

**Image Processing Tool for Quantifying Cell Populations in Fibrotic Cardiac Tissue**

(Technical Paper)

**Machine Learning in Healthcare: How Strict Data Collection Policies Lead to Misrepresentational Models**

(STS Paper)

A Thesis Prospectus

In STS 4500

Presented to

The Faculty of the

School of Engineering and Applied Science

University of Virginia

In Partial Fulfillment of the Requirements of the Degree

Bachelor of Science in Biomedical Engineering

By

Jakub Lipowski

November 1, 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

**ADVISORS**

Shayn Peirce-Cottler, Department of Biomedical Engineering

Bryn Seabrook, Department of Engineering and Society

## **Introduction:**

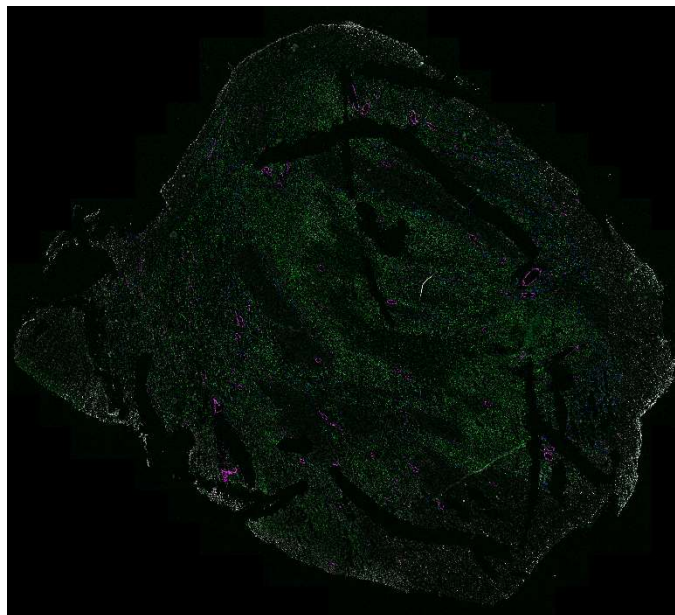
Machine learning's predictive prowess has been slowly seeping in humans' decision-making systems, even ones that decide between life and death. Healthcare, the United States' largest industry is certainly no exception to implementing artificial intelligence. Algorithms are continuously being trained for purposes of drug discovery, disease diagnosis, outbreak prediction, electronic record keeping, and imaging. Regardless of the case, outsourcing labor to a machine learning algorithm limits the autonomy that a human has to make an informed decision. In some cases, like predicting the number of jelly beans in a jar, this is harmless, but in others, like deciding whether a person should be prescribed a drug to curb an autoimmune disease, the results drastically alter a person's quality of life. As a result, the interface of technology and healthcare will always raise ethical dilemmas. The backbone of machine learning is the data that is used to train models so that they generalize to new data points. It is important to ensure that data for artificial intelligence applications in healthcare are representative of entire populations so that the model can generalize well to anyone within a target society. If the model does not generalize well to certain demographics, then misleading predictions could be fatal. The goals of this prospectus are two-fold: a technical project to introduce a novel image processing tool that is an instance of machine learning to draw conclusions about the spatiotemporal dynamics of cell populations within cardiac tissue that has suffered a loss of blood flow and a complementary research paper to use the United States' diverse population to study the relationship between data collection practices and generalizability of resulting machine learning models.

## **Technical Section:**

Machine learning can be used to accelerate research and breakthroughs by analyzing data faster than a human ever could; image processing is a prime instance. The image processing tool

is being developed in tandem with three researchers at Dr. Shayn Peirce-Cottler's microvascular engineering laboratory in the University of Virginia's Department of Biomedical Engineering and the Center for Advanced Biomanufacturing. The tool is being developed in the context of myocardial infarctions (MI), which refer to the loss of blood flow to cardiac tissue, commonly known as a heart attack (*Heart Attack | Johns Hopkins Medicine*, n.d.). A deep understanding of MI is important because heart attacks account for 800,000 deaths in the United States and are expected to create a global economic burden of 918 billion USD by 2030 (Travers et al., 2016). When MI first occurs, cells in cardiac tissue called cardiomyocytes induce the arrival of monocytes, a type of immune cell, at the region of the infarct via perfused blood vessels. Next, monocytes differentiate into macrophages. Macrophages are immune cells that, among other things, secrete chemicals to attract fibroblasts (Duncan et al., 2020). Fibroblasts are cells that secrete proteins that stiffen the cardiac tissue, causing an irregular heartbeat, chamber dilation, and potentially heart failure (Kong et al., 2014). For this reason, it is imperative to understand the spatiotemporal dynamics between cardiomyocytes, monocytes, macrophages, and fibroblasts in the days following MI not only in the entire region of the infarct, but also around key structures such as blood vessels where these cells are originating (Hsieh et al., 2006). This research could lead to identification of new targets for therapeutic intervention to enhance cardiac tissue regeneration and could inform future computational models of multi-cell interactions. The goal of the image processing tool is to automate the analysis of spatiotemporal dynamics between the aforementioned cell populations. The tool will implement supervised machine learning to classify cell types in stained images taken with a microscope and will be implemented through a user-friendly graphical user interface (GUI).

