

UCLA

UCLA Electronic Theses and Dissertations

Title

Computational Methods to Inform Healthcare Decisions at Individual and Population Levels

Permalink

<https://escholarship.org/uc/item/9b02m009>

Author

Rakocz, Nadav

Publication Date

2023

Peer reviewed|Thesis/dissertation

UNIVERSITY OF CALIFORNIA  
Los Angeles

Computational Methods to  
Inform Healthcare Decisions at  
Individual and Population Levels

A dissertation submitted in partial satisfaction  
of the requirements for the degree  
Doctor of Philosophy in Computer Science

by

Nadav Rakocz

2023

© Copyright by  
Nadav Rakocz  
2023

## ABSTRACT OF THE DISSERTATION

Computational Methods to  
Inform Healthcare Decisions at  
Individual and Population Levels

by

Nadav Rakocz

Doctor of Philosophy in Computer Science

University of California, Los Angeles, 2023

Professor Sriram Sankararaman, Chair

At a time when computational power and big data are driving revolutionary changes across various sectors, the healthcare industry is on the verge of a significant transformation. The integration of sophisticated computational techniques promises not only to enhance medical decision-making but also to fundamentally change the delivery of healthcare services. However, the sector grapples with challenges like the underutilization of its abundant data in clinical guidelines, which tend to rely on oversimplified, population-based methods, and the scarcity of annotated and labeled datasets in medical contexts. In this dissertation, we address the challenges impeding the full exploitation of computational capabilities in healthcare. The initial chapters are dedicated to enhancing decision-making at an individual level. Specifically, Chapter One addresses the classification challenges in 3D medical imaging, a task hindered by sparse and labor-intensive annotation processes. Chapter Two introduces a novel approach that leverages transformer models to augment and personalize clinical practice guidelines, thereby enhancing their relevance and applicability to individual patient care. Subsequent

chapters pivot to a population-level perspective, presenting computational techniques that analyze varied datasets, ranging from social media data to records of the COVID-19 pandemic. These methods attempt to identify causal mechanisms and quantify uncertainty to support decision-making that is both data-driven and reliable.

The dissertation of Nadav Rakocz is approved.

Wei Wang

Bogdan Pasaniuc

Eran Halperin

Eleazar Eskin

Sriram Sankararaman, Committee Chair

University of California, Los Angeles

2023

*To my parents Moti and Zehava,  
for instilling in me the virtues of humility and the drive to always strive for excellence.*

*And to my wife, Hili,  
for always being up for an adventure.*

## TABLE OF CONTENTS

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Scope of Research	1
1.2	Contributions and Overview	3
<b>2</b>	<b>Automated identification of clinical features from sparsely annotated 3-dimensional medical imaging</b>	<b>5</b>
2.1	Introduction	5
2.2	Results	8
2.2.1	The SLIVER-net model	8
2.2.2	AMD-Related biomarker prediction	9
2.2.3	Comparison of SLIVER-net to state of the art deep learning approaches	10
2.2.4	Comparison of SLIVER-net with specialist clinician assessments	11
2.2.5	Effect of sample size on the model performance	13
2.2.6	Identifying traces of biomarkers outside of the macula	19
2.2.7	Transfer learning improves model performance	20
2.2.8	The tradeoff between quantity and quality of external data	22
2.2.9	Robustness to the number of slices available in each volume	22
2.3	Discussion	23
2.4	Methods	26
2.4.1	Data	26
2.4.2	3D CNNs	29
2.4.3	SLIVER-net Architecture	29



2.4.4	Training . . . . .	31
2.4.5	Transfer learning . . . . .	32
2.4.6	Model evaluation . . . . .	33
2.4.7	Model explainability . . . . .	33
2.4.8	Simulating model performance with different acquisition parameters .	34
<b>3</b>	<b>Augmenting Clinical Practice Guidelines Using Reinforcement Learning and Causal Transformers . . . . .</b>	<b>35</b>
3.1	Introduction . . . . .	35
3.2	Results . . . . .	38
3.2.1	Data Summary . . . . .	39
3.2.2	Evaluating expected outcomes . . . . .	40
3.2.3	Sensitivity of WIS for behavioral policy estimation . . . . .	41
3.2.4	Permutation testing of temporal importance . . . . .	43
3.3	Methods . . . . .	44
3.3.1	Data . . . . .	44
3.3.2	Modelling . . . . .	45
3.3.3	Model Input . . . . .	45
3.3.4	Architecture and Training . . . . .	47
3.3.5	Evaluation . . . . .	49
3.3.6	Permutation Testing . . . . .	50
3.4	Discussion . . . . .	51
<b>4</b>	<b>The heterogeneous effects of social support on the adoption of Facebook’s vaccine profile frames feature . . . . .</b>	<b>54</b>

