**Analyzing the Dark Data Crisis Within the Framework of Value Sensitive Design**

A Research Paper submitted to the Department of Engineering and Society

Presented to the Faculty of the School of Engineering and Applied Science
University of Virginia • Charlottesville, Virginia

In Partial Fulfillment of the Requirements for the Degree
Bachelor of Science, School of Engineering

**Anish Varma Vegesna**

Spring, 2021

On my honor as a University Student, I have neither given nor received unauthorized aid on this
assignment as defined by the Honor Guidelines for Thesis-Related Assignments

Signature _____ Date _____
Anish Varma Vegesna

Approved _____ Date _____
Sean Ferguson, Department of Engineering and Society

**Analyzing the Dark Data Crisis Within the Framework of Value Sensitive Design**

Due to the explosion of "big data" within the last twenty years, more than 98% percent of the world's information is now stored digitally (Cukier & Mayer-Schoenberger, 2013). The term big data refers to structured and unstructured data that is so complex in quantity and variety that it is too difficult to analyze with traditional methods (Taulli, 2019). Organizations can then use this high volume of data to provide strategic business insights for the advancement of their company. However, as data collection has increased in popularity, so has the need for regulation concerning privacy issues and protection against potential hackers. In a recent study conducted by Internet Society and Consumers International, a survey of tech consumers from around the world found that 75% of people distrust the way data is shared and 63% of people find their connected devices "creepy" in the manner they collect data about people and their behaviors (Mori, 2019). This study shows that consumers are becoming more aware of how companies use their information and strongly believe changes need to be made in how we approach data privacy (Xu et al., 2012).

Despite the harmful consequences data collection can have, the act of data collection itself is not inherently benevolent or malicious. Rather, it is the implications of the data-driven technology that should make us rethink how we approach consumer privacy. For example, social media companies such as Facebook have been using artificial intelligence-powered prediction engines to more accurately advertise to consumers based on demographics and customer preference (Rosalsky, 2020). Predictive engines not only provide a tailored experience for the user, but also reduce marketing costs for the company by easily determining what marketing channels work best for each user (Kaplan, 2016). In this sense, AI and data-driven algorithms are

beneficial for both parties. However, the same technology can lead to "algorithmic amplification of some of the worst ideas." For instance, Facebook has received staunch criticism for recommending misinformation and conspiracy theories to users that have dangerously affected the political and social atmosphere (Rosalsky, 2020).

**Dark Data**

Despite the benefits of a personalized user experience, not all the concerns about data collection and algorithm development are fully recognized. Furthermore, the example above is just one instance of how irresponsible data collection can result in harmful consequences. This issue becomes even more alarming when we learn that 55% of all collected data is "dark data." This term refers to data that organizations may not be able to see or do not know has been captured, but can have a major effect on company decisions and actions (Briggs, 2017). David Hand, author of *Dark Data: Why What You Don't Know Matters*, states that dark data arises in many different "shapes and forms, as well as for many different reasons." Hand lists 15 different types of dark data and argues that an awareness of these gaps in our knowledge can equip organizations to identify when problems occur and protect themselves against danger (Hand, 2020). One type of dark data is "known unknowns," referring to gaps in data that conceal information which could have been recorded, such as missing table values or a low response rate. Another example is "summaries of data," referring to important details that go missing when we generalize or summarize findings in a study (Hand, 2020).

For our purposes, this paper will focus on "unknown unknown" dark data, referring to information that we are not aware is missing or being used due to excessive data collection, which can have harmful consequences if not properly managed (Choi et al., 2019). In the case of rising privacy concerns among consumers, the theft of "unknown unknown" dark data can be

potentially catastrophic for all stakeholders involved. As previously mentioned, dark data itself is not inherently morally right or wrong. The moral issue arises when data-collecting companies are not aware of the harm their practices can cause.

This paper will explore how the framework of Value Sensitive Design can help us understand how we can incorporate human values within our technologies to reduce the unintended negative consequences of dark data. Furthermore, we will also analyze whether there are value conflicts between stakeholders and if it is possible to bridge the gap between the ideals of corporations and consumers.

