

To Reduce Disability Bias in Technology, Start With Disability Data





The **Center for Democracy & Technology (CDT)** is the leading nonpartisan, nonprofit organization fighting to advance civil rights and civil liberties in the digital age. We shape technology policy, governance, and design with a focus on equity and democratic values. Established in 1994, CDT has been a trusted advocate for digital rights since the earliest days of the internet. The organization is headquartered in Washington, D.C. and has a Europe Office in Brussels, Belgium.



This report is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/).

To Reduce Disability Bias in Technology, Start With Disability Data

Ariana Aboulafia

Policy Counsel, Disability Rights in Technology Policy at CDT

Miranda Bogen

Director, AI Governance Lab at CDT

Bonnielin Swenor

Director, Johns Hopkins Disability Health Research Center

The authors of this paper would like to thank Henry Claypool, for sharing his insight and providing feedback on drafts of this document.

This work is made possible through a grant from the Ford Foundation.

Illustrations by Ananya Rao-Middleton.

Art direction by Tim Hoagland.

Contents

Introduction	5
Quantifying Problems With Disability Data Sets	8
Variances In Defining Disability, And How They Can Impact Accurate Data Collection	11
Data Implications Of Disability Stigma	15
Exclusionary Data Collection Practices	17
Concerns With The Disability Data Ecosystem	18
Disability Data Justice & Recommendations	20
1. Collect Disability Data In All Demographic Contexts	21
2. Respect Personal and Data Privacy	22
3. Develop More Inclusive Methods	22
4. Embrace a Growth Mindset	22
5. Include Disabled People in Technologies	23
6. Center Disabled Leaders in Policies	23
7. Make Data Collection and Storage Accessible	23



01

Introduction



When people with disabilities interact with technologies, there is a risk that they will face discriminatory impacts in several important and high-stakes contexts, like employment, benefits, and healthcare.

For example, many jobs use [automated employment decision tools](#) as part of their hiring process. These can include resume screeners and video interview tools that use algorithms to analyze things like vocal cadence or eye movements. These tools can unfairly screen disabled applicants from jobs by, for example, flagging the unusual eye movement of a blind or low-vision individual and removing them from the applicant pool as a result.

People with disabilities have also been [deprived of their benefits](#) when algorithms have been integrated into benefits determination systems, such as those that decide how many hours of home-based care a disabled person can receive through Medicaid. This, then, impacts the ability of those individuals to

live independently. Algorithms are also being incorporated into healthcare decision-making systems, such as playing a role in determining [who stays in a hospital](#) versus being discharged, as well as who [receives opioids](#) as part of post-surgical treatment, and much more. When these algorithmic systems create biased outcomes, people with disabilities can experience reduced health outcomes – and, the impacts of this technology-facilitated discrimination can be amplified for multiply-marginalized disabled people (including disabled people of color, and disabled LGBTQ+ individuals).

[Disability rights](#) and [disability justice](#) activists have a long history of fighting against discrimination that impacts disabled people. While technology-facilitated disability discrimination may represent a newer form of an old injustice, it is not going anywhere. Indeed, as technologies – algorithmic and otherwise – continue to become incorporated into everyday life, and as people with disabilities interact with them more and more, disparate and problematic effects will only increase, both in frequency and in severity.

While it is tempting to write off this bias as the result of the so-called algorithmic “[black box](#),” disparate and discriminatory algorithmic outcomes can often be linked back to problems with the data on which models are trained – and better data is likely to produce better results. Moreover, incomplete or erroneous data sets impact more than just technology. Data that is collected and used to quantify and generate insights about people with disabilities can also inform advocacy efforts for disabled people, including demonstrating the need for and supporting the development of disability-inclusive policies, [allocating funding for public benefits](#), and upholding disability-related civil rights laws. In order to tackle technology-facilitated disability discrimination – and improve the lives of people with disabilities overall – it is first necessary to understand, and then mitigate, the problems endemic to disability-related data.

This paper identifies the various ways in which data sets may exclude, inaccurately count, or be non-representative of disabled

people. It unpacks the factors that result in poor collection and availability of representative data sets, and provides recommendations for how to mitigate these concerns, which we collectively refer to as a “disability data justice” approach.

We highlight several recommendations, including:

1. Disability data should be collected in **all contexts** where other demographic data is collected.
2. Data should be collected and stored in ways that are **respectful of personal and data privacy**.
3. New and **more inclusive methods of both defining disability and collecting disability data** must be developed.
4. Practitioners should **embrace a growth mindset** around disability data.
5. People with disabilities should be included in **the creation, deployment, procurement, and auditing of all technologies**.
6. **Disabled people – particularly disabled leaders and those with technology, disability rights, or disability justice expertise – should be centered** in the creation and implementation of technology and AI policies.
7. Data should be collected and stored in ways that are **accessible to individuals with disabilities**.