4.1	Introduction . . . . .	54
4.2	Results . . . . .	56
4.2.1	VPF adoption exhibits a pattern of complex diffusion . . . . .	56
4.2.2	Pre-existing openness to vaccines requires significantly less social proof for adoption . . . . .	57
4.2.3	Social proof from stronger ties has a greater effect on adoption . . . . .	60
4.2.4	Influencers showed limited effect on adoption . . . . .	61
4.2.5	Modeling the effects of social proof on adoption . . . . .	61
4.2.6	A randomized field experiment provides causal support for social proof and tie-strength effects on the adoption . . . . .	65
4.2.7	Causal machine learning reveals additional heterogeneous treatment effects . . . . .	67
4.2.8	Backfire effects of VPF adoption . . . . .	68
4.3	Methods . . . . .	70
4.3.1	VPF adoption and promotional exposure data (RQ1) . . . . .	70
4.3.2	Influencers' matching and difference in differences (RQ1) . . . . .	71
4.3.3	Logistic regression model (RQ1) . . . . .	73
4.3.4	Randomized field experiment (RQ1) . . . . .	74
4.3.5	Causal Forest for heterogeneous treatment effects (RQ1) . . . . .	76
4.3.6	Backfire effects (RQ2) . . . . .	76
4.3.7	Additional data . . . . .	77
4.4	Discussion . . . . .	80

**5 A Statistical Model for Quantifying the Needed Duration of Social Distanc-**

<b>ing for the COVID-19 Pandemic . . . . .</b>	<b>86</b>
5.1 Introduction . . . . .	86
5.2 Results . . . . .	89
5.2.1 Estimates of $t_{\text{end}}$ from region-specific parameter estimates . . . . .	89
5.2.2 Estimates of $t_{\text{end}}$ using a Bayesian framework . . . . .	92
5.2.3 Sensitivity analysis. . . . .	95
5.3 Methods . . . . .	95
5.3.1 The SEIR Model . . . . .	95
5.3.2 A Bayesian hierarchical model for parameter estimation across multiple locations . . . . .	98
5.3.3 Application to predict the end of social distancing . . . . .	101
5.4 Discussion . . . . .	101
 <b>A Supplementary Material - The heterogeneous effects of social support on the adoption of Facebook’s vaccine profile frames feature . . . . .</b>	 <b>105</b>
 <b>B Supplementary Material - A Statistical Model for Quantifying the Needed Duration of Social Distancing for the COVID-19 Pandemic . . . . .</b>	 <b>109</b>
B.1 Ordinary differential equations . . . . .	109
B.2 Code availability . . . . .	109
 <b>References . . . . .</b>	 <b>110</b>

## LIST OF FIGURES

2.1	SLIVER-net performance. SLIVER-net (dark blue) was compared with a 3D CNN backbone approach (light blue) and 2D CNN (gray). SLIVER-net significantly outperformed both the 3D CNN and 2D CNN in identifying each biomarker in terms of area under the ROC (AUROC) and area under the Precision-Recall curve (Precision-Recall AUC). Top: Precision-Recall AUC for each biomarker. Bottom: ROC AUC for each biomarker. Horizontal bars indicate a significant difference in performance between the two models. Error bars represent 95% confidence interval (CI) calculated using a bootstrapping procedure. . . . .	12
2.2	Comparison of model with clinicians. Our model identified three biomarkers that were annotated by clinicians. We present ROC (left column) and precision-recall (right column) curves for SLIVER-net and the baseline 3d CNN model along with individual annotator performance. For subretinal drusenoid deposits, SLIVER-net appears to outperform retina fellows in terms of both AUC and precision-recall, while the reverse is true for hyporeflective drusen cores. . . . .	14
2.3	Confusion matrices for SLIVER-net and the three retinal specialist annotators. 100 of the 390 test set patients were selected for comparison with clinician performance. The remaining 290 patients were used to compute the SLIVER-net threshold, which was selected to match the mean sensitivity of the annotators. For Subretinal Drusenoid Deposits and Intraretinal HRF, SLIVER-net displays a similar sensitivity to clinicians while operating at fewer false positives. . . . .	15

