

Lessons from Archives: Strategies for Collecting Sociocultural Data in ML

Eun Seo Jo, Stanford University

Timnit Gebru, Google

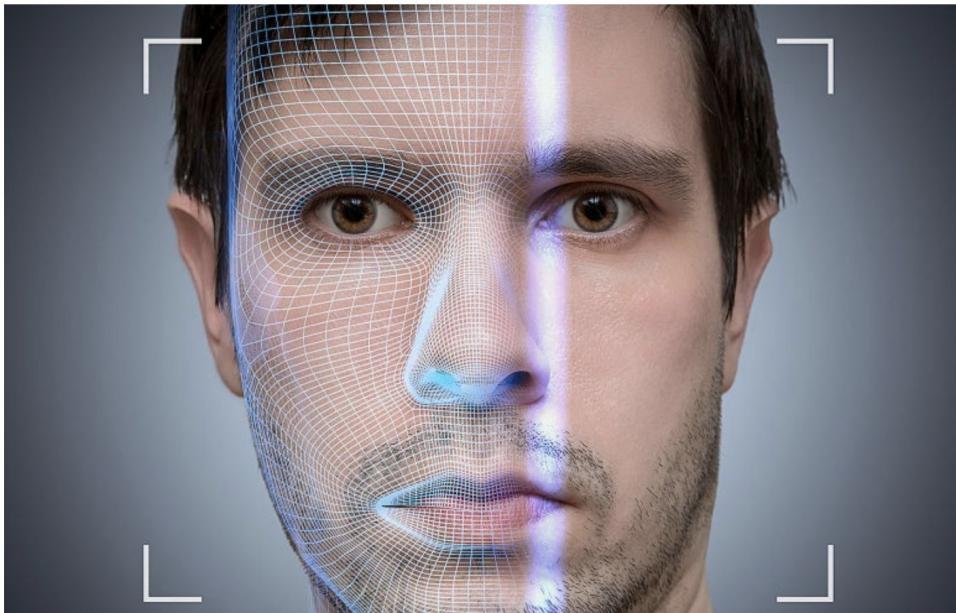
- Many issues in fairness accountability transparency and ethics are rooted in decisions surrounding the data collection, annotation and ownership practices.

The Government Is Using the Most Vulnerable People to Test Facial Recognition Software

Our research shows that any one of us might end up helping the facial recognition industry, perhaps during moments of extraordinary vulnerability.

By OS KEYES, NIKKI STEVENS, and JACQUELINE WERNIMONT

MARCH 17, 2019 • 8:32 PM



Counting the Countless

BY OS KEYES

Reconstructed transcript of a talk I gave at Seattle University earlier this year

MARCH 24, 2019

Good evening everyone! My name is Os, and I'm a PhD student at the University of Washington. According to my website I study gender, data, technology and control; it also says that I'm an inaugural Ada Lovelace Fellow. And I'm here for a variety of reasons, but one of the big ones is that I really enjoy giving talks. Particularly community-oriented talks; remaining grounded in my communities is important to me and for my work to be effective. So I was really pleased when, as a result of my *last* talk here, the Seattle Non-Binary Collective reached out. And they said: "we hear you're a data scientist. Could you do a talk on how trans &/ non-binary people can get involved in data science?"

And I replied: well, to be perfectly honest, I think data science is a profound threat to queer existences. And then for some reason they stopped replying! Who can say why? So when Jodi asked me what I'd like to talk about in this lecture series, I figured I'd do a talk on *that*. Why do I think data science is a profound threat for queer people?

The difficulty of definitions

Error Rate_(1-PPV) By Female x Skin Type



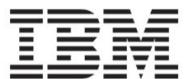
	TYPE I	TYPE II	TYPE III	TYPE IV	TYPE V	TYPE VI
--	--------	---------	----------	---------	--------	---------

Microsoft Face API	1.7%	1.1%	3.3%	0%	23.2%	25.0%
--------------------	------	------	------	----	-------	-------



FACE++

FACE++	11.9%	9.7%	8.2%	13.9%	32.4%	46.5%
--------	-------	------	------	-------	-------	-------



IBM	5.1%	7.4%	8.2%	8.3%	33.3%	46.8%
-----	------	------	------	------	-------	-------

Buolamwini & Gebru FAT* 2018, Slides from Joy Buolamwini

Microsoft improves facial recognition technology to perform well across all skin tones, genders

June 26, 2018 | [John Roach](#)



DIGITAL

Amazon Rekognition May Finally Be Audited and Ranked Alongside Other Vendors

A more universal test for facial recognition systems is needed

By Lisa Lacy | February 19, 2019



Start building apps today with
25+ free services and a \$200 credit

[Try Azure free](#)



[IBM Research Blog](#) Topics ▾ Labs ▾ About

AI

IBM Research Releases ‘Diversity in Faces’ Dataset to Advance Study of Fairness in Facial Recognition Systems

[NEWS](#) [MANAGEMENT](#) [OVERSIGHT](#) [DEFENSE](#) [TECH](#) [CONTRACTING](#) [PAY & BENEFITS](#)

Senators Are Asking Whether Artificial Intelligence Could Violate U.S. Civil Rights Laws

By Dave Gershman | [Quartz](#) | September 21, 2018 | [5 Comments](#)

RELATED



Shutdown Roundup:
NTSB Isn't
Investigating
Accidents, Warren
Warns of Decreased
Financial
Investigations and
More

January 25, 2019 | 1
[Comment](#)

Democratic Senator
Questions Mulvaney's
Hatch Act Compliance

September 21, 2018 | 32
[Comments](#)

Senators Kamala Harris and Cory Booker both signed letters to federal agencies asking about AI bias. J. Scott Applewhite/AP

Why Philadelphia Is
on the Federal
Government's

Google using dubious tactics to target people with 'darker skin' in facial recognition project: sources



By GINGER ADAMS OTIS and NANCY DILLON
NEW YORK DAILY NEWS | OCT 02, 2019 | 6:56 PM



A person takes part in a Google facial recognition project. (Obtained by Daily News)

ADVERTISEMENT

The advertisement is for Sling TV. It features a blue background with white text. At the top, it says "sling". Below that is a "LIMITED TIME OFFER" section. The main offer is "40% off" in large orange text, followed by "for the first month". At the bottom, there is a "WATCH NOW >" button and the small print "RESTRICTIONS APPLY".