While significant changes in data collection are needed to inclusively design algorithmic systems, these changes are possible – and necessary.



02

Quantifying Problems With Disability Data Sets



There are many ways that data sets can be non-representative of disabled populations. These include:

- **Non-inclusive.** Datasets should be considered non-inclusive when they do not include or properly label any data related to disability as a demographic. This often occurs when those building out the dataset(s) choose not to collect disability-related data, perhaps because they do not consider disability to be a demographic or they believe disability data is too difficult or sensitive to collect. When these datasets are used as inputs to a system or analysis, those systems likely will include information related to individuals with disabilities anyway, as some disabled people likely wind up as part of a dataset regardless. However, this data would not be properly categorized or tagged as belonging to those individuals, or being disability-related data, because respondents were unable to properly identify as disabled in the survey.
- **Underinclusive.** Datasets can be underinclusive – meaning that they do not include sufficient data about disability – for several reasons.

- » Some data sets undersample – that is, fail to include a sufficient number of disabled people in the dataset.
- » Lack of knowledge of information about whether people in a dataset have a disability can also create an underinclusive dataset.
 - Either of these can occur, for example, if data collectors obtained some data related to disability, but inadvertently exclude significant portions of disabled populations. These related issues can occur as a result of misunderstandings as to the definition of disability, as well as inability to reach certain portions of the disability community, including disabled people living in institutions, facilities, etc.
- **Inaccurate.** The data is inaccurate, meaning that the data that is collected is incorrect. There are also at least two ways that data can be inaccurate.
 - » Data about people’s disability or disabilities can be inaccurate; and/or
 - » Other data about a person can be inaccurate, notwithstanding whether data about their disability is accurate.
 - Either of these can occur, for example, due to mistakes in data collection or processes after collection, including consolidation into reports. While these are two distinct issues, they lead to similar challenges when attempting to create datasets that are properly representative of disability.

The above are all contributory factors in the creation of non-representative data sets on disability. When these faulty data sets are then used to create algorithms, those systems are more likely to work poorly, result in errors, or lead to inequity for both disabled and nondisabled people. Furthermore, when datasets used to design or train algorithmic systems reflect an underlying deficiency regarding people with disabilities and their experiences, the results can further marginalize these same groups, particularly when those

systems are integrated in areas like employment and healthcare. While there are ways to mitigate some impacts of biased systems, a lack of disability data also limits the ability to identify algorithmic systems' bias or discrimination on the back end, rendering proposed solutions like auditing minimally useful.

In addition to identifying common issues with disability-related data, and data collection practices, it is vital to consider why these issues exist and provide recommendations to minimize them.



03

Variations In Defining Disability, And How They Can Impact Accurate Data Collection



A root barrier to improving the status quo when it comes to disability data is that disability is defined in varying ways. Four different methods of defining and conceptualizing disability – **legal, medical, social, and identity** – each have important meaning, but complicate approaches to collecting and collating disability-related information. That is, because there are such significant variances in the identity, social, medical, and legal constructions of disability, accurately estimating the number of disabled people is inevitably difficult. Understanding the variances in defining disability is a vital step towards accurate disability data collection and organization.

Legal and statutory definitions of disability form the basis of many disability rights claims. The [Americans with Disabilities Act \(ADA\)](#), for example, defines individuals with disabilities as those with any physical or mental impairment that impacts one or more major life activities (including, but not limited to, walking, sleeping, and eating). The ADA also includes within its definition of disability individuals with a “history or

record” of such impairment, as well as individuals who would be “perceived by others as having such impairment.”

While the ADA is a profoundly important statute, it is not the only piece of legislation to define disability. For example the [Social Security Administration](#) defines a person as “disabled” (and thus eligible for benefits) if they are “[unable] to engage in any substantial gainful activity because of a medically determinable physical or mental disability (ies) that is either expected to result in death, or has lasted or is expected to last for a continuous period of 12 months.” These two legal definitions of disability illustrate some of the difficulties in finding a common definition for the purpose of data collection and analysis. These waters are further muddied by the existence of additional models and methods of defining disability.