2.4 Examples of discordant cases. B-scans of example cases where SLIVER-net’s determination disagreed with the expert human graders, with heat map overlay highlighting the most informative regions of the image as determined by the algorithm. In A,B virtually no separation can be seen between retinal pigment epithelial (RPE) band and the drusen, which presumably made it difficult for the algorithm to determine that these were intraretinal hyper-reflective foci (IHRF). In fact, on post-hoc review, the senior retina specialist sided with the algorithm. In C, the heat map highlights the relevant features, but the algorithm failed to identify these tiny conical or spike-like elevations as subretinal drusenoid deposits (SDD). It should be noted that no clear distinction in reflectivity is observed between the SDD and the underlying RPE. In D, the heat map highlights a drusen but there are no apparent IHRF. However, there are occasional tiny bright dots in the Henle’s layer which are due to retinal capillaries but may have been confused as IHRF. This is a true false positive. In E, the algorithm detected a drusen with hyporeflective core, but the drusen was small  $< 40\mu\text{m}$  in height. By definition, graders do not assess the internal reflectivity in lesions this small. In F, the algorithm also determined hDC to be present, but the internal reflectivity of the drusen, while reduced, is not dark enough to be called hyporeflective. This is also a true false positive. . . . . 17

2.5	<p>Composite of B-scan Images of Example Cases with Disagreement between Multiple Graders. Top row: IHRF, Middle row: SDD, Bottom row, hDC. In A. A an IHRF is clearly visible (white circle) but is in a region of atrophy. Some graders excluded consideration of the feature as a result. This finding was correctly detected by the algorithm. B. A tiny brighter dot (arrow) is observed in the ELM band. This was interpreted by some graders as a possible IHRF. However, the feature is too small and the reflectivity is not as bright as the RPE band. This finding was correctly excluded by the algorithm. C and D. The EZ has a slightly “wavy” profile suggestive of possible underlying subretinal drusenoid deposits (within the white circles). In both these cases, the algorithm correctly identified the presence of these subtle SDD. E. The drusen (white arrow) is relatively small and its height is borderline for being <math>\geq 40\mu\text{m}</math>, which is the minimum threshold set by the grading protocol in order to be able to assess internal reflectivity. Graders disagreed with regards to whether the lesion met the size criterion. F. The internal reflectivity of the drusen is slightly reduced but is clearly brighter than the vitreous overlying the retina. The reflectivity is not sufficiently reduced to be confident that a hDC is present, and hence the disagreement between graders. IHRF: intraretinal hyper-reflective foci; ELM: external limiting membrane; EZ: ellipsoid zone; SDD: subretinal drusenoid deposits; hDC: hypo-reflective drusen core. . . . .</p>	18
2.6	<p>3D CNN was trained on full data. Top: Mean precision-recall AUC across all biomarkers. Bottom: Mean ROC AUC across all biomarkers. Error bars represent 95% confidence interval (CI) calculated using a bootstrapping procedure. . . . .</p>	19

2.7	Comparison of external datasets and their effect on performance. SLIVER-net trained from scratch (light blue), pre-trained using ImageNet (blue), and pre-trained using Kermany (dark blue). Top: Precision-recall AUC for each biomarker. Bottom: ROC AUC for each biomarker. Horizontal bars indicate statistically significant differences ( $p < 0.05$ ). Error bars represent 95% confidence interval (CI) calculated using a bootstrapping procedure. . . . .	21
2.8	Left: Mean precision-recall AUC across all biomarkers. Right: Mean ROC AUC across all biomarkers. Error bars represent 95% confidence interval (CI) calculated using a bootstrapping procedure. . . . .	23
2.9	Left: Mean precision-recall AUC across all biomarkers. Right: Mean ROC AUC across all biomarkers. Error bars represent 95% confidence interval (CI) calculated using a bootstrapping procedure. . . . .	24
2.10	Our model operated on a 2d tiling of the OCT volume. Resnet18 served as the abstract feature extractor, and the representations for each slice were aggregated using slice integration and a 1D CNN. Finally, biomarkers were predicted using fully connected layers. . . . .	30
3.1	Sensitivity Analysis of expected outcomes under different modelling assumptions for the observed policy. . . . .	42
3.2	High level overview of the IT architecture . . . . .	46
4.1	The probability of adopting a VPF is conditioned on a number of friends who have adopted and b the number of friend’s adoption posts seen. A pattern of complex diffusion is evident, in which as the number of social proof exposures increases, so does the likelihood of the user adopting the frame. . . . .	58

