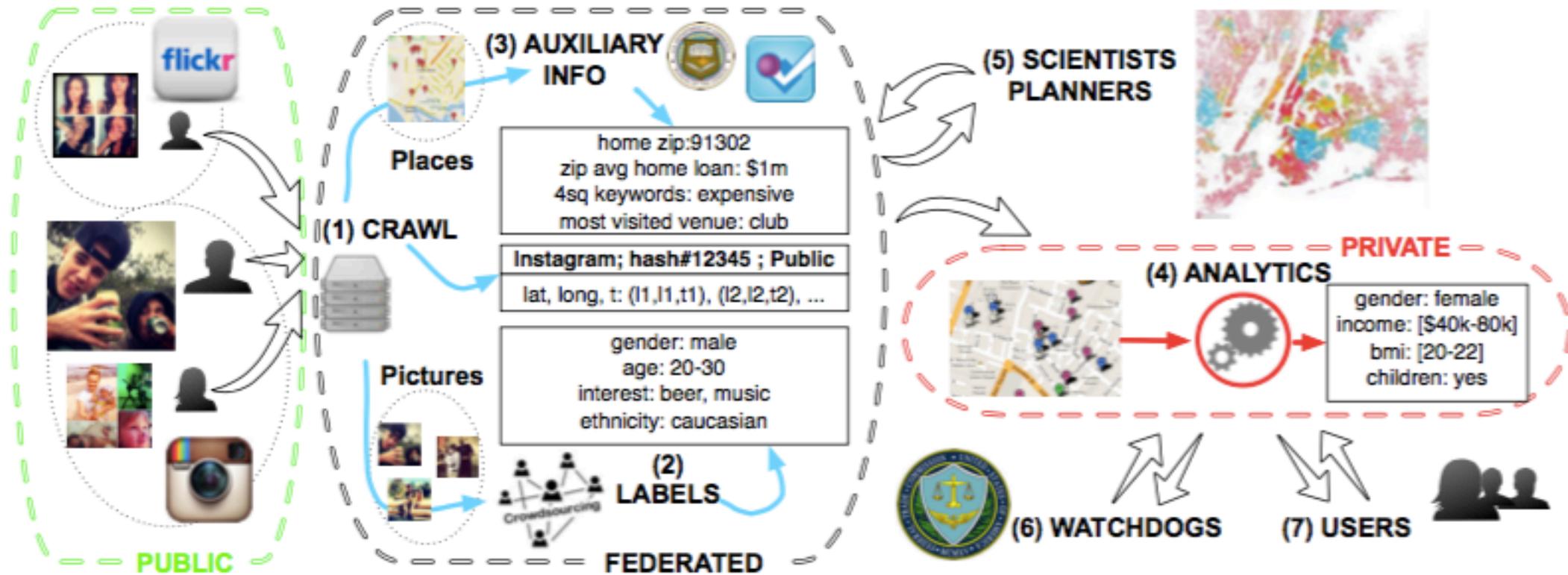
REALITY TEST IN 3 INCREASINGLY HARD EXAMPLES