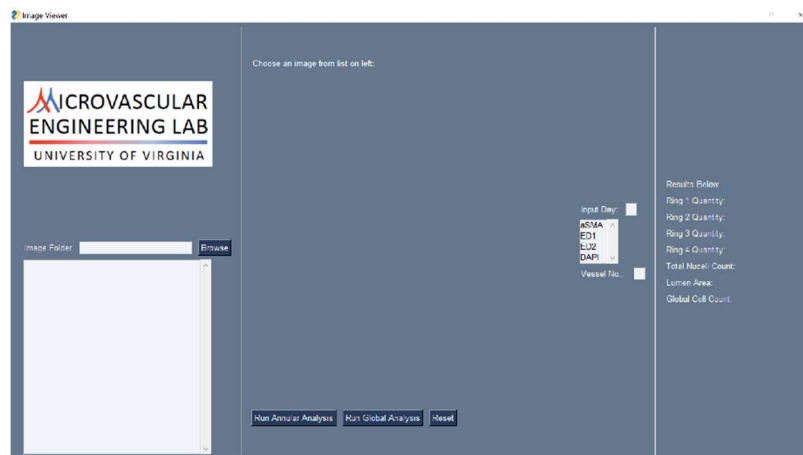
The image processing tool depends on images of the heart that have already been procured using the following procedure. First, a heart attack was surgically induced in rats. These injured hearts were removed from the rats 1, 3, 4, 5, and 6 days following MI. The heart walls were then sliced into thin micrographs that were imaged. Prior to imaging, the sections of cardiac tissue were immunostained, which means that fluorescently labeled biological markers were used to tag specific cell types. The montaged micrographs (about 12mm by 11 mm in size) of entire immunostained cross-sections of cardiac tissue (generously provided by collaborators at the University of Virginia, Laura Caggiano, Ph.D. and Jeffrey Holmes, M.D., Ph.D.) were obtained with a 10x objective on a Leica Thunder microscope. Figure 1 shows an image of one of these micrographs.



**Figure 1:** Image of immunostained cardiac tissue from a rat that has suffered MI, taken four days after the injury.

The image processing program will be developed using Python (3.7) and Ilastik (1.4.0b5), a Python based supervised machine learning tool for image segmentation. The program will be implemented to analyze images of immunostained cardiac tissue obtained 1, 3, 4, 5, and 6 days following MI. Ilastik will be used to classify different cell types based on the

stain using image segmentation algorithms. A user-friendly GUI will be developed, which will prompt the user to select an image, count different cell types, compute cell densities, and perform a novel “annular analysis” to quantify macrophage extravasation from blood vessels. Annular analysis will allow the user to select any region in the image, outline a blood vessel, and collect information about the types, quantities, and densities of cells occurring at several distances away from the vessel. These distances will be demarcated by equidistant rings radiating from the initial user-defined outline of the vessel. All the data collected from analysis of an image will automatically be written to an Excel spreadsheet. For the annular analysis, the GUI will also save a bar plot of the cell densities surrounding the preselected blood vessel. SQL, a database query tool, will be used to organize tables of data based on cell types and variables of interest. These data will then be used to develop heatmaps of cell densities across entire tissue sections. A preliminary model of the GUI has already been proposed and is depicted in Figure 2. Ultimately, the technical project will deliver a functional image processing software accompanied by a technical report describing its functionalities.



**Figure 2:** Proposed graphical user interface for image processing tool.

## **STS Topic:**

Data, an incredibly general term that refers to pieces of information that unlock nonobvious relationships, push research forward, and help solve society's most pressing challenges. Predictive machine learning algorithms are heavily dependent upon reliable data in order to deliver insights about the future. Collection of such data is critical to a successful machine learning model, mainly because data needs to be unbiased and should represent the populations for which the algorithm is applied to. A model trained using data for a specific demographic cannot generalize to other demographics without a high probability of producing fatal errors (Gianfrancesco et al., 2018). Despite a need for representative data samples to train and validate machine learning models in healthcare, data privacy laws around the world are extremely variable, so neither the same type nor volume of data can be collected from all regions of the globe, severely limiting the potential of machine learning models. The purpose of this research will be to reconcile the need for representative training data during machine learning model development with the limitations that the United States' public policies place on data collection and privacy.