ARGUMENT

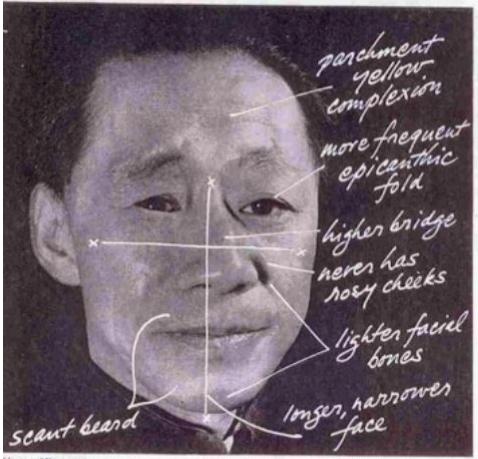
Beijing's Big Brother Tech Needs African Faces

Zimbabwe is signing up for China's surveillance state, but its citizens will pay the price.

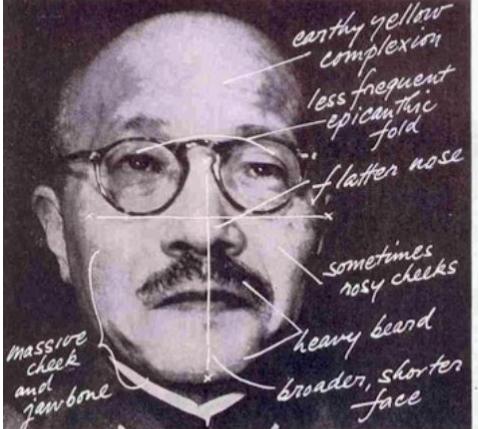
BY **AMY HAWKINS**

JULY 24, 2018, 10:39 AM

- Sometimes you need to make sure there aren't disparate error rates across subgroups (e.g. melanoma detection).
- Sometimes the task just should not exist (e.g. automatic gender recognition).
- Sometimes the manner in which the tool is used is very problematic because of who has the power (data) and ability to train powerful models, vs who is subjected to those models.



HONORABLE PUBLIC SERVANT. Ong Wen-han, is representative of the Chinese anthropological group with long, fine-boned men and scant beard. Epicanthic fold of skin above eyelid found in 88% of Chinese. Southern Chinese have round,



PEASANT WARRIOR. General Hideki Tojo, current Premier, is unusual, closer to type of humble Jap than highbred relatives of Imperial Household. Typical are his heavy beard, massive cheek and jaw bones. Peasant Jap is squat Mongol-

HOW TO TELL JAPS FROM THE CHINESE

ANGRY CITIZENS VICTIMIZE ALLIES
WITH EMOTIONAL OUTBURST AT ENEMY

In the first discharge of emotions touched off by the Japanese assaults on their nation, U.S. citizens have been demonstrating a distressing ignorance on the delicate question of how to tell a Chinese from a Jap. Innocent victims in cities all over the country are many of the 75,000 U.S. Chinese, who are here to increase their numbers. So serious were the consequences threatened that the Chinese themselves last week prepared to tag their nationals with identification buttons. To dispel some of this confusion, LIFE here adds these rules-of-thumb from the anthropometric conformations that distinguish friendly Chinese from enemy alien Japs.

To physical anthropologists, devoted debunkers of race myths, the difference between Chinese and Japs is measurable in millimeters. Both are related to the Fakimo and North American Indians. The modern Jap is the descendant of Mongols who invaded the Japanese archipelago back in the days of prehistory, and of the native aborigines who possessed the islands before them. Physical anthropology, in consequence, finds Japs and Chinese as closely related as Germans and English. It can, however, set apart the special types of each national group.

The typical Northern Chinese, represented by Ong Wen-han, Chinkiang's Minister of Economic Affairs (left, above), is relatively tall and slender built. His complexion is parchment yellow, his face long and delicately boned, his nose more finely bridged. Representative of the Japanese people as a whole is the peasant of Tojo (left, below), who betrays aboriginal antecedents in a more heavily built build, a broader, more massively boned head and face, flat, often pig, nose, yellow-ocher skin and heavier beard. From this average type, aristocratic Japs, who claim kinship to the Imperial Household, diverge sharply. They are proud to approximate the patrician lines of the Northern Chinese.



Chinese journalist, Joe Chiang, found it necessary to advise his nationality to gain admittance to White House press conference. Under Immigration Act of 1924, Japs and Chinese, as members of the "yellow race," are barred from immigration and naturalization.

One Month, 500,000 Face Scans: How China Is Using A.I. to Profile a Minority

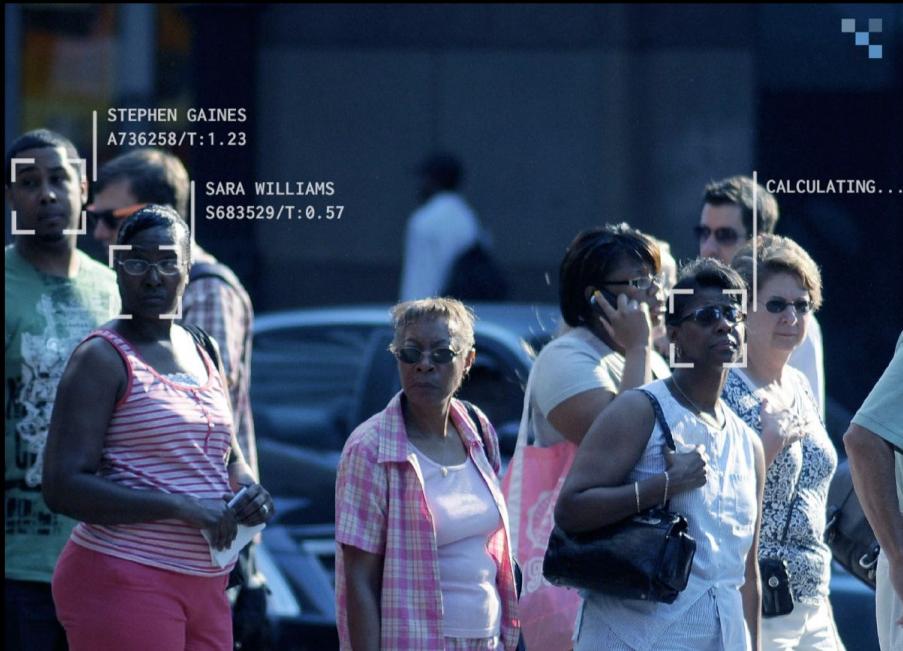
In a major ethical leap for the tech world, Chinese start-ups have built algorithms that the government uses to track members of a largely Muslim minority group.