OUTLINE

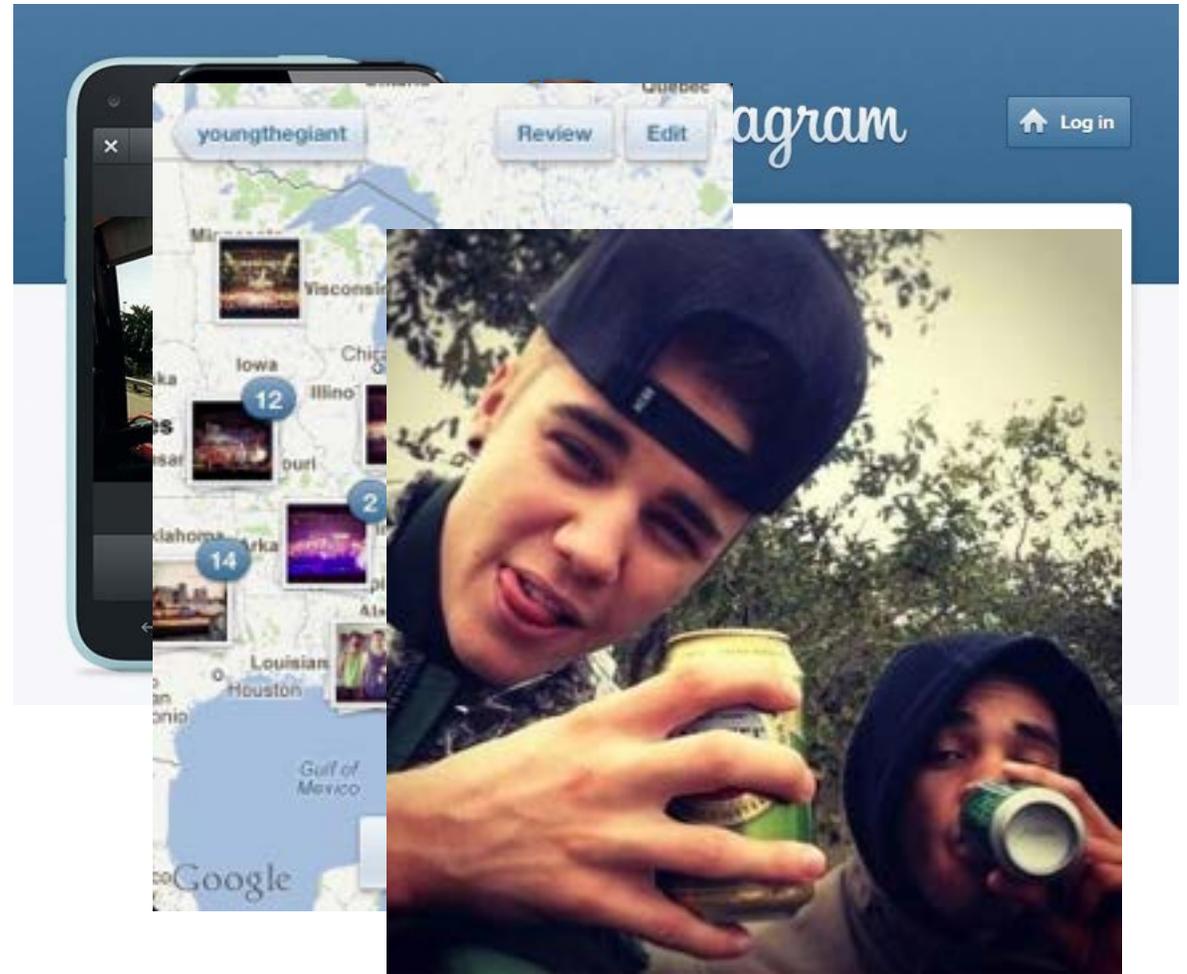
- How I came to care about demographics in big data
 - And why we need new tools for mobile social media
- Risks: Identity Reconciliation using location data
- Opportunities: Demographic study of Urban mobility
- How to scale?
- Conclusion

OUR PROPOSITION: MOBILITY DATA UNCHAINED



BUT WHERE DO YOU START?