4.2	The probability of adopting a VPF, conditioned on the number of friend’s adoption posts seen, and segmented by a Authoritative health (AH) pages followed by the users, and b the COVID-19 vaccination rate in the user’s home county (binned by top and bottom 25 percent quartiles). These cuts provide proxies for pre-existing vaccine attitudes and show that significantly less social proof is required to reach comparable adoption rates as we move up in the levels (representing more openness to vaccination, in aggregate). . . . .	59
4.3	The probability of adopting a VPF segmented by the levels of tie strength with prior adopters and conditioned on the (a) number of friends who have adopted and (b) the number of friend’s adoption posts seen. Users with strong ties to prior adopters seem to be more likely to adopt the VPF when social proof exposure increases compared with users that have weaker ties with prior adopters. . . . .	60
4.4	<b>(a)</b> A difference in differences (DID) approach was taken to estimate the effect of an influencer’s adoption decision on the decisions of their followers. In this illustrative example, we show an adopting influencer and a matched non-adopting control (in practice, we use 10 matched controls per influencer). We estimate the effect of the adopting influencer’s decision on her followers by looking at the departure from the counterfactual provided by the non-adopting influencer’s followers’ behavior. <b>(b)</b> A permutation method enables deriving an empirical null distribution of DID values per influencer, allowing determination of P-values (adjusted for multiple hypotheses testing), and revealing that only about 5% of influencers show a significant effect at alpha=0.05. . . . .	62



4.5	Logistic regression coefficients for adoption conditioned on the discovery channel (different QPs=“quick promotions”), tie strength of any included social proof, and a broad set of confounders, including those where heterogeneity in adoption response was observed. The coefficients (log odds scores) imply that VPF adoption is strongly affected by social aspects such as seeing a vaccination post from a close friend or seeing a promotion informing users that their friends have adopted the frame. FB-age is defined as the number of days since the user has signed up to the Facebook platform. “l28-” is an inactivity variable defined as the number of days within the last four weeks at which the user has not been active on the FB platform. (*) marks a log transformation. (†) marks a standardization transformation for zero mean and unit variance. . . . .	64
4.6	Users in the treatment arm received the friend aggregation post as a means of social proof for VPFs. The three friends for this format were selected at random, enabling estimation of conditional average treatment effects conditioned on approximated tie strengths to the friends in the aggregation. The findings show an increasing trend in CATE correlated with increasing levels of tie strength.	66
4.7	Each bar represents the importance of the associated feature in maximizing the heterogeneous treatment effect. . . . .	68
4.8	Heterogeneous treatment effects grow with user age (a). This pattern is likely driven by strong ties, as older cohorts tend to have proportionally more stronger tie friends who have adopted (b). As the age cohort’s strong tie friend proportion exceeds 0.25, we see increasing CATE, significantly different from 0 (c) text annotation shows select cohort age ranges). . . . .	69
4.9	(a) Shows the average negative actions per day across all adopters and controls in the sample. (b-f) Show the average negative actions per day, across adopters and controls in each stratum matched based on the propensity score to receive treatment (VPF adoption). . . . .	70

5.1	Comparison of the observed trajectory of the number of cases in United kingdom, New York, Spain, France, Germany, and Denmark (prior to the date where social distancing was imposed). We provide fits based on region-specific parameters (we choose sets of parameters that all lie within the 95% confidence set). The different sets of parameters diverge significantly in the subsequent dates showing the under-determination of this model. . . . .	89
5.2	The time until social distancing ends (in months) based on the SEIR model, using different $R_0$ and $\tau$ values. For each of the regions (Spain, United Kingdom, New York, France, Germany, and Denmark) we also marked the parameters that provided a good fit as shown in Figure 5.1. . . . .	90
5.3	Simulating the number of cases under the social distancing regime where social distancing is turned on when the number of cases exceeds 35 per 10,000 and is turned off when it drops below 5 per 10,000. We show 3 different sets of parameters matching data taken from France as seen in Table 2. . . . .	91
5.4	The range of trajectories for the number of cases predicted using samples from distribution implied by the global parameters estimated on all regions. . . . .	93
5.5	The distribution for the time until social distancing will end implied by the global parameters $R_0^{(0)}, \tau^{(0)}, \sigma_R^2, \sigma_\tau^2$ . The median is October 2020, the mode is September 2020 and the variance is 16 months. . . . .	94
5.6	(a) Different times until social distancing will end based on different choices of the regions. Each sample was generated by choosing four regions out of the six (UK, France, Spain, Germany, New York, and Denmark), estimating their global parameters, and then measuring the median for the time social distancing will end implied by these parameters. (b) End of social distancing regime as a function of the percentage of $R_0$ during the regime. . . . .	96