The so-called **medical model** of disability views disability as caused by individual limitations, and divorces disability from social, economic, and political contexts. Most importantly, the medical model views a “cure” (that is, the elimination of the disability) as remediation for any disability-related concerns. This model can perpetuate ableism by presupposing that all individuals with disabilities need and want to change – that, if given the choice, they would not be disabled. When combined with [technosolutionism](#) – the idea of technology as a panacea – these notions of disability can lead to [technoableism](#) – the idea that disability is a problem to be solved, that the best solution is elimination of the disability through technology.

Alternatively, the **social model** of disability considers disability as something caused by individual differences (as opposed to inherent limitations) and views whatever limitations one may experience as a result of disability as inextricably connected to the social, political, and economic systems of that disabled people interact with on a daily basis. Rather than focusing on elimination of the *disability* as remediation for hardships, the social model of disability focuses on reforming these systems to make it easier for individuals with disabilities to live within them.

The **identity** model represents an evolution even past the social model, and recognizes that, for many individuals in the disability community, disability is part of their personal and cultural identity. Increasingly, individuals with disabilities consider [disability as part of who they are](#). It is this view of disability that has contributed to the idea that disability data should be collected as part of demographic data, and that disability (much like race, gender, or ethnicity) should be considered a demographic category.

The medical model generally defines disability through the lens of official diagnoses. However, other models of disability disagree, and allow individuals to determine themselves whether or not they identify as someone with a disability (also known as self-identification). A reflection of this can be found, for example, via the [Autistic Self Advocacy Network](#), a disability rights organization whose website specifically states that it considers “[a]utistic people who were diagnosed by a doctor, and autistic people who figured out they were autistic on their own” to be equally part of the autistic community. These constructions of disability generally create more inclusive environments – for example, by welcoming individuals who may not have access to healthcare resources that would allow them to receive a diagnosis – and allow for disability to be considered part of an identity, or a demographic, in a way that the medical model does not.

While the social and legal models may seem worlds apart from each other, the social model of disability is (at least, most closely) the model of disability recognized by the ADA. The ADA, after all, recognizes as disabled (and protects from discrimination as a result) not only individuals with particular impairments, but also those who are “perceived by others” as having such an impairment. As [Andrew Pulrang wrote](#) for *Forbes*, “the ADA defines disability not so much as a population, but in terms of a kind of experience... the experience of ableism.” Following this, it is the medical model that is an outlier, to a certain extent.

Furthermore, inaccuracies will naturally arise from the confusion among these multiple definitions of disabilities, as data gatherers (who are likely unfamiliar with subtleties inherent to defining

disability) use data sets derived from one definition (like one based on the medical model) when another would be more appropriate or accurate (like one based on demographic, or the social model). Nevertheless, underinclusive or noninclusive disability-related data can have many real-life implications, from obscuring who may need services, to limiting the ability to measure bias by a particular algorithmic system. For these reasons, despite its difficulty, it is vital to come to a consensus as to how to best create data sets that are properly inclusive of disabled people in any given population.

Discourse on defining disability came to a head recently, when the Census Bureau [proposed changing who counts as disabled](#) in its data-gathering processes. Specifically, the Census Bureau proposed changing questions on functional difficulties from yes or no answers (e.g., does this individual have difficulty dressing or bathing) to graded scales that ask participants to rate their level of difficulty in completing certain activities (e.g., regarding dressing or bathing does this person have no difficulty, some difficulty, a lot of difficulty, or are they not able to do the task at all?). With this change, only individuals reporting having “a lot of difficulty” or “cannot do at all” would be counted as disabled – a change which, [according to disability rights and justice advocates and experts](#), would have artificially reduced the official number of disabled people in the U.S. by nearly 40%. Following an outcry from researchers and activists, the Census Bureau abandoned its plan. Nonetheless, this situation illustrates how variation in disability definitions can have significant implications for disability estimates – which, again, can impact funding for benefits and services, and significantly (albeit inadvertently) contribute to tech-facilitated disability discrimination.



04

Data Implications Of Disability Stigma



The path toward accurate collection of disability-related data is also impacted by the stigma surrounding disabled identity.

In certain communities – both place-based and cultural – there is significant [social stigma](#) associated with disability (including, but not limited to, mental health conditions) that impacts disability data collection efforts. This stigma can lead individuals who may have a disability (either according to legal definitions, according to definitions used by data-gatherers, or in people’s private self-assessment) to avoid identifying as disabled in any sort of public forum.

In addition, factors such as grief and denial may play a role in not wanting to identify as disabled, particularly for those who become disabled later in life. People may also feel guilt around being “disabled enough” to qualify for limited resources, or fear discrimination – which may be particularly relevant when asking individuals to identify as disabled on job applications or the like. There is so much stigma surrounding disability, in fact, that for years even the word “disabled” was replaced by euphemisms, most of which have since [fallen out of favor](#).