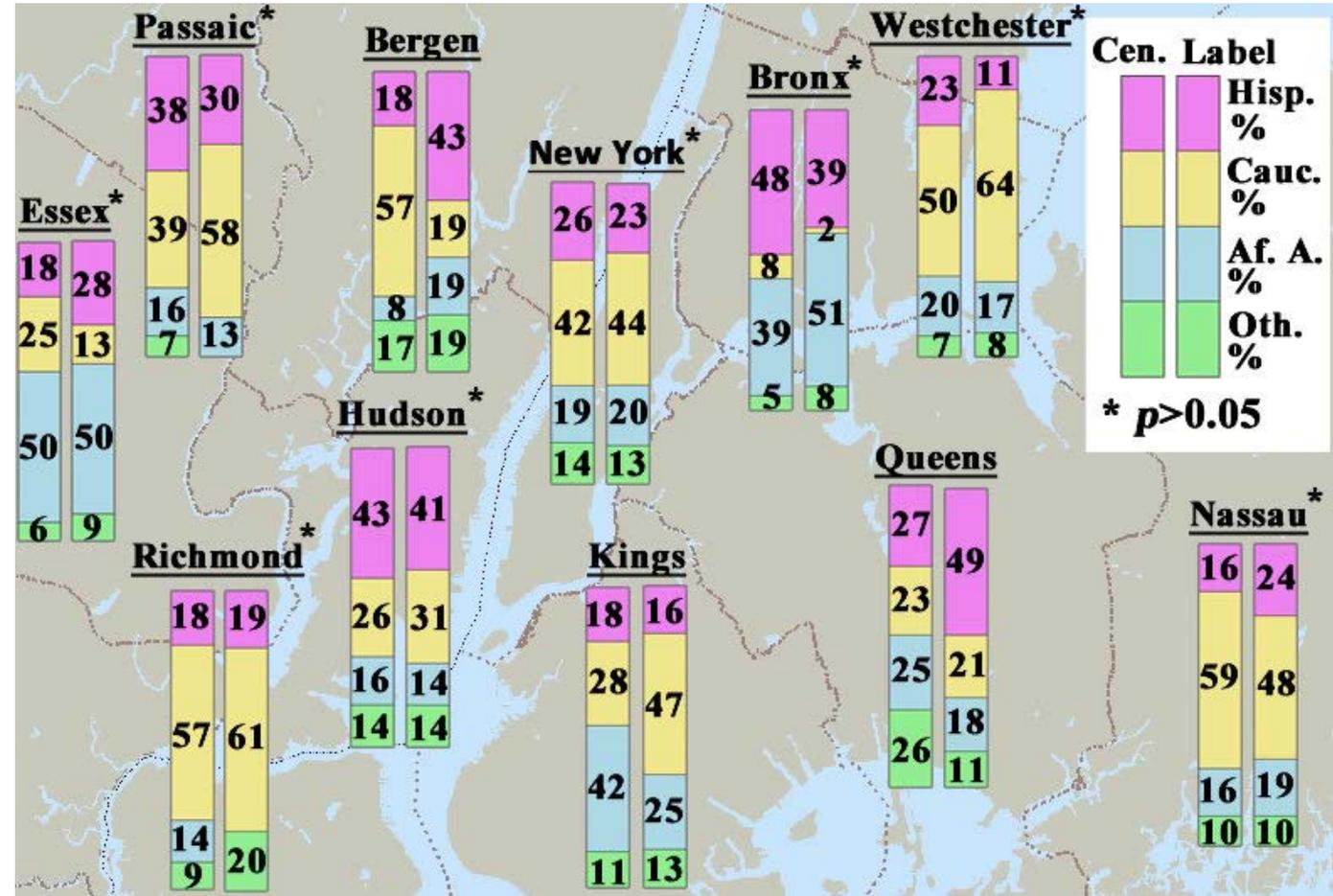
- Geotagged photos
 - 118,000 users
 - 35M geotags
 - Did **not** download or save any photos
 - Only public profiles
- Labeled in 2 ways:
 - Benchmark group
 - Crowdsourcing for scale