**Framework of Value Sensitive Design and Applied Case Studies**

Value Sensitive Design (VSD) is a design methodology that advocates for the integration of human values when planning and designing new technologies (Friedman, 1996). First developed by Dr. Batya Friedman, the framework was created when she noticed that designers were too focused on their technologies without considering the social implications of their creations. VSD seeks to be proactive by highlighting the root values at play by all stakeholders, such as businesses and their consumers, and determines if the technology in question supports or constrains their values (Umbrello, 2019). Specifically, this framework employs a tripartite methodology consisting of conceptual, empirical, and technical investigations. The first step in implementing VSD in the design process is to conduct a conceptual investigation, which consists of defining what values both stakeholders support in the design process and deciding how they can engage in tradeoffs among competing values. Empirical investigations then focus on if there are differences between espoused practices (what people say) compared with actual practices (what people do). This second step can help us understand what an organization's motivations, reward structures, and incentives truly are. Thus, empirical investigations may be the most

important step in understanding how designers can generate increased revenue, employee

satisfaction, and customer loyalty while still accounting for values. Finally, technical

investigations involve the proactive design of systems to support the values originally identified

in the conceptual investigation (Friedman et al., 2001).

Within the conceptual investigation, it is often beneficial to formulate a "value hierarchy"

that allows the moral values of stakeholders to be visualized (see Figure 1).
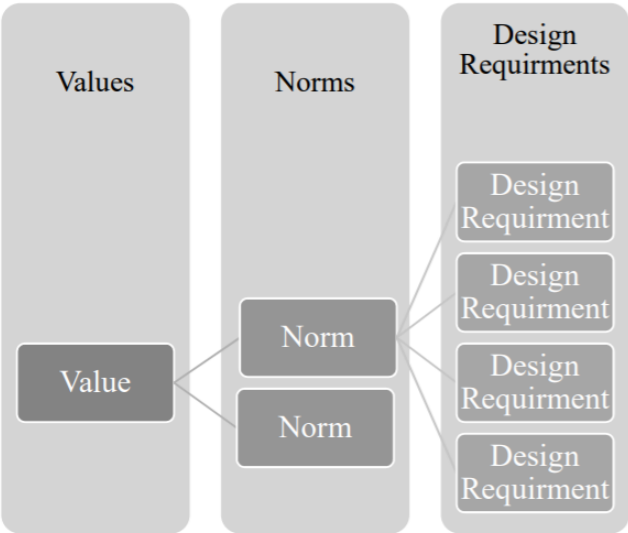


Figure 1. Bi-Directional Values Hierarchy

The newest addition to our framework of analysis is the idea of a norm. Unlike values, norms are

not objective and are sensitive to context and situation (Umbrello, 2019). The values hierarchy

provides a method to define stakeholder values, recontextualize those values within the desired

technology, and establish a set of design requirements based off these morals. We will now look

at examples of how two organizations have approached incorporating values within their data

collection process very differently, and how it has helped or hindered their designs.

**VSD in Artificial Intelligence Systems**

So far we have seen that technology is not value-neutral, but values are consistently integrated within design. For instance, artificial intelligent (AI) systems require actors and decision-making processes that direct and constrain developmental pathways towards certain areas rather than others (Umbrello, 2019). The Institute for Ethics & Emerging Technologies (IEET) was one of the first groups to evaluate the merits of Value Sensitive Design in relation to AI to reduce future AI risks. In their study, the UK employed the "House of Lords Artificial Intelligence Committee" to conduct a VSD conceptual investigation to determine a common list of moral values most important to universities/academics, NGO's/think tanks, governmental bodies, and for-profit industries. Among all the organizations, the top three values were identified as transparency, control, and data privacy. The bottom three values were identified as inclusivity, accessibility, and explainability (Umbrello, 2019).

The study was able to effectively make value conflicts apparent and it allowed designers to reduce the likelihood of uninformed decisions. However, the committee was unable to decipher how certain criteria, such as transparency, could easily be implemented into the design process. There were concerns expressed about how transparency regarding design choices or design principles between the creators and their users would be difficult, if not impossible (Umbrello, 2019). There are many instances of dark data-collecting technologies in which transparency is imperative, even if it comes at the cost of power and accuracy for companies. Regardless of if the IEET was able to successfully determine a value-oriented data collection process or not, applying the framework to their AI and identifying core values among stakeholders will help identify moral design requirements, which can later be expressed on a technical level. Moving forward, we will see how transparency between stakeholders is a value

conflict that arises often and one that can cause a rift in the company-consumer relationship, especially in the context of smart home devices.