The United States' Health Insurance Portability and Accountability Act (HIPAA) applies data protection to personally identifiable health-related information stored by healthcare professionals, but does not extend this protection to other sources of health data, like businesses developing devices and phone applications. On the other hand, the European Union (EU) has introduced the European General Data Protection Regulation in April of 2016, which applies to all EU citizens and mandates informed consent requirements for collection of self-identifying data, regardless of where a citizen is in the world (Vayena et al., 2018). Meanwhile, over half of the African continent is yet to introduce data privacy laws (*Slowly but Surely, Data Protection*

*Regulations Expand throughout Africa*, n.d.). Bias caused by such variable data collection policies lends itself to using data that may be more skewed on the basis of geography as well as the type of source collecting the data. This in turn alters the reliability of machine learning models trained using this data. Although examining this topic on a global scale is beyond the scope of this research, the United States serves as a great case for studying usage of data in machine learning models due to its diverse population. There are three key topics that must be understood to determine the relationship between data collection and reliability of machine learning models within the United States: unequal accessibility to healthcare, discrepancies in data collection practices, and the concept of risk analysis.

It is undisputable that the United States has implemented a series of discriminatory regulations in its history, disadvantaging African-Americans, Hispanics, and other racial groups. Although the Civil Rights Act of 1964 ended legal discrimination, there are still disparities noticeable within the healthcare system, which are strongly connected to socioeconomic inequalities (Data et al., 2003). Of the 27.5 million uninsured people within the U.S., 45 percent cite cost as the cause for being uninsured. Additionally, in the case of the 42 million African Americans in the U.S., the average annual cost for healthcare premiums is about 20 percent of their average household income (*Racism, Inequality, and Health Care for African Americans*, 2019). Such inequalities lead to a decreased access to healthcare, which in turn lead to a misrepresentation of entire populations in aggregated data that can go on to be used in large-scale precision medicine models. Socioeconomic inequalities are a key instigator of discrepancies in training data used for machine learning models, but the methods and policies that dictate how health data is obtained within the United States are flawed as well.

HIPAA's Privacy Rule protects personally identifiable health information obtained by covered entities, which include health plans, providers, clearinghouses, and business associates who collect data on behalf of the aforementioned covered entities (Rights (OCR), 2008). However, individuals can access information that is "deidentified," meaning that it has removed 18 specified personal identities from a given dataset (Nass et al., 2009). Censorship of such information can deem a dataset to be useless for research purposes, including development of machine learning models. Since this issue has been brought to the attention of Health and Human Services, a compromise was developed to allow entities to only strip direct identifiers, like a Social Security Number, and require permission from subjects to disclose data for purposes beyond research (Price & Cohen, 2019). Still, this compromise only brought down the requirement for deidentification from 18 identities to just 16, making it hard to develop useful algorithms in healthcare.

HIPAA policies, along with cautious use of machine learning models in healthcare, can be explained through the lens of risk analysis. Risk analysis is a key framework for reconciling the discrepancy between misrepresentational training data and strict data collection policies within a biased healthcare system. Risk analysis is a systemic way of dealing with specific events and their perceived consequences and hazards. As modernization introduces novel technologies, policies are often initiated to mediate any risks that could set social structures back. However, in the process of risk analysis, society examines itself, changing itself within the process, a concept coined "reflexivity" (*Beck's Sociology of Risk: A Critical Assessment - Anthony Elliott, 2002, n.d.*). Risk analysis is criticized for defending harmful policies, which is a criticism that will be considered in this research (Eid, 2003).



## **Methodologies:**

Research Question: How can the discrepancy between a need for representative training data and strict data collection polices within an unequal healthcare system be reconciled to minimize the risks of biased machine learning models?

The research can be partitioned into two overarching tasks:

1. Researching collection of health data in the United States. Key terms including HIPAA, Privacy Rule, electronic health records, and privacy will be used to learn about data collection practices within the United States. Since this topic is closely related to a national policy, there are many online resources that can be analyzed independently, so no interviews will be necessary. Gathering information about data collection within the United States will help to identify imperfections among data collection policies.
2. Performing case studies of machine learning algorithms in healthcare. A case refers to any instance where machine learning is used in healthcare. The research will pay specific attention to collection of data within each case along with the results of the machine learning model. Cases will be collected from online resources and faculty in the Biomedical Engineering Department at the University of Virginia. These cases will be used to identify flawed machine learning algorithms in healthcare to be used as supporting evidence for the research.