Male, Caucasian, Likes beer

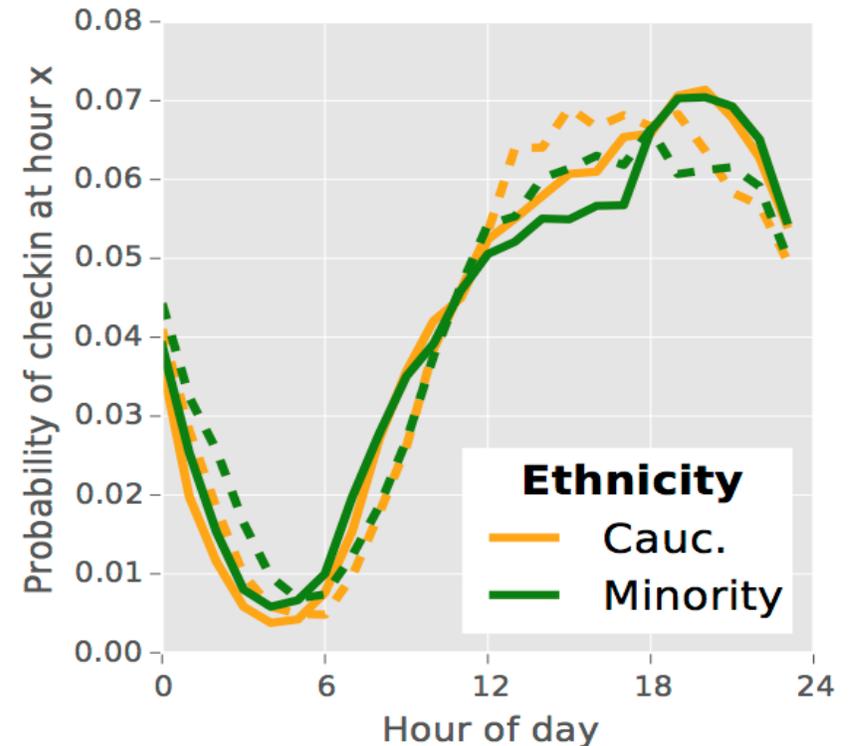
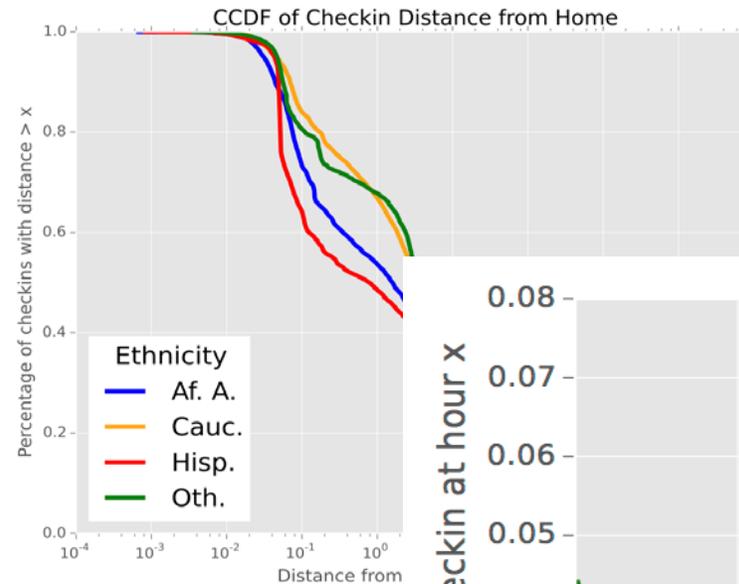
PROMISE #1: REDUCE COST OF SURVEYS

- Survey are (very) costly and slow, Call Details Records too.
- Reproduce CDR trends
 - Weekly maximum distances
 - Median distances traveled
 - Average distance traveled from home



PROMISE #2: ENRICH OUR UNDERSTANDING

- Distance traveled vary greatly across races
- Activity at night vary greatly across gender
- Even weekly temporal activity patterns vary across races