US ADULTS INDEXED 130 MILLION

One in two American adults is in a law enforcement face recognition network used in **unregulated** searches employing algorithms with **unaudited accuracy**.

The Perpetual Line Up
(Garvie , Bedoya, Frankle 2016)



© 2016 Center on Privacy & Technology at Georgetown Law



Clare Garvie
@ClareAngelyn

INTRODUCTION

"Real-time video surveillance appears to be a simple question of supply and demand. As the technology improves, we anticipate that real-time face recognition systems will become commonplace."¹

—*The Perpetual Line-Up*, 2016

Authorities in Guiyang have eyes everywhere. Thanks to a vast, sophisticated camera system blanketing this Southwest Chinese city, police are purportedly able to locate and identify anyone who shows their face in public—in a matter of minutes. They can trace where you have been over the past week. If you are a citizen they can “match your face with your car, match you with your relatives and

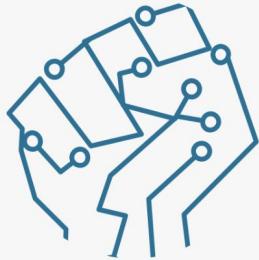
¹ See, e.g., *Surveillance Nation* (2012).

² See, e.g., *Surveillance Nation* (2012).

Amazon Pushes Facial Recognition to Police. Critics See Surveillance Risk.







Black in AI

[Donate](#)

Black in AI (BAI)

Black in AI is a place for sharing ideas, fostering collaborations and discussing initiatives to increase the presence of Black people in the field of Artificial Intelligence. If you are in the field of AI and self-identify as Black, please fill out [this Google Form](#) to request to join and we will add you to various platforms that we maintain. We also welcome allies to join our group using the Google form. Allies will be added to our email lists, where we send out group updates and requests for assistance.

Like our [Facebook Page](#) and follow us on [Twitter](#) to learn about our members and various activities!

Tweets by @black_in_ai

Black in AI Retweeted



Thabo Malete
@NerdBw

Highlight of last week: Got the opportunity to present my work on EEG based control at the [@black_in_ai](#) conference which ran alongside [#NeurIPS18](#) [@biustbw](#)

124
Shares



'Bias deep inside the code': the problem with AI 'ethics' in Silicon Valley

As algorithms play a growing role in criminal justice, education and more, tech advisory boards and academic programs mirror real-world inequality



Chad Loder @chadloder



Stanford just launched their Institute for Human-Centered Artificial Intelligence ([@StanfordHAI](#)) with great fanfare. The mission: "The creators and designers of AI must be broadly representative of humanity."

121 faculty members listed.

Not a single faculty member is Black.

People



Stefan Wager



Jim Breyer



Steve Denning



John Hennessy



Robert "Bob" King



Marissa Mayer



Samuel J.
Palmisano



Heidi Roizen

1,246 8:39 PM - Mar 20, 2019



1,501 people are talking about this



...while the fair ML literature has largely focused on “de-biasing” methods and viewed the training data as fixed, most of our interviewees report that their teams consider data collection, rather than model development, as the most important place to intervene

- We need a whole specialty in ML just dealing with data. We don't teach that in classes or treat it as a fundamental part of our scientific publications.

- Issues to do with data are complex. E.g.
 - Consent
 - Power
 - Inclusivity
 - Representation
 - Subgroup classification, annotation
 - Privacy & ethics
 - Transparency

- We need to have an interdisciplinary approach specifically focused on data. We can learn from disciplines such as anthropology and history

Anthropological/Artificial Intelligence & the HAI

26 MARCH 2019

Last week Stanford launched the [institute for human-centered artificial intelligence](#), and to kick things off [James Landay posted about the roles AI could play in society, and the importance of exploring smart interfaces.](#)

I've followed the HAI's development in passing, and I watched the inaugural event in the background on Monday last week while I was doing other work. I study algorithmic systems that make important decisions about us - which I call "street-level algorithms" in reference to Michael Lipsky's [street-level bureaucracies](#) - and some of the work I've done in the past has taken a more careful look at historical parallels between things we see today (like [quantified self](#) and [piecework](#)) to see if we can learn anything useful either for making sense of phenomena from a sociological perspective, and sometimes for informing the design of systems from an engineering perspective. James is a professor in the Human-Computer Interaction group at Stanford, and I'm a PhD student in that group.

So I was worried to find James leave details out from a series of anecdotes - details that would seriously undermine the point James seemed to be trying to make in his post. I started writing notes to call out how a more cynical perspective might describe the future or remember the past that James writes about; but with the launch of the HAI, the reaction from people around the world, and specifically *the responses from people in the HAI*, it seems like a more serious point that needs to be made.

The voices, opinions, and needs of disempowered stakeholders are being ignored today in favor of stakeholders with power, money, and influence - as they have been historically; our failure to listen promises to doom initiatives like the HAI.