The majority of these concerns can be traced back to the presence and proliferation of ableism, which is often internalized by disabled people. Internalized ableism can lead or contribute to feelings of grief and denial, as well as feelings of shame towards oneself, as well as to other people with disabilities. Ableism more generally leads and contributes to the societal conditions that lead to the very real discrimination and stigma that people with disabilities experience once they begin to identify as disabled. [In an essay for *Vogue*](#), Katie Baskerville wrote of both the impact of internalized and external ableism and how they impacted her decision to identify as a person with a disability. On internalized ableism, she quoted disability activist Rachel Charlton-Dailey, who said that “What we’re taught in society is that to be disabled is the worst thing you can possibly be – so, of course people are not going to want to identify as disabled.” Baskerville then candidly wrote of how disabled people experience ableism more generally, stating that “Being openly labeled as disabled could call into question my capabilities, mental or otherwise. In all honesty, I do not want to risk diminishing my perceived social value in exchange for further stigmatization and prejudice.”

These social forces lead individuals to avoid identifying as disabled in the first place. [One study](#) found that while 64% of surveyed adults had a health condition or impairment, only 12% identified as a person with a disability – numbers which are statistically significant enough to undermine confidence in the accuracy of data collection efforts related to disability, particularly those that do not do anything to control for underidentification. While truly and fully addressing disability-related stigma may require a societal paradigm shift and a reduction of ableism – and this is a worthy goal, to be sure – there are ways to combat exclusionary data practices that do not require such significant changes.



05

Exclusionary Data Collection Practices



Standard data collection approaches were not designed with disabled people in mind, which limits generalizability of findings and can lead to bias. Unaddressed barriers to the inclusion of disabled people exist at every step of the data collection process. Data collection efforts often do not include settings like group homes, or jail or prison facilities within their sample populations ([a disproportionate number of incarcerated individuals are disabled](#)). Data collection tools and materials are also frequently inaccessible: for example, survey questionnaires that are not provided in accessible formats, including compatibility with screen readers, exclude people who have low vision or are blind. These exclusionary data collection practices are of particular concern among technology companies, who tend to lag in their attention to the needs of disabled users and face countervailing pressure to minimize data collection of sensitive characteristics, including disability-related data. Despite the broad societal implications of biased data from the tech sector, laws and policies prohibiting discriminatory exclusion of disabled people from data collection opportunities are scant.

06

Concerns With The Disability Data Ecosystem



Various aspects of data collection practices have contributed to the creation of an unfavorable data ecosystem, which is partially responsible for the conditions that cause the inaccurate and under-inclusive collection of disability-related information.

In addition to considering the many dimensions of bias in the data that is collected, it is essential to consider the bias created by the data we do not have, and are not collecting. This applies to disability data, as far too few surveys, studies, and data collection efforts include questions to estimate the number of people with disabilities or to quantify accessibility needs or barriers. [These phenomena](#) – which collectively contribute to the systemic issues with the data collection ecosystem, particularly as they relate to disability – result in several areas of concern. Some of these concerns are heightened when data collected is then used to train or create algorithmic systems.

First, if data collection efforts fail to take active measures to include disabled people – by prioritizing the inclusive design of the collection approaches,

including centering the accessibility of collection mechanisms — practitioners will inadvertently exclude disabled people throughout the collection process, and likely will have no way of knowing they've done so. Excluding disabled people, and information about them, from these efforts, will result in data collectors producing datasets with no way of knowing whether those datasets are inclusive of people with disabilities.

Second, just as non-representative data sets create biased systems when used to train algorithms, when used to *test* algorithms those same sets make measuring bias difficult or impossible. Without data about disabled people's experiences, practitioners (including developers of algorithmic systems and technologies) will face challenges proactively testing how disabled people will be impacted by a policy change or a technology. Lack of data can make it more difficult to uncover whether seemingly anecdotal instances of harm and exclusion are in fact systemic and whether interventions to mitigate those harms have been effective.

Furthermore, this failure model feeds itself: the less inclusive data collection practices are of individuals with disabilities, the more likely it is that those individuals will simply choose not to participate. This, then, only amplifies the exclusionary and discriminatory outputs of the models that are trained with these inherently underinclusive data sets. This issue does not only impact disabled people. Many disabled people occupy several demographic categories – for example, disabled people of color, or LGBTQ+ people with disabilities. If people with disabilities choose not to participate in certain data collection efforts entirely, this could also impact the accuracy of data collection for the other demographic groups to which disabled people belong. This makes it more difficult to accurately understand and address algorithmic bias and discrimination generally.