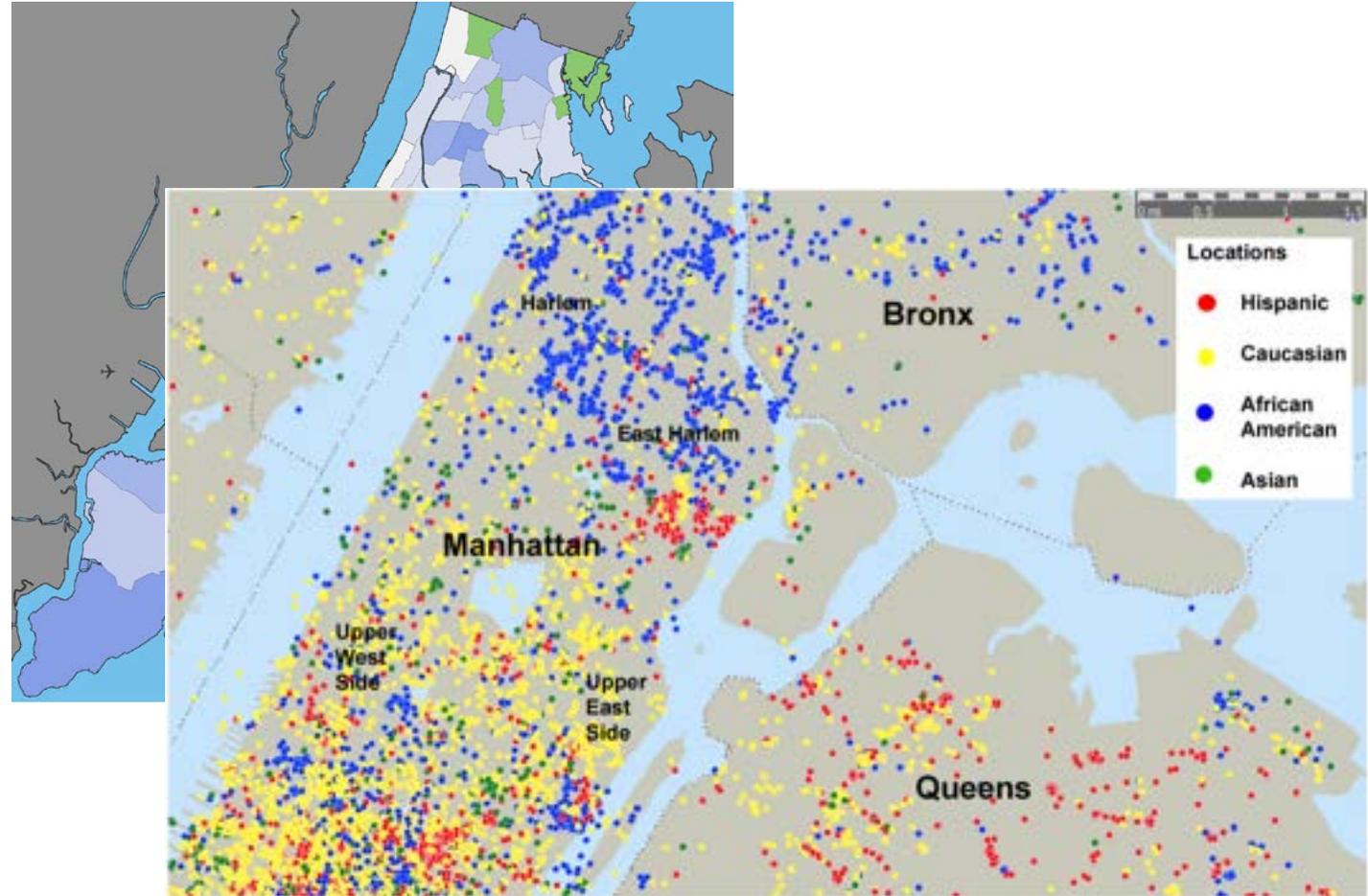


PROMISE #2: ENRICH OUR UNDERSTANDING

New Study of Segregation

- **Interaction Segregation:** how much people of different groups meet.
- **Typically studied with home location zipcode**

Using public social media we found more **interaction**



OUTLINE

- How I came to care about demographics in big data
 - And why we need new tools for mobile social media
- Risks: Identity Reconciliation using location data
- Opportunities: Demographic study of Urban mobility
- How to scale?
- Conclusion

WHY WE NEED A SECOND GENERATION?

- Collecting accurate labels from Crowdsourcing is bottleneck
 - Does not leverage the scale of information widely available
 - Remains slow and/or the most significant cost
 - Crowd not reproducible, and is it neutral?
- Can we leverage decade of **Computer Vision** progress?
 - Face++ API : Asian/Black/White, Female/Male, confidence
 - What if multiple faces? In multiple pictures?
 - Majority Rules/Confidence Majority Rules/Normalized Picture/Profile Picture