[home](#)

[research](#)

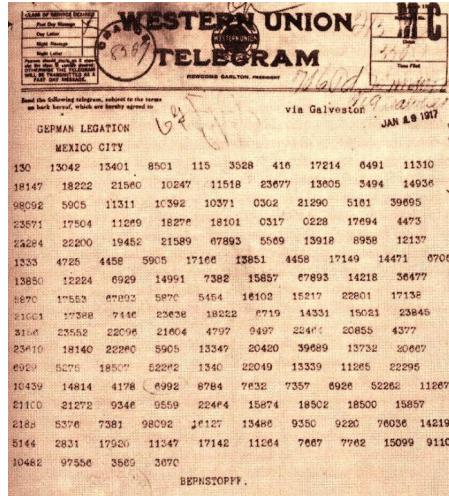
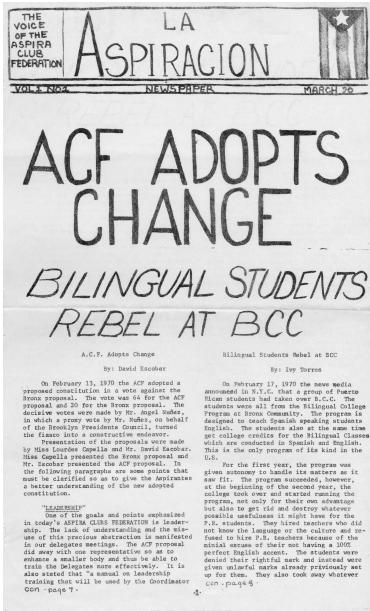
[blog](#)

[contact](#)

[CV](#)

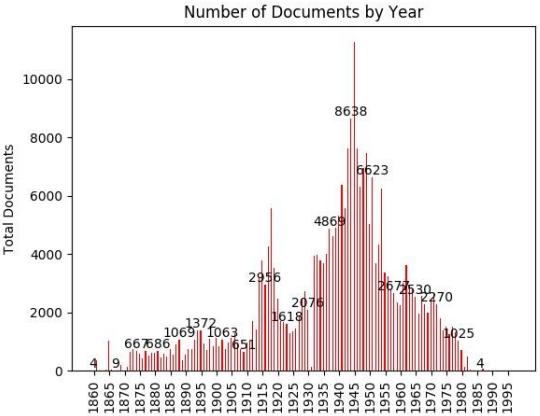
- Archives are the oldest attempt to gather human data

What are Archives?



- Institutional Missions /Agenda
- “Balanced and Representative Record”
- Private (Rockefeller Archives) & Public (UN Archives)

Case Study 1: U.S. State Department



DECLASSIFIED
October 29, 1957.
Volume 26, Page 166-167
By [Signature] Date [Signature]

TOP SECRET

I was visited by Professor Rabi, Admiral Strauss, Gordon Gray, and one or two others. The purpose was to bring to me their conclusions reached by Professor Rabi's Committee, called the Scientific Advisory Committee to the Director of Defense Mobilization.

Briefly, their conclusion was that we now enjoy certain advantages in the nuclear world over the Russians and that the measure of these gaps can be closed only by our own team on the part of the Russians. Admiral Strauss has therefore reached the conclusion that we should, as a matter of self-interest, agree to a suspension of all tests subject only to the installation of inspectional systems that would almost surely reveal the occurrence of any unauthorized test. He further believes that such tests should be conducted without any knowledge reaching the outside world, but the Rabi Committee believes that with a half dozen or so properly equipped inspectional posts inside of Russia, any significant explosion could be detected.

While the Rabi Committee agreed that the design of Russian weapons could be an advancement of ours, they say that the reported advantage would be as nothing compared with maintaining the particular scientific gap that exists in the design of the Russian H-bomb as compared to ours.

The nature of this gap is that Russian bombs are unshielded against certain types of *delta* radio activity that could be placed around them as they approach. The effect of this would not be to destroy the bomb but to reduce its effect by something like 99%.

Admiral Strauss and his group of scientists do not believe some of the assumptions made by the Rabi Committee. They feel that if we disclosed information on our tests, the Russians would, by stealing all of our secrets, equal and eventually surpass us. So Admiral Strauss and his associates believe we should continue all of our experiments and let them out in the open, refusing to be victimized by Russian technology. They are quite firm in their belief that we could not protect ourselves adequately against that duplicity.

The outcome was that Gordon Gray, Admiral Strauss and General Cutler are going to try to get (if possible) an agreement of scientific opinion in this whole matter to see what we should do about it.

Incidentally, I learned that some of the mutual antagonisms among the scientists are so bitter as to make their working together almost an impossibility. I was told that Dr. Rabi and some of his group are as antagonistic to Drs. Lawrence and Teller that communication between them is practically nil.

D.D.E.

TOP SECRET



- Historians at the State Department
- Culled from State/Presidential/CIA etc. Archives
- Declassified Documents, Hand Selected
- Procedure determining relevance + scope

Case Study 2: Internet Archive



TV Ad From Hillary Clinton Campaign

- Crawling Webpages
- 30+ petabytes/330 billion webpages
- 20+ years of internet
- No institutional mission (bigger → better)

Internet Archive is a non-profit library of millions of free books, movies, software, music, websites, and more.



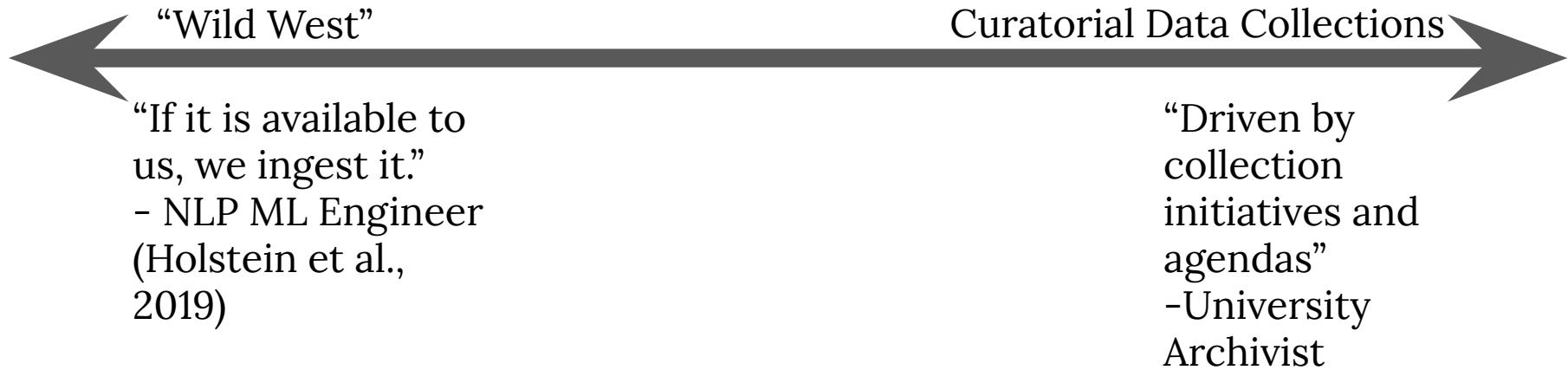
Data Collection Spectrum



Data Collection Spectrum



Data Collection Spectrum



Data Collection Spectrum

Poor Supervision of Data Collection and Handling

Supervision

Professional and Selective Curation



Data Collection Spectrum

Poor Supervision of Data Collection and Handling

Commercial Products/Research Output, Accuracy (minimal public outcry, consumer complaints, Holstein et al., 2019)