**VSD in Dark Data Within Amazon Smart Devices**

As previously mentioned, this paper focuses on dark data that often results from excessive data collection. In a medical study exploring the attitudes towards video-based monitoring systems in smart home environments for the disabled and elderly, the most explicit fears expressed were the fear of data delivery without consent and fear of illegal access (Ziefle et al., 2011). Both fears can be categorized into the dangers of dark data, and are representative of how smart device users feel towards the most popular smart home device on the market: the Amazon Echo.

Amazon has received massive scrutiny over the last seven years since the release of their Amazon Echo, a brand of smart speakers that is also connected to the voice-controlled personal assistant, Alexa. The speaker has been advertised to include voice interaction, music playback, alarms, and even act as a home-automation hub (Fowler, 2019). For instance, the voice-controlled personal assistant activates when the user announces a "wake word," to prompt the device to listen for a command. Amazon claims that no audio is stored on the device or sent to the Cloud unless the device detects a wake word (Fowler, 2019). However, there are many instances where the speaker begins recording by accident although no user purposefully activated the wake word. Users also often report concerns about where the data is stored (on the device itself or on the Cloud), how long the data on the device remains, and whether data is ever fully erased (Privacy International, 2019). Overall, there is little transparency between Amazon and their customers regarding their data collection, which adds to the catalog of dark data involving consumer information.

Groups such as the Institute for Ethics & Emerging Technologies and Amazon have approached data collection in very different manners. On one hand, IEET has kept the Value Sensitive Design framework in mind before developing artificially intelligent systems, whereas Amazon has received public backlash because of their inability to do so. The next section will explore catastrophic consequences dark data can have, answer if value sensitive design can limit these consequences through a conceptual, empirical, and technical investigation, and discuss how we can re-align stakeholder values to create a best-case scenario.

## Bridging the Gap Between Corporations and Consumers

### Dark Data in the Current Landscape

In worst-case scenarios, irresponsible data collection and the mishandling of dark data can have catastrophic effects on consumer privacy. The credit bureau known as Equifax recently paid more than $650 million to resolve claims stemming from a 2017 data breach that exposed sensitive information on more than 147 million Americans (Cowley, 2019). In the largest settlement of a data breach in terms of dollar amount and number of victims, the Equifax breach highlighted how little autonomy consumers have over their private information. Consumers are able have some control, such as freezing files to prevent new credit lines from being opened, but they do not have the ability to stop credit bureaus from collecting information and selling it to auto loan, mortgage, and credit card issuers (Golinger, 1998). Furthermore, credit reporting is the second most frequent source of consumer complaints due to the number of inaccurate information about an individual on their credit report, with more than one in five consumers having a potentially material error in their reports (Klein, 2017). The dangerously large and incorrect amount of customer dark data collected and sold by credit bureaus can have catastrophic effects on consumer privacy, as seen in the Equifax breach. Of course, Equifax is

not the only company within the last decade to compromise user information. Timehop, Under Armour, and Yahoo have also been involved data breaches that leaked the information of millions, sometimes billion, of users (Mele, 2018).

**Applying VSD to Dark Data-Collecting Technologies**

Because dark data itself is not a technology, our analysis within the Value Sensitive Design framework will be applied to systems that collect excess user information, also called "unknown unknown" dark data. For our purposes, we will mainly refer to Equifax's data collection systems to conduct conceptual, empirical, and technical investigations. However, the conclusions made at the end of the investigation can be applied to all technologies that collect excess data.