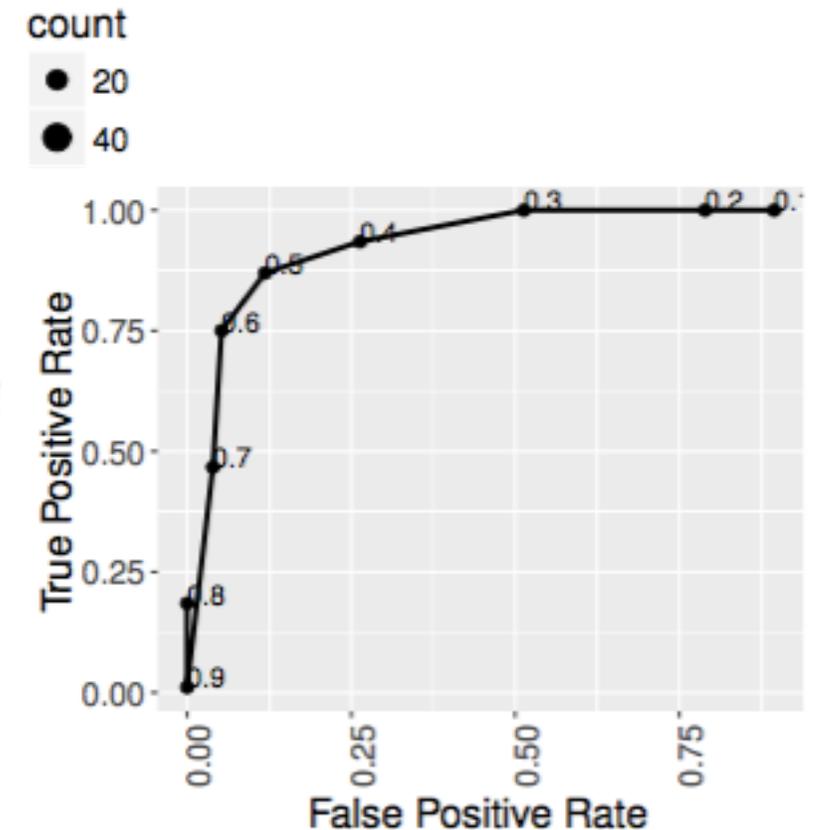
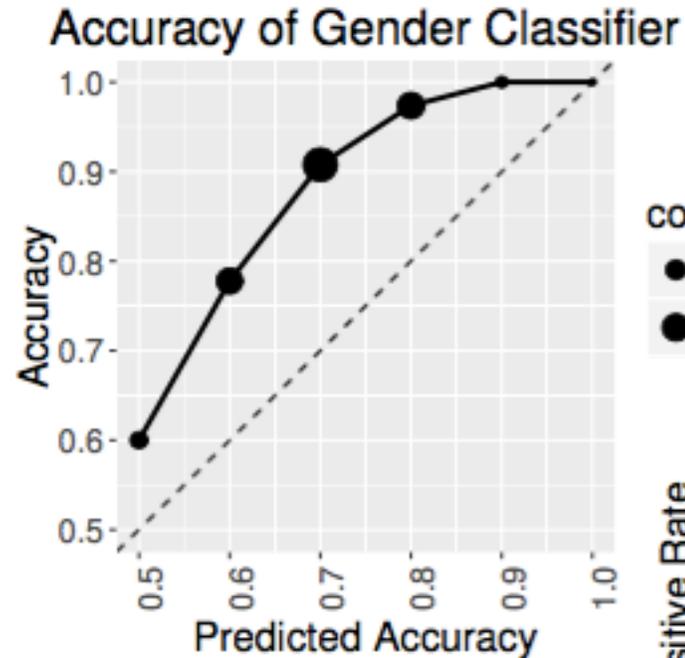


<i>Name</i>	<i>Cost</i>	<i>Speed</i>	<i>Accuracy</i>
Manual	Expensive	Moderate	High
Face Recognition	Free / Cheap	Slow	Moderate
Census	Free	Fast	Low

