

**Reliable Analytics for Disease Indicators: Leveraging Smart Devices to Predict Health**  
(Technical Paper)

**Instagram, Amazon, and Machine Learning: Ethical Implications of Collecting and  
Analyzing Commercial User Data**  
(STS Paper)

**A Thesis Prospectus Submitted to the**

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

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

Tucker Wilson  
Spring, 2020

Technical Project Team Members

Erin Barrett  
Cameron Fard  
Hannah Katinas  
Charlie Moens  
Lauren Perry  
Blake Ruddy  
Shalin Shah  
Ian Tucker

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

## Introduction

Artificial intelligence and machine learning is perhaps the most rapidly expanding field of computer science today, and one of the most ethically impactful application of these technologies is on personal data. Most internet-based technologies generate large amounts of user data through logging our online activities, searches, views, and interactions. The amount of data we generate is multiplied even further when considering the popularity of the smartphone or other smart technologies, such as the Amazon Echo, which have the capability to collect location, motion, health, photo, and audio data (Zhou & Gurrin, 2012). As users generate this massive amount of data on themselves, companies are rushing to use their databases to create insights.

The uses of this data are many, but one clear application of user data is in the field of digital advertising. Companies, such as Google, are logging user activity on owned websites, such as YouTube, and browsers, such as Google Chrome, and using that data to determine a user's interests. Social media sites, such as Facebook, Twitter, and Instagram, also employ interaction data between users, such as follower lists, to make predictions on a user's interests based on the interests of users they interact with. Once a user's interests have been determined, advertising services can predict which products a user might be interested in purchasing and display targeted online advertisements accordingly.

Another field employing user data, particularly data generated by smartphones and smart wearable devices, is that of healthcare. Apple's website openly advertises the health implications of the Apple Watch, for example, which stores heart rate data, medical ID's for first responders to access, and records from healthcare facilities for users to access (Apple, Inc., n.d.). iPhones

are also capable of tracking additional health data, such as sleep patterns and daily exercise. With this amount of data available at all times, it may be possible to predict when a user is sick without a doctor visit, or even predict when they are at higher risk of disease before showing any noticeable symptoms.

In both of these cases, user-generated data can have a significant impact on a user's life, and fully exploring the potential uses of this data and the implications of doing so is necessary to make informed decisions about the future of user predictive technologies. The following technical topic will explore and improve upon methods of using user-collected smartphone data to predict illness with applications to the military, while the STS topic seeks to understand the landscape of user-generated data from social media and smart devices and explore the potential dangers of collecting such data.

### **Technical Topic**

The technical project is a part of ongoing research conducted for the Defense Advanced Research Projects Agency (DARPA) to design and develop reliable disease detection analytics through data collected from smartphones, specifically applied to military personnel stationed in combat zones. The ultimate goal of the research is to create “a mobile application that passively assesses a warfighter's readiness immediately and over time” (Patel, n.d.). By building predictive health analytics that utilize smartphone sensors, the onset of illnesses, concussions, or even mental health issues can be noticed in real time. In the current stage of research, the technical team will research and develop strategies to optimize the tradeoff between data collection frequency and battery life. By gaining a better sense of these limitations, accurate predictive models can be built without unneeded noise in the data.

Mobile sensing data used in this research will be collected through the Sensus Application. This application, developed at the University of Virginia (UVA) (Lockheed Martin & Advanced Technology Laboratories, 2017), uses “event-driven architecture that triggers actions in response to changes to the device or network state.” This data will be utilized to create context recognition models, which determine what ambulatory state the user is in, like walking, running, or sitting. Additionally, the Sensus application will push surveys as notifications to participants’ mobile phones to create additional context around the data collected. These surveys will ask questions about the user’s activities immediately before answering the survey, such as the user’s location, length of activity, phone position, and more. The team can then build a model off of this data to better predict when a user’s phone is actually in use. The main goal of this prediction is to deploy a model of *adaptive sensing*, in which a user’s phone can determine when it is being used and turn sensors off and on accordingly.

The technical project group consists of nine undergraduate Systems Engineering students. The team is divided into three subteams: the Data Modeling Team, the Data Visualization Team, and the Data Collection Team. These teams were constructed for the current needs of the project, and are subject to change and overlap depending on the need in each area. The Data Modeling Team will work to prove the efficacy of adaptive sensing in an attempt to find a balance between data collection and battery usage. Ultimately, this team will develop an algorithm as a potential alternative to the adaptive sensing model currently being used. The Data Visualization Team will make significant improvements to the web-based visualization platform used by the researchers to increase understanding and context of the data they are collecting. Improvements to this platform will allow better insights to be easily accessible. The Data Collection Team is

designated to complete the IRB so that the data collection among the student cohort can begin. Once the IRB is completed and approved, the team will be responsible for organizing the participants in the study.