### Conceptual Investigation

A conceptual investigation defines what values should be supported in the design process and answers questions such as if moral values (i.e., right to privacy) should have greater weight than non-moral values (i.e., power, efficiency, aesthetic preferences) (Friedman, 2001). As previously discussed, transparency and control were listed as the most important values among various academics, NGO's, governmental bodies, and corporations as design requirements for AI technology. The study regarding medical technology in smart home showed similar results in which participants stated they were comfortable with the technology only if personal data was anonymously transferred to physicians under their jurisdiction (Ziefle et al., 2011). Thus, it seems for now as if there is no discrepancy in what values all stakeholders deem important: data transparency and user control. Although, we will later explore why this is not necessarily the case. Because values are objective in nature, we can then recontextualize these values within the specific Equifax data collection system to define what norms may be required. In this case,

transparency may mean providing the customer with details on the credit data collected, free credit score updates, information sold to various agencies, and more. User control in this case may mean allowing customers to identify possible information errors and choose what organizations receive their information. Finally, these norms can be translated into applied practice and design requirements by providing immediate phone or email alerts to customers after any activity involving their credit information and by developing a liability structure that reduces the amount of inaccurate customer information within their systems.

**Empirical Investigation**

If we had previously concluded that there is no discrepancy between stakeholder values, we must ask why it is that Equifax's data collection system does not allow for the norms and design requirements just described: where does the misunderstanding lie? Empirical investigations highlight human activity and point to what true motivations an organization may hold and what reward structures and economic incentives are in place (Friedman, 2001). The IEET points to "power and accuracy" as common motivations for not seeing the established values all the way through the design process (Umbrello, 2019). This is especially apparent within social media systems in the case of data-driven algorithms that rely on developing accurate algorithms to increase customer engagement, increase revenue, and compete with other companies. Unfortunately, harmful dark data is often a result of an attempt to improve algorithms and achieve success within the market.

In the case of smart home devices such as the Amazon Echo, companies rarely have an economic incentive to collect less data about their users. Thus, moral values such as transparency, user control, and privacy are not seen as imperative as non-moral values such as efficiency and accuracy. Similarly, Equifax and other credit bureaus rely on users' financial data

and credit reports to aggregate and re-sell to lenders (Sweet, 2017). Furthermore, as the

gatekeepers to whether their can customers receive a credit card, affordable house loan, or a car,

credit bureaus have the ability to limit user autonomy by charging them for receiving information

about their credit score on an annual basis (Sweet, 2017). Identifying the true motivations of dark

data-collecting companies may help us understand if it is possible to incentivize companies to

align their values with those of users while also meeting non-moral values.

### Technical Investigation

A technical investigation consists of two parts: the first step determines how the

technology supports or constrains human values and the second step determines how the values

identified in the conceptual investigation can be integrated within the technology (Friedman,

2001). As we have seen in technologies that produce dangerous dark data, such as the Amazon

Echo and the Equifax data collection system, the design choices made for the technologies do not

inherently allow them to align with values of transparency and user control. In the case of the

Echo, it is imperative the device lets users know what data is being collected, when it is being

collected, and where it is being stored. This is currently not possible because a wake word may

accidentally activate the speaker and begin collecting dark data without user consent. Similarly,

Equifax customers are rarely provided information about what credit data is collected, when it is

collected, and who it is sold to. Without an alert system or an affordable method to double check

the information being collected, the identified values can never be successfully integrated into

the design.

To limit the amount of unnecessary user information being collected there clearly needs

to be a realignment of values between stakeholders, in this case corporations and consumers.

Dark data is not the result of a single specific technology; therefore, it is difficult to develop a

standard set of rules that will limit dark data collection for all devices. However, we can learn

from corporations that have improved data collection through Value Sensitive Design to develop

a general set of principles that will limit consumer privacy concerns.

### *Re-Aligning Values to Limit Dark Data*

Firstly, in the IEET's study on integrating human values within AI technologies, it was

concluded that transparency is an important but extremely difficult, if not impossible, value to

incorporate into technologies (Umbrello, 2019). Although it may be difficult to incorporate

values into the technology or device itself, the designer can require interventions at each design

stage to increase transparency via technical explainability. Frequent pauses within the design

process can allow designers a chance to communicate working principles and the reasoning

behind certain design choices with potential users. Periodic communication between

stakeholders promotes technical transparency and allows the designer to keep human values in

mind for the duration of the design process. Recently, Facebook appointed 20 people from

around the world to serve on their "Supreme Court" for speech that will issue rulings on what

kinds of posts will be allowed on their website (Ingram, 2020). As previously mentioned,

Facebook has struggled to contain the spread of misinformation to their users in recent years.