## **Conclusion:**

The research proposed in this document consists of two projects. The technical project aims to develop, validate, and publish an image processing tool that automates the analysis of spatiotemporal dynamics of cell populations in cardiac tissue that has suffered a loss of blood

flow. The sociotechnical project will reconcile the discrepancy between a need for representative training data and strict data collection policies within an unequal healthcare system. By virtue of being automated, the image processing software will increase the rate at which researchers are able to generate conclusions and hypotheses. The tool will also quantify presently unknown relationships between cell populations in scarred cardiac tissue in order to help identify targets for therapeutic intervention. The result of the sociotechnical project will identify potential areas for change within current data collection policies in the United States. The goal is to find a compromise between strict data collection laws applied to a diverse population and the need for data that represents the diverse American population.

## Works Cited

- Beck's Sociology of Risk: A Critical Assessment—Anthony Elliott, 2002*. (n.d.). Retrieved October 31, 2021, from <https://journals.sagepub.com/doi/10.1177/0038038502036002004>
- Data, N. R. C. (US) P. on D. C. of R. and E., Melnick, D., & Perrin, E. (2003). IMPROVING RACIAL AND ETHNIC DATA ON HEALTH. In *Improving Racial and Ethnic Data on Health: Report of a Workshop*. National Academies Press (US). <https://www.ncbi.nlm.nih.gov/books/NBK222062/>
- Duncan, S. E., Gao, S., Sarhene, M., Coffie, J. W., Linhua, D., Bao, X., Jing, Z., Li, S., Guo, R., Su, J., & Fan, G. (2020). Macrophage Activities in Myocardial Infarction and Heart Failure. *Cardiology Research and Practice*, 2020, 4375127. <https://doi.org/10.1155/2020/4375127>
- Eid, M. (2003). Reflexive Modernity and Risk Society. *International Journal of the Humanities*, 1. <https://doi.org/10.18848/1447-9508/CGP/v01/58162>
- Gianfrancesco, M. A., Tamang, S., Yazdany, J., & Schmajuk, G. (2018). Potential Biases in Machine Learning Algorithms Using Electronic Health Record Data. *JAMA Internal Medicine*, 178(11), 1544–1547. <https://doi.org/10.1001/jamainternmed.2018.3763>
- Heart Attack | Johns Hopkins Medicine*. (n.d.). Retrieved October 31, 2021, from <https://www.hopkinsmedicine.org/health/conditions-and-diseases/heart-attack>
- Hsieh, P. C. H., Davis, M. E., Lisowski, L. K., & Lee, R. T. (2006). Endothelial-Cardiomyocyte Interactions in Cardiac Development and Repair. *Annual Review of Physiology*, 68, 51–66. <https://doi.org/10.1146/annurev.physiol.68.040104.124629>

- Kong, P., Christia, P., & Frangogiannis, N. G. (2014). The pathogenesis of cardiac fibrosis. *Cellular and Molecular Life Sciences: CMLS*, 71(4), 549–574.  
<https://doi.org/10.1007/s00018-013-1349-6>
- Nass, S. J., Levit, L. A., Gostin, L. O., & Rule, I. of M. (US) C. on H. R. and the P. of H. I. T. H. P. (2009). HIPAA, the Privacy Rule, and Its Application to Health Research. In *Beyond the HIPAA Privacy Rule: Enhancing Privacy, Improving Health Through Research*. National Academies Press (US). <https://www.ncbi.nlm.nih.gov/books/NBK9573/>
- Price, W. N., & Cohen, I. G. (2019). Privacy in the Age of Medical Big Data. *Nature Medicine*, 25(1), 37–43. <https://doi.org/10.1038/s41591-018-0272-7>
- Racism, Inequality, and Health Care for African Americans*. (2019, December 19). The Century Foundation. <https://tcf.org/content/report/racism-inequality-health-care-african-americans/>
- Rights (OCR), O. for C. (2008, May 7). *The HIPAA Privacy Rule* [Text]. HHS.Gov. <https://www.hhs.gov/hipaa/for-professionals/privacy/index.html>
- Slowly but surely, data protection regulations expand throughout Africa*. (n.d.). Retrieved October 3, 2021, from <https://iapp.org/news/a/slowly-but-surely-data-protection-regulations-expand-throughout-africa/>
- Travers, J. G., Kamal, F. A., Robbins, J., Yutzey, K. E., & Blaxall, B. C. (2016). Cardiac Fibrosis. *Circulation Research*, 118(6), 1021–1040.  
<https://doi.org/10.1161/CIRCRESAHA.115.306565>
- Vayena, E., Blasimme, A., & Cohen, I. G. (2018). Machine learning in medicine: Addressing ethical challenges. *PLoS Medicine*, 15(11), e1002689.  
<https://doi.org/10.1371/journal.pmed.1002689>