At the end of the study, the team will deliver a recommendation for smartphone data collection that effectively accounts for a user's battery life and critical predictive data and a recommendation for intuitive data visualizations for the researchers' web platform. The technical project is funded through a grant provided by DARPA. Additional resources include test phones and desktop computers to run software and view data. The technical project will produce a conference paper for the Systems Information Engineering Design Symposium (SIEDS) that will take place in May, 2020.

### **STS Topic**

The proposed STS project will focus on the first use case of machine learning algorithms, user-generated data collection and analysis for the purposes of online advertising. While this research will cover user data collection and predictive technologies generally, it will focus on two specific use cases: Instagram and the Amazon Echo. Both of these use cases heavily employ predictive technologies to perform user advertising, and particularly with the Amazon Echo the 'always-on' nature of these technologies has serious implications for user privacy.

First, the current landscape of data collection. Social media sites and web browsers collect user data through site activity as well as information coming from their devices, such as geographic and time data. If a user has an account on a site, demographic data such as age, sex, race, etc., is added into this collection (Malthouse et al., 2018). Both smartphones themselves and smart home devices collect data in analogous ways, primarily relying on microphone

recordings, which are either used directly, as with the Amazon Echo responding to audio commands, or stored for potential future use (Rediger, 2017). On social media sites in particular, another technique known as sentiment analysis, which analyzes user text posts to determine the post's positive, negative, or neutral sentiment, can be used to describe a user's feelings towards a particular topic of interest or product (Katsurai & Satoh, 2016). These data sets combine to create an incredibly rich data landscape on individual users. On Instagram, for example, this information creates a unique advertising opportunity. Instagram thrives off of constant relatable content creation. Companies are recognizing such an approach and using user data to target their customers, either by sponsoring popular Instagram personalities that are well-liked by their customer base or by promoting a branding image that is popular among their customers (Carah & Shaul, 2015).

However, this constant collection of rich data poses security and privacy issues. The Amazon Echo in particular has been the subject of many safety concerns. For example, a 2017 study on breaking the Echo's security (Haack et al.) showed the ease at which Amazon user security PINs could be guessed, and stated that there are perceivable cases where "listeners may be able to recover personal details, including payment information" from the information the Echo is constantly sending to Amazon's servers. When critical user information is involved, including but not limited to credit card information, such lax security presents major privacy issues. Additionally, government organizations have a history of monitoring cell phones and smart devices (Rediger, 2017), and audio recordings from the Amazon Echo have already been used as evidence in law enforcement investigations, even though the Constitutionality of such use has not yet been determined (Jackson & Orebaugh, 2018). Finally, these types of

technologies can be dangerous by design. An example is Russian interference in the 2016 US Presidential Election by running targeted advertisements on Facebook, but other examples exist. For example, through data and sentiment analysis on Instagram posts supplemented by user surveys on said posts, a 2017 experiment with 166 participants was able to predict clinical depression among Instagram users with higher accuracy (lower false diagnosis rate) than a practitioner (Reece & Danforth).

These implications are coupled with the confounding issue that most users are unaware of or apathetic to the extent to which their data is being collected and leveraged. While a 2014 study of Instagram users (Talib et al.) reported that 67% of users surveyed “know that SNS [Social Network Site] sells personal user information to other organization[s]”, 63% do not consider Instagram’s privacy policy when using the site. Similarly, while 68% of Amazon Echo users surveyed from a 2018 study (Manikonda et al.) reported some privacy concerns about the device, only 6% mentioned that “the device should respect their privacy.” These disconnects imply lack of knowledge or concern from users, either from general apathy or purposeful obscurity from the companies collecting their data.

In studying a large, widely distributed network of companies, algorithms, databases, governments, and users, this project will employ Actor-Network Theory to map these complex relations and model their change over time. Actor-Network Theory, or ANT, is a method of formally describing a complex network of stakeholders and contributors to a technological system. It consists of defining actors, which can be people, companies, technologies, or any other entity that affects the system, and that exist within a network, with unspecified and constantly changing relationships connecting the actors. Defining intermediaries, the languages through

which actors communicate, is how ANT seeks to explain the complex relationships within the network. A common critique of this theory is that it is purely descriptive and does not seek to explain the impact any actors have or why the network takes its current form. However, in the STS project ANT will only be used as a method of mapping the landscape of user-generated data, not as a framework for explaining the actions of any stakeholders. Therefore, this limitation will not be an issue. Instead, the theory of Technological Momentum will be employed to study how these systems of data gathering and executing on that data grew out of small-scale experiments and will be much harder to change in the modern era. Technological Momentum is a theory that technologies, when first created, are done so because of and subsequently shaped by the society they were created in and the stakeholders that created and used them. However, at some point, a technology grows large enough that further shaping is incredibly difficult and the technology begins to shape the society. An example is the invention of man-made electricity that transformed into the electrical grids of major cities. One critique of this framework comes from David Nye (2007), who stated that “cultures select and shape technologies, not the other way around,” and that no technology is ever “taking humanity somewhere in particular.” Nye argues that no technology, no matter the size, is ever deterministic and unshaped by humans. In order to recognize this critique, the STS topic will seek to use Technological Momentum to explain why changes to the user-generated data landscape will be more difficult because of its size and avoid the claim that change is impossible.