Although this issue is not a direct result of dark data, it is representative of how frequent

communication with users can promote transparency between stakeholders and ultimately benefit

both parties. Facebook, like Amazon or Equifax, does not benefit from user distrust. By allowing

this oversight board to essentially have the last word, it limits Facebook executives from

controlling speech rules and gives control to users.

Secondly, organizations can limit the amount of dark data produced by mapping out a

complete view of the data at the beginning of the design process. In a recent interview with the

Chief Marketing Officer of Bazaarvoice, Sarah Spivey, she states that the rapid rise in our ability to collect data has not been matched by our ability to filter and manage the information (Whitler, 2018). She points to a company known as KnowledgeHound that specializes in saving companies money by eliminating repeat data sets within organizations. Because there is rarely a systematic approach to data collection and organization, redundant information and saturated dark data is a common problem among companies (Whitler, 2018). As we saw in the empirical investigation, mapping out available and redundant data can offer companies the financial incentive needed for them to follow user values through the design process. This process would be especially useful for companies such as Equifax, who are managing the information of millions of Americans. Users can regain control and confidence their data is correct once the information is managed systematically.

Lastly, Spivey encourages companies to start with a clear business objective and to be purposeful in the problems they want to solve and key questions they want to ask (Whitler, 2018). By focusing on the data collection effort, companies can avoid unnecessarily collecting information and identify exactly why they need the data. Following this idea while keeping in mind the value of transparency between corporations and users will build a trusting relationship between stakeholders, limit privacy concerns, and decrease the amount of dark data that is collected.

## Conclusion

Value Sensitive Design can be an extremely useful framework to understand why many technologies result in dark data and excessive data collection in general. By conducting a conceptual, empirical, and technical investigation, we were able to conclude that the values of corporations often do not align with those of consumers. Despite transparency, user control, and

privacy being expressed as important moral values for all organizations, data-collecting

companies may regard values such as power, efficiency, and competition to be more important.

Although it may be difficult to incorporate all moral values within the technologies we

design, VSD can point to what values we deem most important. If it is too difficult to transform

these values into design requirements, companies can conduct regular interventions to increase

transparency with potential users, map out the data collection process, and be purposeful in the

problems they want to answer to bridge the value conflict gap and limit the negative

consequences of dark data. Of course, large corporations may be hesitant to adopt these measures

because of an already established data collection process that would be too complicated or

expensive to change. However, it is imperative that they look to the disastrous consequences

their technology can have on privacy and consumer safety if values of users are not considered.

## References

Briggs, B. (Ed.). (2017). Dark analytics: Illuminating opportunities hidden within unstructured data. Retrieved November, 2020, from

https://www2.deloitte.com/us/en/insights/focus/tech-trends/2017/dark-data-analyzing-unstructured-data.html

Choi, J. P., Jeon, D., & Kim, B. (2019). Privacy and personal data collection with information externalities. *Journal of Public Economics, 173*, 113-124.

doi:10.1016/j.jpubeco.2019.02.001

Cowley, S. (2019, July 22). Equifax to Pay at Least $650 Million in Largest-Ever Data Breach Settlement. Retrieved from https://www.nytimes.com/2019/07/22/business/equifax-settlement.html.

Cukier, K., & Mayer-Schoenberger, V. (2013, May/June). The Rise of Big Data. Retrieved November, 2020, from

http://cs.brown.edu/courses/cs100/lectures/readings/riseOfBigData.pdf

Fowler, G.A. (2019, May 6). Alexa has been eavesdropping on you this whole time. Retrieved from https://www.washingtonpost.com/technology/2019/05/06/alexa-has-been-eavesdropping-you-this-whole-time/.

*Friedman, B. (1996, December). Value Sensitive Design. ACM Interactions. Retrieved from* https://dl.acm.org/doi/pdf/10.1145/242485.242493?casa_token=Sgue9i2FFdgAAAAA:n9l1PfO1mhavxSBRDEl7qecMZkbPQ3XGD5072COhNvFexgWZRT5dC3tmTP9gyG8n-gIfuGpoy7Q.