Table 2: Demographic Labeling Techniques

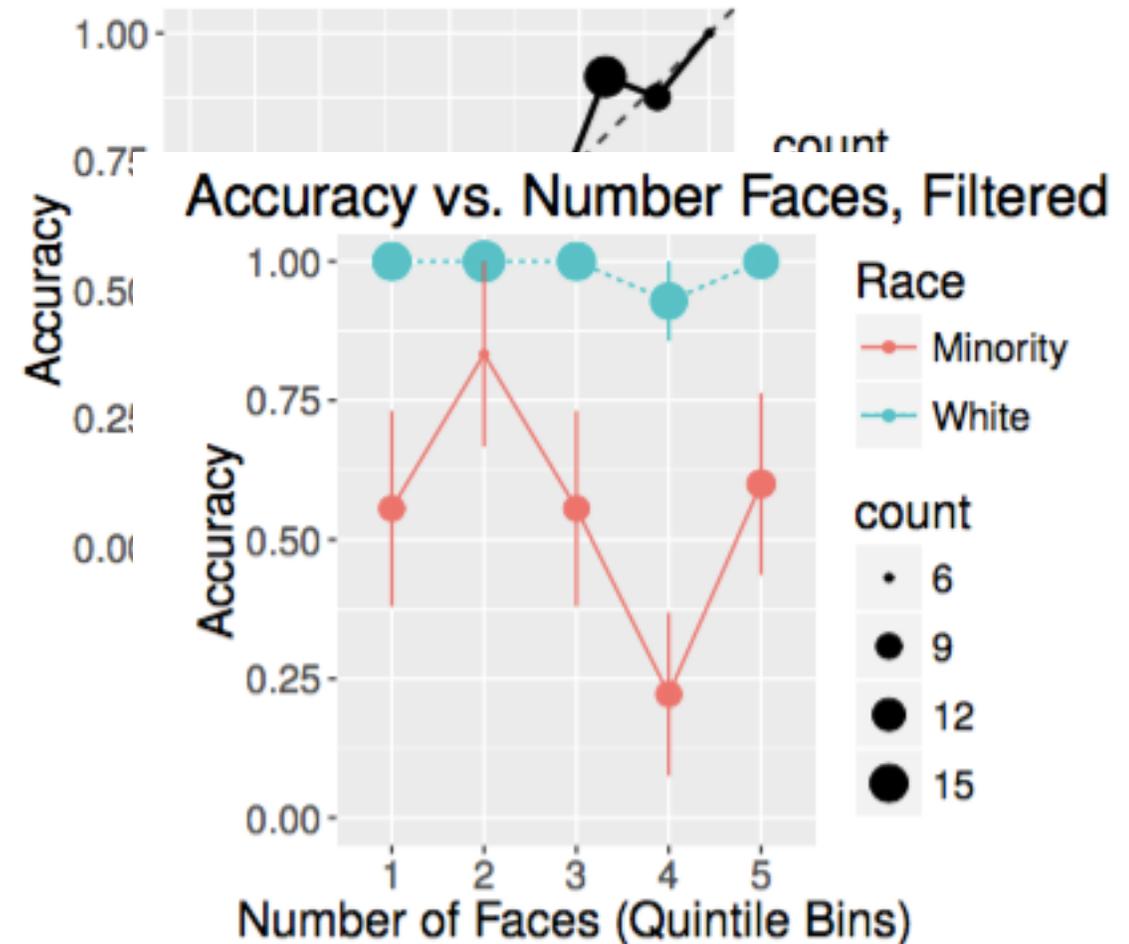
LESSONS LEARNED SO FAR: THE GOOD

- Promising: gender
 - Accuracy: 87.5%
- Conservative Pred.
- ROC curve
 - 0.6-threshold finds 75% female with only 7% misclassified male



LESSONS LEARNED SO FAR: THE UGLY

- Race indicates challenges
 - Accuracy $\sim 82\%$
 - Confidence may be overstated!
- But, wait there is worse!
 - Minority: more data means less accuracy! (tyranny of majority)
 - Simple “neutral” algorithms may under-represent specific groups!



LESSONS LEARNED SO FAR: THE LESS UGLY

- **Balanced Error Rate (BER)**
 - Average of errors made across classes (over-represents small classes in effective error rates)
 - Optimizing for this metric can remove some of the bias