### **Research Question and Methods**



The main question being studied in this research will be whether user data collection and analysis poses security, privacy, or safety risks to the users, and whether there exist potential methods to mitigate such damage.

This paper will primarily employ documentary research on primary and secondary sources to perform these explorations. These sources include studies performed on users of this technology, court cases raised against companies that employ these data collection techniques (Rediger, 2017), the European Union's recent update to its cybersecurity standards for private companies (2016), and case studies on instances where these types of data collection methods have been exploited. These case studies include a showcase of how data gathered from Amazon devices at home and at work can be used to track movements of those near the devices (Do et al., 2018), as an example. These sources should provide a view of the landscape of user data collection for machine learning, and hopefully illustrate problems with the current systems. In addition to these documentary sources, a broad survey of UVA students will be used to explore a selection of user's perspectives on what data is being collected on them and how it is being used. This survey should strengthen the claim that users are generally unaware of the extent to which their data is being collected and the risks they have accepted by using these technologies.

### **Conclusion**

With a showcasing of healthcare data being used to predict disease among users in the technical portion, and many examples and implications of user data analysis and exploitation in the STS portion, this thesis should provide a strong foundation for policy recommendations on improving the user privacy and data security while maintaining the data's predictive capabilities. While policy recommendations are not the main focus of this thesis, it should still deliver

principles by which to construct policies in an ethical way with adequate consideration given to the users these technologies can impact.

## References

- Apple, Inc. (n.d.). Healthcare - Apple Watch. Retrieved from <https://www.apple.com/healthcare/apple-watch/>.
- Carah, N., & Shaul, M. (2015). Brands and Instagram: Point, tap, swipe, glance. *Mobile Media & Communication*, 4(1), 69–84. doi: 10.1177/2050157915598180.
- Do, Q., Martini, B., & Choo, K.-K. R. (2018). Cyber-physical systems information gathering: A smart home case study. *Computer Networks*, 138, 1–12. doi: 10.1016/j.comnet.2018.03.024.
- General Data Protection Regulation (2016) *Official Journal* L119, 4 May 2016, p. 1-88
- Haack, W., Severance, M., Wallace, M., & Wohlwend, J. (2017). Security Analysis of the Amazon Echo. Retrieved from <https://pdfs.semanticscholar.org/35c8/47d63db1dd2c8cf36a3a8c3444cdeee605e4.pdf>.
- Jackson, C., & Orebaugh, A. (2018). A study of security and privacy issues associated with the Amazon Echo. *International Journal of Internet of Things and Cyber-Assurance*, 1(1), 91. doi: 10.1504/ijitca.2018.10011257.
- Katsurai, M., & Satoh, S. (2016). Image sentiment analysis using latent correlations among visual, textual, and sentiment views. *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. doi: 10.1109/icassp.2016.7472195.
- Lockheed Martin & Advanced Technology Laboratories (2017). DARPA warfighter analytics using smartphones for health (WASH) ReADI technical section. Cherry Hill, NJ: Lockheed Martin and Advanced Technology Laboratories.
- Malthouse, E. C., Maslowska, E., & Franks, J. U. (2018). Understanding programmatic TV

- advertising. *International Journal of Advertising*, 37(5), 769–784. doi: 10.1080/02650487.2018.1461733.
- Manikonda, L., Deotale, A., & Kambhampati, S. (2018). What's up with Privacy?: User Preferences and Privacy Concerns in Intelligent Personal Assistants. *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society - AIES 18*. doi: 10.1145/3278721.3278773.
- Nye, D. E. (2007). Not Just One Future. In *Technology matters: Questions to Live With*. Cambridge, MA: MIT Press.
- Patel, T. (n.d.). Warfighter Analytics using Smartphones for Health (WASH). Retrieved from <https://www.darpa.mil/program/warfighter-analytics-using-smartphones-for-health>.
- Rediger, A. M. (2017). Always-Listening Technologies: Who Is Listening and What Can Be Done About It. *Loyola Consumer Law Review*, (2), 229–252.
- Reece, A. G., & Danforth, C. M. (2017). Instagram photos reveal predictive markers of depression. *EPJ Data Science*, 6(1). doi: 10.1140/epjds/s13688-017-0118-4.
- Talib, S., Razak, S. M. A., Olowolayemo, A., Salependi, M., Ahmad, N. F., Kunhamoo, S., & Bani, S. K. (2014). Perception analysis of social networks privacy policy: Instagram as a case study. *The 5th International Conference on Information and Communication Technology for The Muslim World (ICT4M)*. doi: 10.1109/ict4m.2014.7020612.
- Zhou, L. M., & Gurrin, C. (2012). A survey on life logging data capturing. *SenseCam Symposium*. Retrieved from <http://doras.dcu.ie/17533/>.