Friedman, B., Kahn, P.H., Borning, A. (2001, February 12). Value Sensitive Design: Theory and

    Methods. Retrieved from

    http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.11.8020&rep=rep1&type=pdf.

Golinger, J. (1998, March). PIRG: Mistakes Do Happen: Credit Report Errors Mean Consumers

    Lose. Retrieved from http://cdn.publicinterestnetwork.org/assets/UDV-

    IZWNPfc9VDuHy9q4wA/mistakesdohappen3_98.pdf.

Hand, D. J. (2020, January 28). *Dark Data: Why What You Don't Know Matters*. *Princeton*

    *University Press.*

Ingram, D. (2020, May 6). Facebook names 20 people to its 'Supreme Court' for content

    moderation. Retrieved from https://www.nbcnews.com/tech/tech-news/facebook-names-

    20-people-its-supreme-court-content-moderation-n1201181

Kaplan, M. (2016, January 13). 6 Benefits of Using Predictive Marketing Tools, for Ecommerce.

    Retrieved from https://www.practicalecommerce.com/6-Benefits-of-Using-Predictive-

    Marketing-Tools-for

    Ecommerce#:~:text=Algorithms%20can%20predict%20a%20shopper's,best%20mix%20

    of%20marketing%20communications.

Klein, A. (2017, September 28). The real problem with credit reports is the astounding number

    of errors. Retrieved from https://www.brookings.edu/research/the-real-problem-with-

    credit-reports-is-the-astounding-number-of-errors/.

Mele, C. (2018, August 1). Data Breaches Keep Happening. So Why Don't You Do Something?

    Retrieved from https://www.nytimes.com/2018/08/01/technology/data-breaches.html.

Mori, I. (2019, May 1). The Trust Opportunity: Exploring Consumer Attitudes to the Internet of

    Things. Retrieved from https://www.internetsociety.org/resources/doc/2019/trust-

    opportunity-exploring-consumer-attitudes-to-iot/

Privacy International. (2019, April 17). The mystery of the Amazon Echo data. Retrieved from

    https://privacyinternational.org/news-analysis/2819/mystery-amazon-echo-data.

Rosalsky, G. (2020, August 4). Are Conspiracy Theories Good for Facebook. Retrieved from

    https://www.npr.org/sections/money/2020/08/04/898596655/are-conspiracy-theories-

    good-for-facebook

Sweet, K. (2017, October 6). Equifax makes money by knowing a lot about you. Retrieved from

    https://www.usatoday.com/story/money/personalfinance/2017/10/06/equifax-makes-

    money-knowing-lot-you/738824001/.

Taulli, T. (2019, October 27). What You Need To Know About Dark Data. Retrieved from

    https://www.forbes.com/sites/tomtaulli/2019/10/27/what-you-need-to-know-about-dark-

    data/?sh=55016d4d2c79

Umbrello, S. (2019). Beneficial Artificial Intelligence Coordination by means of a Value

    Sensitive Design Approach. *Big Data and Cognitive Computing, 3*,

    doi:10.3390/bdcc3010005. Retrieved from https://www.mdpi.com/2504-2289/3/1/5.

Whitler, K. A. (2018, March 17). Why Too Much Data Is A Problem And How To Prevent It.

    Retrieved from https://www.forbes.com/sites/kimberlywhitler/2018/03/17/why-too-

    much-data-is-a-problem-and-how-to-prevent-it/?sh=2bed6208755f

Xu, H., Crossler, R. E., Belanger, F. (2012, June 19). A Value Sensitive Design Investigation of

    Privacy Enhancing Tools in Web Browsers. *Decision Support Systems*, *54*, 424-433.

    doi:10.1016/j.dss.2012.06.

Ziefle, M., Rocker, C., Holzinger, A. (2011, July). Medical Technology in Smart Homes:

Exploring the User's Perspective on Privacy, Intimacy, and Trust. *Institute of Electrical*

*and Electronics Engineers.* Retrieved from

https://ieeexplore.ieee.org/document/6032273/citations#citations.