“Wild West”

Supervision

Objectives

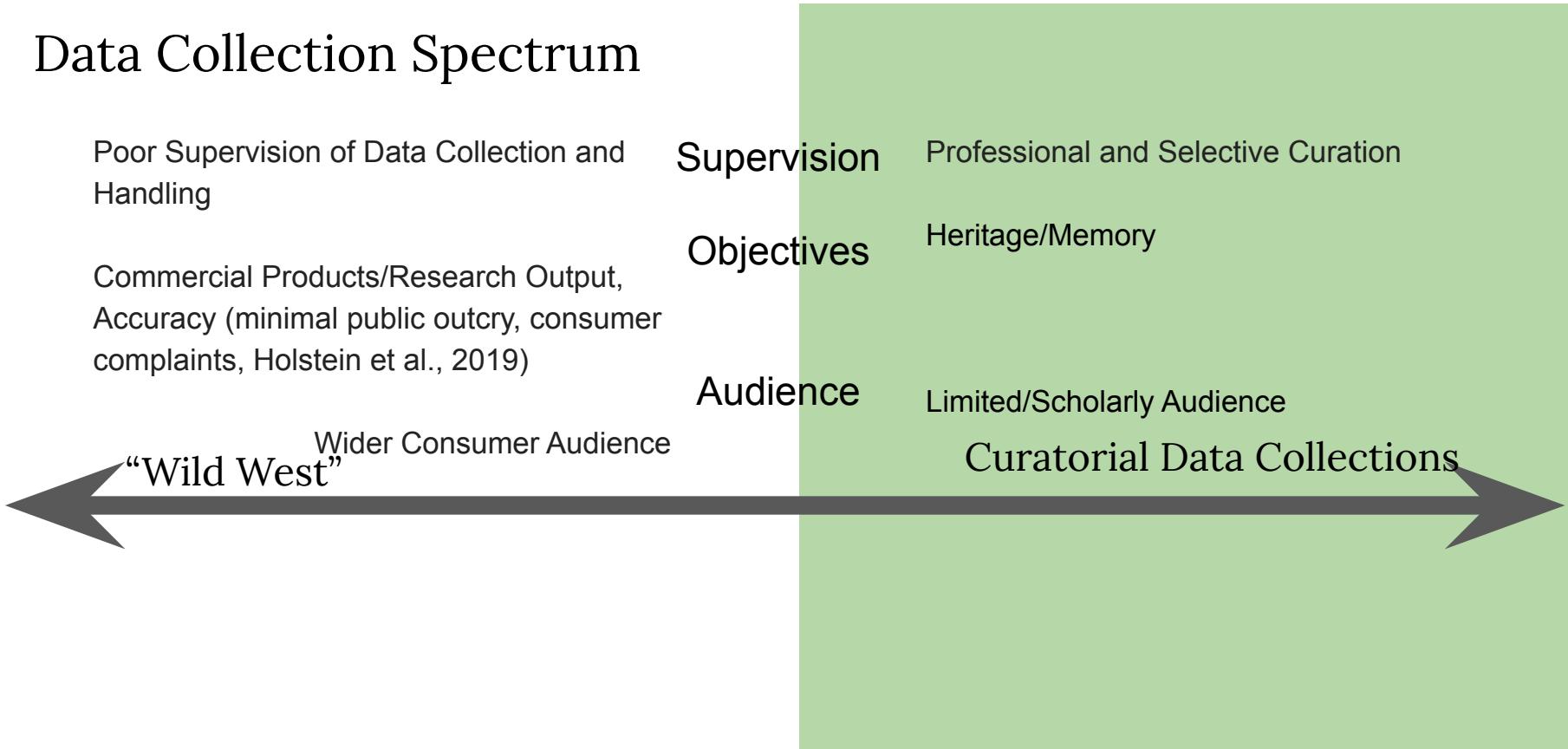
Professional and Selective Curation

Heritage/Memory

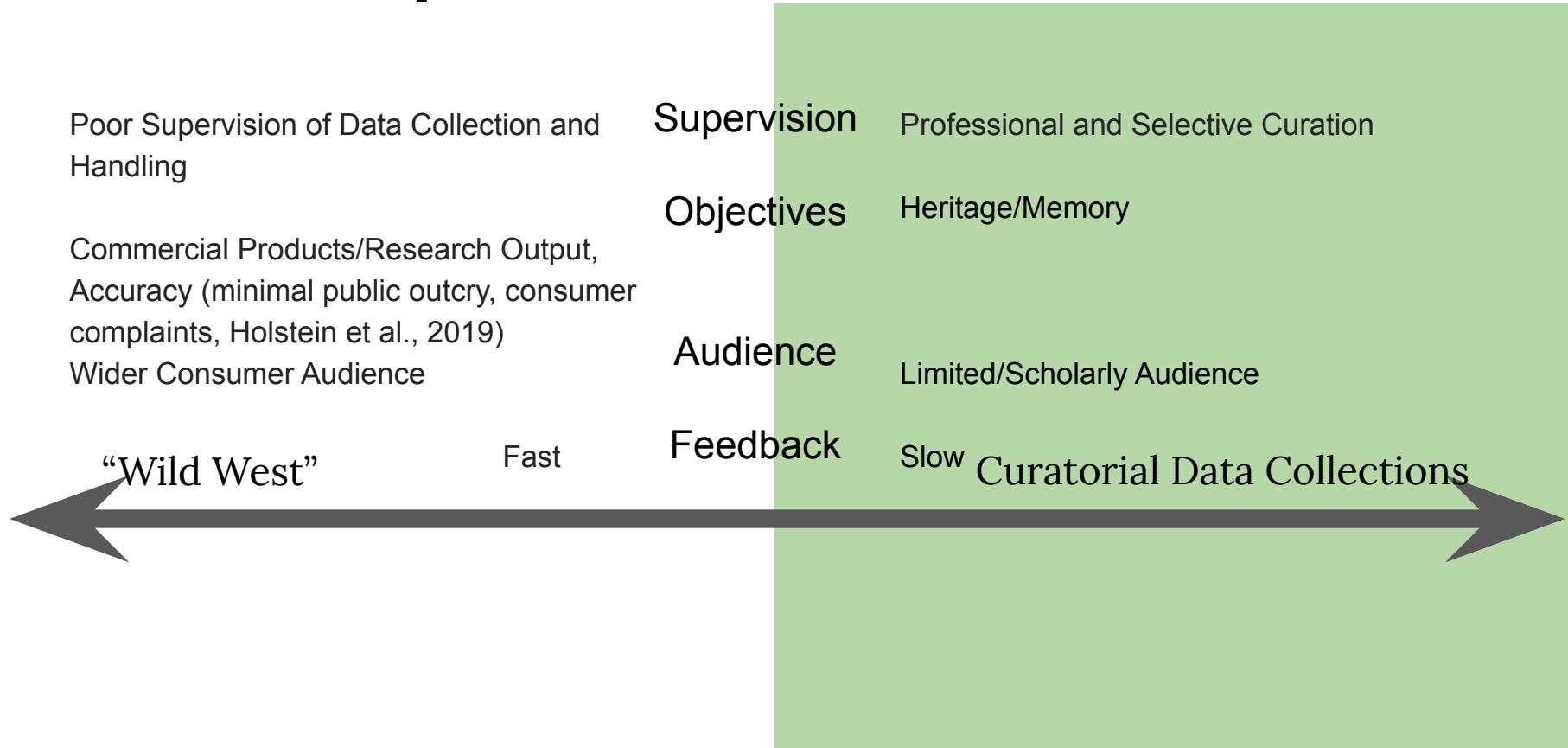
Curatorial Data Collections



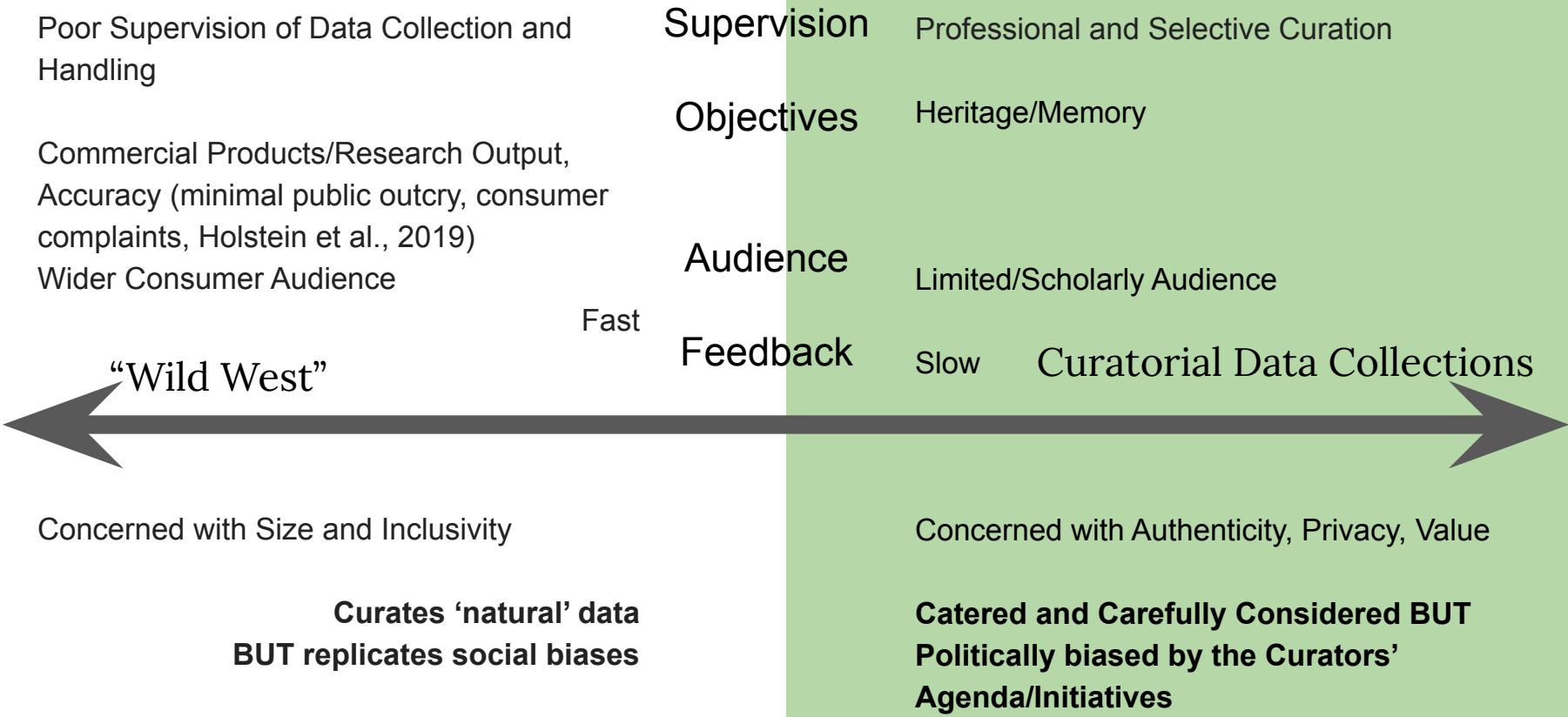
Data Collection Spectrum



Data Collection Spectrum



Data Collection Spectrum



Data Collection Spectrum

Poor Supervision of Data Collection and Handling

Commercial Products/Research Output, Accuracy (minimal public outcry, consumer complaints, Holstein et al., 2019)

Laissez-Faire

Wider Consumer Audience

“Wild West”

Domain/Regional Expertise and Labor Cost

Concerned with Size and Inclusivity

**Curates ‘natural’ data
BUT replicates social biases**

Supervision

Professional and Selective Curation

Objectives

Heritage/Memory

Audience

Limited/Scholarly Audience

Feedback

Slow

Curatorial Data Collections

Constraints

Space

Concerned with Authenticity, Privacy, Value

**Catered and Carefully Considered BUT
Politically biased by the Curators’
Agenda/Initiatives**

Interventionist

From MACRO to MICRO

Mission Statement

Collection Development Policy

Appraisal

Processing

From MACRO to MICRO

Mission Statement

Collection Development Policy

Appraisal

Processing

From MACRO to MICRO

Mission Statement

Collection Development Policy

Appraisal

Processing

From MACRO to MICRO

Mission Statement

Collection Development Policy

Appraisal

Processing

From MACRO to MICRO

Mission Statement

Collection Development Policy

Appraisal

Processing

1. Mission Statement

- Highest level of agenda formulation determining topic/concepts of concern.