5.7	SEIR model schema. Each individual in the population begins at susceptible state $S$ , and will enter into the exposure state $E$ with transition rate $\beta I$ in each time unit, where $\beta = R_0\gamma$ . In exposure state $E$ , an individual will go into infectious state $I$ with transition rate $\nu$ . Of all the people who arrive at state $I$ , $p_M$ of them will recover (state $R$ ), $p_H$ will be hospitalized but will never reach critical care (state $H_H$ ), and $p_C$ will be hospitalized to later be in critical care (state $H_C$ ). All transitions from the $I$ state will occur with transition rate $\gamma$ . People in $H_H$ will enter into $R$ with a transition rate $\delta_H$ ; people in $H_C$ state will enter into critical state $C_C$ with a transition rate $\delta_C$ , and then enter into $R$ state with a transition rate $\epsilon_C$ . We set parameters $p_M = 0.956$ , $p_H = 0.0308$ , $p_C = 0.0132$ , $\nu = 1/4.6$ , $\delta_C = 1/6$ , $\delta_H = 1/8$ , $\epsilon_C = 1/10$ as were estimated by Kissler <i>et al.</i> All states are normalized with respect to population size $N$ . . . . .	97
5.8	Parameter estimation diagram: We assume that the parameters $R_0^{(k)}, \sigma_k^2, \tau^{(k)}$ are drawn from a distribution which is defined by the parameters $R_0^{(0)}, \sigma_R^2, \tau^{(0)}$ , and $\sigma_\tau^2$ . We then assume that the cumulative case number curve $y_k(t)$ is generated by the process defined by these parameters. We estimate the most likely values of $R_0^{(0)}, \sigma_R^2, \tau^{(0)}$ , and $\sigma_\tau^2$ using maximum marginal likelihood approach. . . . .	98
A.1	Flowchart illustrating user and data selection for RQ1. . . . .	106
A.2	Flowchart illustrating user and data selection for RQ2. . . . .	107
A.3	Different messages promoting adoption of a Vaccine Profile Frame (VPF) <b>A)</b> VPF Post - A post that is automatically generated upon adoption and displayed to friends within their newsfeed. <b>B)</b> Friend Aggregation Post - A newsfeed post informing users that three of their friends have adopted a VPF. <b>C)</b> Profile/Newsfeed Notification - A non-social notification presented on either the user's profile page or in their newsfeed encouraging them to adopt the frame. . . . .	108

## LIST OF TABLES

2.1	Discordant cases were reviewed by a senior retina specialist grader (SS). Upon re-review the senior retina specialist disagreed with the original ground truth grading in some cases, but in all discordant cases the findings were borderline. Observations with regards to the cause for difficulty in ground-truth assessment are provided. FP: false positive; FN: false negative; IHRF: intraretinal hyper-reflective foci; SDD subretinal drusenoid deposits; hDC hyporeflective drusen core; RPE: retinal pigment epithelium. . . . .	16
2.2	The biomarkers used for this study and their prevalence throughout the three datasets. . . . .	28
3.1	Encounter Statistics . . . . .	39
3.2	Demographics Statistics . . . . .	39
3.3	Prevalence of Disease Severity Levels . . . . .	40
3.4	Summary of estimated expected outcomes for different models across different rewards with 95% confidence intervals (best values highlighted). . . . .	41
3.5	P-values indicating the significance of sequential data in modelling next best actions.	44
4.1	Regression Models Performance . . . . .	63
4.2	Descriptive statistics of features used in the regression analysis. . . . .	72
4.3	Descriptive statistics of features used in the randomized field experiment analysis.	75
4.4	Descriptive statistics of features used in the backfire effect analysis. . . . .	78
4.4	Descriptive statistics of features used in the backfire effect analysis. . . . .	79
5.1	The maximum-likelihood estimates for $R_0$ and $\tau$ for every region . . . . .	92