07

Disability Data Justice & Recommendations



The disability rights and justice movements were built on the rallying cry of **“nothing about us without us,”** which has importantly evolved to simply **“nothing without us”** – in this final section, we provide recommendations for how to create representative datasets by following precepts of “disability data justice.”

In keeping with this philosophy, and in order to achieve justice and equality for everyone, people with disabilities must be involved in every level of decision-making, in every arena where those decisions can impact our lives, including technology. People with disabilities are already being impacted by technologies – indeed, the harms of AI and algorithmic systems for this population are not in the future, they are here already, and affecting disabled people in every aspect of their lives. These issues will only continue as technology becomes even more integrated in everyday life. Combating these harms, and mitigating algorithmic bias starts with ensuring that disability data is as inclusive as possible.

One of the goals of the disability rights and justice movements is the empowerment of disabled people, and the means by which to empower disabled people

in this context is by using a [disability data justice approach](#) when attempting to engage with disability data, particularly in the context of using this data to reduce technology-facilitated disability discrimination.

The authors recommend the following in keeping with utilizing a “**disability data justice**” approach to disability data and its collection.

1. Collect Disability Data In All Demographic Contexts

Disability data should be **collected in all contexts** — physical locales, data collection platforms, and otherwise — where data is collected related to age, race, ethnicity, and gender identity, as part of core demographic data collection. Disability data should not be limited to collection only in disability-specific places, and should be included when data is collected from systems related to systems including, but not limited to, healthcare, voting, employment, education, technology, and artificial intelligence. This can be referred to as a “disability data in all places” policy.

- Ensuring that **disability data is collected in more places** will provide larger data sets with at least some information about disability included, which can help contribute to more representative data sets overall.
- This recommendation **should not be construed to condone mass surveillance or any data-gathering done without knowledge or consent**, particularly of marginalized individuals. It is written in acknowledgement that core demographic data collection sometimes occurs in places, or under circumstances, where it should not, and in no way approves of data-gathering done in these contexts. However, in arenas wherein data collection is done responsibly, and that data collection includes the ethical gathering of information on demographics such as race and gender, disability should be included.

2. Respect Personal and Data Privacy

Data should be collected and stored in ways that are **respectful of personal and data privacy**, including policies for data minimization (collecting only the data that is necessary for a particular survey or project), purpose limitation (limiting the use of data collected for only the purpose originally agreed upon by subjects), and data deletion (both in the regular course and by request).

3. Develop More Inclusive Methods

New and **more inclusive methods of both defining disability and collecting disability data** must be developed. These new measures should include the diverse perspectives of disabled people, and the disability community and experts with lived experience must be centered in the development of these measures. The resulting measures must also more accurately estimate and capture information from a wider range of disabled people.

4. Embrace a Growth Mindset

Practitioners should **embrace a growth mindset** around disability data, which encourages auditors and data collectors to consistently work on improving their methodologies without discouraging them from using the information that they currently have. Underinclusive data (while problematic) may be preferable than fully noninclusive data for purposes of algorithmic fairness.

5. Include Disabled People in Technologies

People with disabilities should be **included in the creation, deployment, procurement, and auditing of all technologies**, particularly those that utilize algorithmic systems or any sort of biometric process, or that will be integrated into systems with high likelihood of significant consequences for people with disabilities, like education, employment, benefits determinations, or healthcare.

6. Center Disabled Leaders in Policies

Disabled people - particularly disabled leaders and those with technology, disability rights, or disability justice expertise - should be centered in the creation and implementation of technology and AI policies to ensure that these policies adequately address the potential threats (and benefits) these technologies pose to disabled people.

7. Make Data Collection and Storage Accessible

Because people with disabilities should be included and centered in every form of data gathering and technological use of data, **it is important that data be collected and stored in a way that is accessible to individuals with disabilities**, such as ensuring compatibility with screen readers and other forms of assistive technology.

Algorithmic systems, and many other technologies, are being created, deployed, and even evaluated and audited without disabled people in mind, and – just as concerningly – with data that does not recognize or reflect the presence or experiences of people with disabilities in society. These are both directly contributing to tech-facilitated disability discrimination and to entrenching ableism. When it comes to AI, algorithmic systems, and technology more generally, data is power. It is vital that people with disabilities be able to harness that power and use it to its fullest potential. More inclusive data is an important place to start.



 cdt.org

 cdt.org/contact

 **Center for Democracy & Technology**

1401 K Street NW, Suite 200

Washington, D.C. 20005

 202-637-9800

 @CenDemTech