2. Collection Development Policy

- A more specific policy drawn from the mission statement about what is collected, what is not, where to search for sources.

3. Appraisal

- Evaluation based on criteria of whether a given selection of sources is worth collecting.
- Asking whether this collection fits the outlines of the mission statement, evaluating the rarity and authenticity of the provenance, and its value for future generations.

4. Processing/Indexing (Micro-Appraisal)

- Processing the sources individually or at the folder/document level, including indexing them and updating the finding aid.
 - Sources may be discarded on the grounds of privacy concerns/irrelevance.
-

Shared Concerns

Digitization is seen as democratization of the archive. But we must think critically about what we are digitizing. Digitization programs could reinforce cultural stereotypes and canonicities. What are the criteria used to select manuscripts?

- Andrew Prescott, "Why do we Digitize? The Case for Slow Digitization"

Shared Concerns

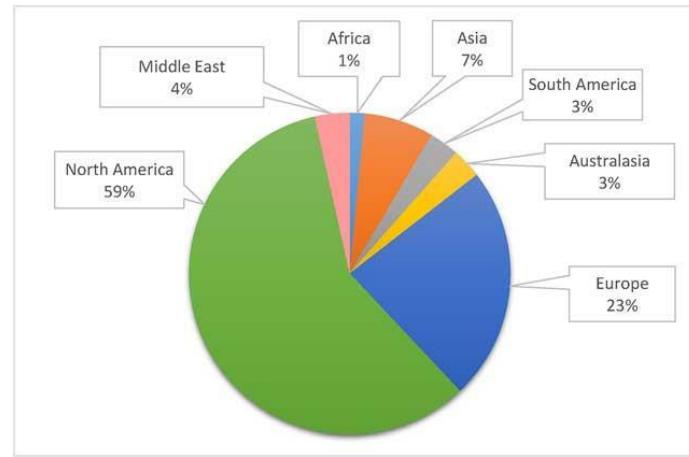
"It seems clear that the archive is primarily the product of a judgment, the result of the exercise of a specific power and authority, which involves placing certain documents in an archive at the same time as others are discarded. The archive, therefore, is fundamentally a matter of discrimination and of selection, which, in the end, results in the granting of a privileged status to certain written documents, and the refusal of that same status to others, thereby judged 'unarchivable.' The archive is, therefore, not a piece of data, but a status."

-Achille Mbembe, ``The Power of the Archive and its Limits"

Model 1: Consortium Models

“Ideally, we would have area specialists for each region”

- Internet Archive Archivist
- 1980s onwards → Consortium Agreements
- HathiTrust, Ivy Plus, OhioLINK
- International Coalition of Library Consortia (ICOLC)
- Human Genome Project
- Minimize overlap in collection efforts
- Open access to wider users
- Share resources such as human expertise



Regional distribution of ICOLC member consortia (Feather 2015)

Model 2: Participatory Archives/ Community Archives

- Allowing locals/grassroots to define their own collections
- Oral Histories
- Local Infrastructure
- Preservation Outreach (eg. Community History Toolkit)
- Mukurtu: Tool to Provide Access to Native Languages/Cultures - “using their own protocols”, “manage and share their digital cultural heritage in their own way”



Model 3: Codes of Conduct/Acquisition Policy

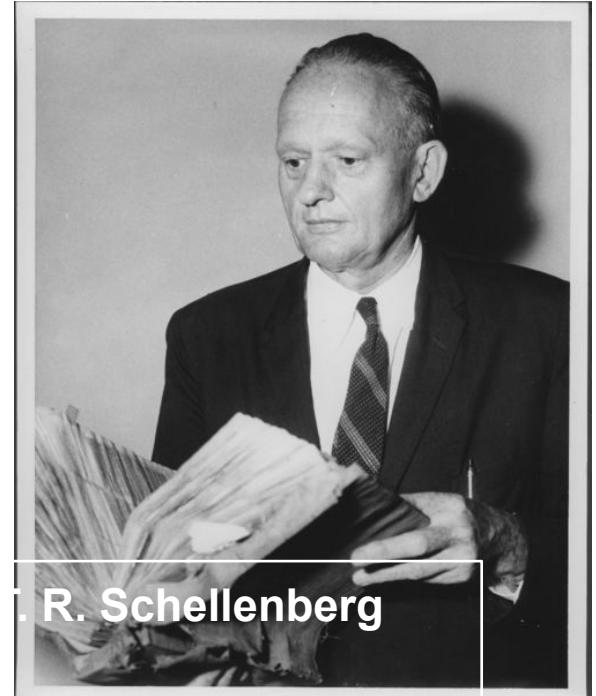
- International Council of Museum (ICOM): Code of Ethics for Museums
- "...minimum standards of professional practice and performance..."
- International Council on Archives (ICA): Code of Ethics
- Ethics, Privacy, Preservation, Access Standards
- "Right to be Forgotten": European Union General Data Protection Regulation 2016
- Society for American Archivists (SAA)



Model 4: Appraisal

- 5 Major Concepts & Methodologies
1. “Unbiased” Record Creator
 2. Organizational History
 3. Patterns of history/Speculation about Future Usage
 4. Professional Archivists
 5. Archival Darwinism: Random or Natural Selection