5.2 The date of the end of social distancing for different sets of parameters that fit  
the data (as shown in Figures 5.1 and 5.2) . . . . . 104

## ACKNOWLEDGMENTS

First and foremost, I wish to express my profound gratitude to my PhD advisors, Eran Halperin and Sriram Sankararaman. The impact of their mentorship on my PhD journey has been immeasurable. Both Eran and Sriram have been outstanding mentors, generously sharing their extensive knowledge and expertise. I am deeply grateful for their guidance, support, and the freedom they granted me to explore my research interests, as well as the opportunities to collaborate with others.

I owe a debt of gratitude to Eleazar, who introduced me to the world of computational biology. His mentorship throughout my PhD has been invaluable. Eleazar is a remarkable example in fostering a collaborative, interdisciplinary research community. One key lesson from him: great relationships often start with great food.

I would like to extend a special thanks to my dissertation committee members: Sriram Sankararaman, Eleazar Eskin, Bogdan Pasaniuc, Wei Wang, and Eran Halperin, for their guidance, patience, and unwavering support throughout my PhD journey.

My collaboration with SriniVas Sadda have been pivotal. His willingness to share clinical insights has been a great asset, for which I am thankful.

The support and camaraderie of my friends and colleagues at UCLA have been invaluable. My lab mates, Jeff Chiang, Brian Hill, Oren Avram, Akos Rudas, Brandon Jew, Ella Petter, Liat Shenhav, Alec Chiu, Chris Robles, Ruthie Johnson, and Boyang Fu, have greatly enriched my lab experience, making it both enjoyable and educational.

I am grateful for my collaborations with Amit Bahl, Udi Weinsberg, and Sindhu Ernala at Meta, and with Dominik Dahlem at Optum Labs, all of whom have been fantastic collaborators.

My appreciation extends to Sim-Lin Lau, Stephania Kay, Joseph Brown, and Helen Tran for their assistance and patience with administrative tasks, often going beyond their duties

to help.

Finally, the unwavering support of my family has been the cornerstone of my journey. My parents, Moti and Zehava, have been steadfast in their support of my education, providing everything necessary for my success. Their unconditional love, encouragement, and patience have been invaluable. They embody the ideals of kindness, hard work, perseverance, and striving for excellence, and will always be my role models.

Chapter 2 is a version of Nadav Rakocz, Jeffrey N. Chiang, Muneeswar G. Nittala, Giulia Corradetti, Liran Tiosano, Swetha Velaga, Michael Thompson, Brian L. Hill, Sriram Sankararaman, Jonathan L. Haines, Margaret A. Pericak-Vance, Dwight Stambolian, Srinivas R. Sadda, Eran Halperin “Automated identification of clinical features from sparsely annotated 3-dimensional medical imaging.” *NPJ digital medicine* (2021). Chapter 3 is a version of a work that was submitted for publication by Nadav Rakocz, Dominik Dahlem, Vijay Nori, Zahra Mahmoodzadeh, Jeffrey Hertzberg, David Cook, Eran Halperin “Augmenting Clinical Practice Guidelines Using Reinforcement Learning and Causal Transformers”. Chapter 4 is a version of Rakocz Nadav, Ernala Sindhu, Nir Israel, Weinsberg Udi, Bahl Amit “The heterogeneous effects of social support on the adoption of Facebook’s vaccine profile frames feature.” *Humanities and Social Sciences Communications* (2023). Chapter 5 is a version of a work that was submitted for publication by Nadav Rakocz, Boyang Fu, Eran Halperin, and Sriram Sankararaman titled “A Statistical Model for Quantifying the Needed Duration of Social Distancing for the COVID-19 Pandemic”.