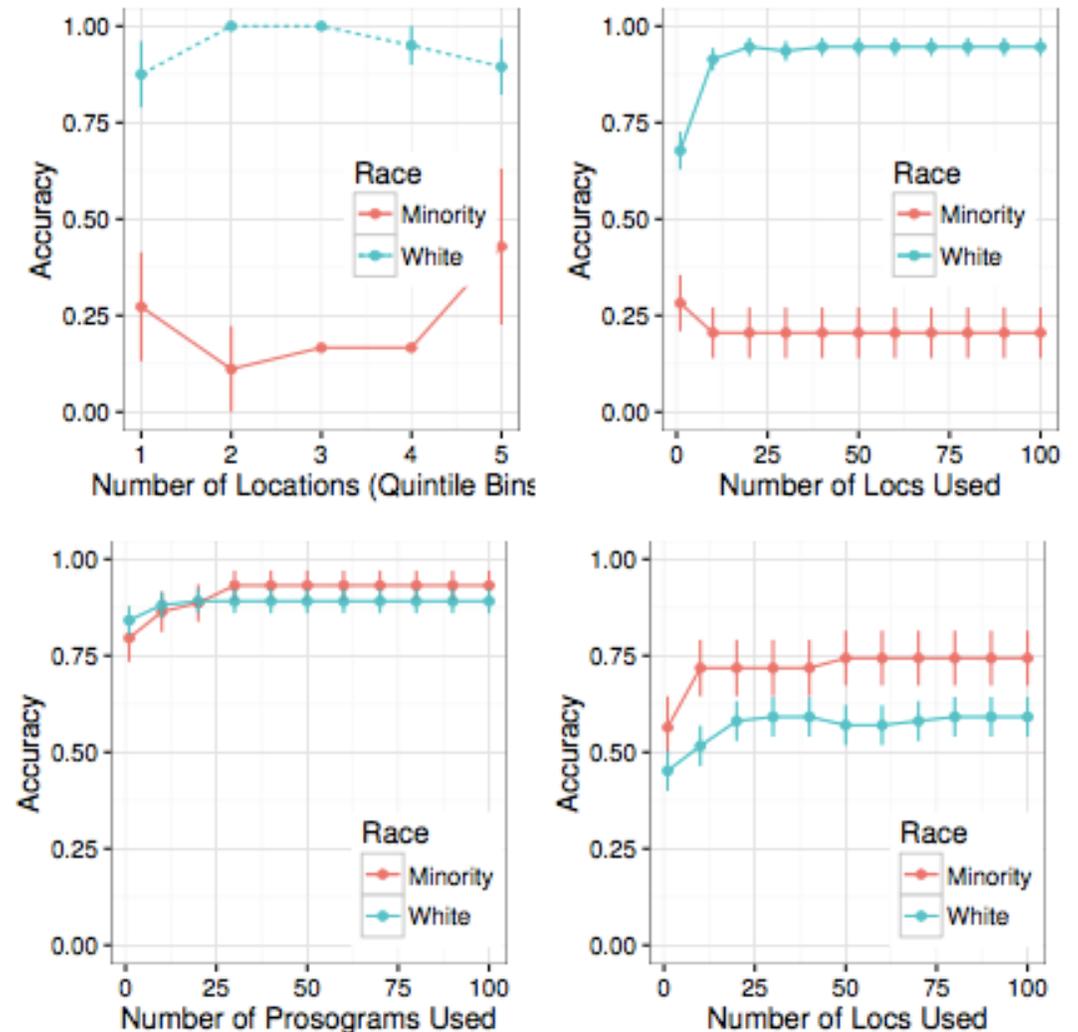


Figure 20: Location or faces vs. accuracy, with a choice of threshold determine by Balanced Error Rate.

POINTS FOR DISCUSSION

- Should we use Social Media for demographics research?
 - Yes, because it is feasible
 - Yes, if we want to protect users against new discrimination form
 - Yes, if we want to have equal representation
- But experience brings additional caution
 - Crowdsourcing raises cost/neutrality issues?
 - Even innocuous neutral algorithms can under-represent minorities
 - Critical to understand those effects when informing policies

THANK YOU FOR YOUR TIME!

- [1] C. J. Riederer, S. Zimmeck, C. Phanord, A. Chaintreau, and S. M. Bellovin, “I don't have a photograph, but you can have my footprints.: Revealing the Demographics of Location Data,” presented at the COSN '15: Proceedings of the third ACM conference on Online social networks, 2015, pp. 185–195.
- [2] C. Riederer, D. Echikson, S. Huang, and A. Chaintreau, “FindYou: A Personal Location Privacy Auditing Tool (demo),” WWW '16: Proceedings of the 25th international conference on World Wide Web, Apr. 2016.
- [3] C. J. Riederer and A. Chaintreau, “Scaling Up the Census with Social Media,” upcoming (presented at Proceedings of AAAI ICWSM Workshop on The Social Media and Demographic Research Workshop), May 2016.