F. Gerald Ham, *Selecting and Appraising Archives and Manuscripts*



Appraisal Report Example

Background

Description

Evaluation

Figure 7-1 Appraisal Report

New York State Archives and Records Administration

Appraisal Report No: 90-25
Agency: Governor's Office of Employee Relations (GOER)
Subdivision: Child Care Unit (CCU)
Contact: Sandra Koss, tele. 473-3075
Action Originated: Submission of Schedule for Office
Title, Date and Volume of Records: New York State Labor/Management Child Care Advisory Committee Official Files, 1982-current, 6 cubic feet
Location of Records: Unit office, 1 Commerce Plaza
Condition: Good

BACKGROUND
The Child Care Unit (CCU) of the Governor's Office of Employee Relations is responsible for the development of New York State employee child care programs. The unit initiates and/or provides assistance in the establishment and expansion of State worksite child care centers under the aegis of not-for-profit corporations, provides seed monies and some operational funding for the centers, assists them in securing other sources of funding, renders general technical assistance to the centers, and monitors their performance in an advisory capacity. The unit also assists the New York State Labor/Management Child Care Advisory Committee (CCAC) in developing policies related to State Employee child care programs. CCAC also approves funding for child care centers and other associated projects through GOER.

DESCRIPTION
This series consists of 6 cubic feet of standard-sized paper records documenting CCAC monthly meetings. These include minutes, agenda, a copy of the information packet on issues to be discussed, resolution authorizing funding of programs and other activities, and summary fiscal reports from child care centers and various special projects. The series is arranged by meeting date. It is in good condition, presents no access problems, and contains no confidential materials. The series contains very good information on both CCAC and CCU functions and activities, policy development and decisions, individual State Employee child care centers, funding sources for these centers and other CCU programs, and various topics related to child care in general such as educational levels of child care staff, endemic high rates of staff turnover, and implications of evening and sick-child care.

EVALUATION
The wholesale movement of women into the workforce and the proliferation of single parent families has made child care a major social issue during the later 20th century in the United States. Relatively few families across the nation are untouched by the serious problem of how to ensure that their children receive proper care while parents earn a living. According to GOER staff members, the New York State employee child care program is an innovative and groundbreaking answer to this problem for State employees. This series provides excellent primary evidential documentation for this program and will allow policy researchers and historians

Title, Date, and Volume of Records:

New York State Labor/Management Child Care Advisory Committee Official Files, 1982-current, 6 cubic feet

Model 5: Cultural/Demographic Surveys & Guidelines

“We don’t know what we
don’t know”

- Internet Archive Archivist

- Uneven Digital Access
- eg. Native American groups running their own radio stations
- Cultural Ban/Taboo
- Critical Languages
- Organizations eg. Lighthouse for the Blind

Table 1

	N	%	Cumulative downwards %	Cumulative upwards %
more than 100 million	8	0.13		99.9
10–99.9 million	72	1.2	1.3	99.8
1–9.9 million	239	3.9	5.2	98.6
100,000–999,999	795	13.1	18.3	94.7
10,000–99,999	1,605	26.5	44.8	81.6
1,000–9,999	1,782	29.4	74.2	55.1
100–999	1,075	17.7	91.9	25.7
10–99	302	5.0	96.9	8.0
1–9	181	3.0	99.9	

David Crystal,
Language Death (2000)

Transparency and Inclusivity:

Mission Statements, Collection Policies, Community Archives

- Setting data collection agenda based on digital availability inevitably produces misrepresentation of data
- Left without active management of data composition, these methods can lead to poor or limited scope in demographics
 - Eg. what kind of users use Reddit?
 - Resulting models are heavily reliant on the specific source (ie. all of this was trained on Reddit data thus reflecting all users of Reddit).

- Archives have faced similar critiques of exclusivity. Traditional archives had focused on state and government documents, to preserve the documents of the governing and social elites
- To further inclusivity and data collection from the long tail, archives have expanded mediums for data input through
 - active collection policies
 - prioritizing data collection from minority groups
 - democratizing the data collection process via community archives.
- Rather than starting with data sets by availability, data collection in archives starts with a statement of concept or topics that often signal a commitment to collecting information about underrepresented groups.

- Community archives are projects of data or document collection in the ownership of the group that is being represented.
 - Projects such as historypin provide platforms for local communities to define and contribute their own cultural and heritage collections.
 - Widen the channel of user input in data collection.
- While many of the initiatives have grassroots beginnings, foundations and institutions have been actively funding and promoting archiving in the periphery.

- This model of decentralization also enables minority groups to consent to and define their own categorization.
- Some cultures demand non-Western systems of representation.
- Mukurtu is an example of a content management system built to allow indigenous communities to house their own materials “using their own protocols.”

Power: Data Consortia

- To increase parity in data ownership, archival and library sciences have developed a consortial model.
- Consortia have several mutual benefits for participating groups the main being the ability to gain economies of scale.
- Groups of libraries can make expensive purchases such as subscriptions to academic journals, pool resources for large scale projects

Ethics & Privacy: Codes of Ethics and Conduct

- Many organizations (e.g. Partnership in AI, IEEE Global Initiative on Ethics of Autonomous & Intelligent Systems)
- But these measures have yet to address the ethics associated with data collection let alone provide enforcement mechanisms
- Handling human information exposes archivists to various ethical dilemmas:
 - selecting which documents to toss or keep
 - granting access to sensitive content
 - dealing with intellectual property
- Overlapping organizations across archives, libraries, and museums independently have codes of ethics and conduct.

Ethics & Privacy: Codes of Ethics and Conduct

- Enforcing ethical regulation in data collection in ML faces challenges because of the lack of an incentive mechanism for researchers and practitioners
- Most archivists are full-time professional data collectors.
- **Archives work by a membership system whereby breaching the code of conduct could result in losing professional membership (ARA 2018).**
- Many suborganizations of archival and records collectors have ethics panels or committees that evaluate each alleged violation case by case (ARA 2018; ICRM 2016).

Conclusion

Societal biases enter when we

- Formulate what problems to work on
- Collect training and evaluation data
- Architect our models and loss functions
- Analyze how our models are used in society

Conclusion

Our papers, classes and incentive structures only focus on

- **Architect our models and loss functions**

Conclusion

We need to learn from other disciplines who have looked at data much more rigorously & work with them to draw lessons for ML.