## VITA

- 2020-2023 Ph.D. Candidate Computer Science  
University of California, Los Angeles, CA
- 2022 Research Intern  
Optum Labs, Minnetonka, MN
- 2020-2022 M.S. Computer Science  
University of California, Los Angeles, CA
- 2021 Research Intern  
Meta, Manlo Park, CA
- 2018-2020 Statistician  
University of California, Los Angeles, CA
- 2017-2018 Machine Learning Engineer  
Ultime Genomics, Israel
- 2012-2016 B.S. Electrical and Electronic Engineering, summa cum laude  
Tel Aviv University
- 2013-2016 Logic Design Engineer  
Intel, Israel

## PUBLICATIONS

\* denotes equal contribution



**Prediction of activity in eyes with macular neovascularization due to age-related macular degeneration using deep learning** Rakocz Nadav\*, Corradetti Giulia\*, Chiang Jeffrey N, Avram Oren, Udi, Nittala Muneeswar Gupta, Halperin Eran, Sadda Srinivas; Eye (2023)

**The heterogeneous effects of social support on the adoption of Facebook's vaccine profile frames feature** Rakocz Nadav, Ernala Sindhu, Nir Israel, Weinsberg Udi, Bahl Amit; Humanities and Social Sciences Communications (2023)

**Automated identification of incomplete and complete retinal epithelial pigment and outer retinal atrophy using machine learning** Chiang Jeffrey N, Corradetti Giulia, Nittala Muneeswar Gupta, Corvi Federico, Rakocz Nadav, Rudas Akos, Durmus Berkin, An Ulzee, Sankararaman Sriram, Chiu Alec, Halperin Eran\*, Sadda Srinivas R\*; Ophthalmology Retina (2023)

**Automated large-scale prediction of exudative AMD progression using machine-read OCT biomarkers** Rudas Akos, Chiang Jeffrey N, Corradetti Giulia, Rakocz Nadav, Avram Oren, Halperin Eran\*, Sadda Srinivas R\*; PLOS Digital Health (2023)

**Automated large-scale AMD progression prediction using machine-read OCT biomarkers** Rudas Akos, Chiang Jeffrey N, Corradetti Giulia, Rakocz Nadav, Halperin Eran\*, Sadda Srinivas R\*; medRxiv (2022)

**Automated Identification of Incomplete and Complete Retinal Epithelial Pigment and Outer Retinal Atrophy Using Machine Learning** Corradetti Giulia, Chiang Jeffrey N, Corvi Federico, Rakocz Nadav, Rudas Akos, Durmus Berkin, Chiu Alec, An Ulzee, Sankararaman Sriram, Halperin Eran\*, Sadda Srinivas R\*; Investigative Ophthalmology &

Visual Science (2022)

**Methylation risk scores are associated with a collection of phenotypes within electronic health record systems** Thompson Mike\*, Hill Brian L\*, Rakocz Nadav, Chiang Jeffrey N, Geschwind Daniel, Sankararaman Sriram, Hofer Ira, Cannesson Maxime, Zaitlen Noah, Halperin Eran; NPJ genomic medicine (2022)

**Imputation of the continuous arterial line blood pressure waveform from non-invasive measurements using deep learning** Hill Brian L, Rakocz Nadav, Rudas Ákos, Chiang Jeffrey N, Wang Sidong, Hofer Ira, Cannesson Maxime\*, Halperin Eran\*; Scientific reports (2021)

**Automated identification of clinical features from sparsely annotated 3-dimensional medical imaging** Rakocz Nadav\*, Chiang Jeffrey N\*, Nittala Muneeswar G, Corradetti Giulia, Tiosano Liran, Velaga Swetha, Thompson Michael, Hill Brian L, Sankararaman Sriram, Haines Jonathan L, Halperin Eran\*, Satta Srinivas R\*; NPJ digital medicine (2021)

**Non-Invasive and Continuous Blood Pressure Monitoring Using Deep Convolutional Neural Networks** Hill Brian L, Rakocz Nadav, Chiang Jeffrey N, Hofer Ira, Halperin Eran, Cannesson Maxime; ANESTHESIA AND ANALGESIA (2020)

**A Statistical Model for Quantifying the Needed Duration of Social Distancing for the COVID-19 Pandemic** Rakocz Nadav, Fu Boyang, Halperin Eran, Sankararaman Sriram; medRxiv (2020)

**An automated machine learning-based model predicts postoperative mortality using readily-extractable preoperative electronic health record data** Hill Brian

L, Brown Robert, Gabel Eilon, Rakocz Nadav, Lee Christine, Cannesson Maxime, Baldi Pierre, Loohuis Loes Olde, Johnson Ruth, Jew Brandon, Uri Maoz, Aman Mahajan, Sriram Sankararaman, Ira Hofer, Halperin Eran; British journal of anaesthesia (2019)

# CHAPTER 1

## Introduction

### 1.1 Scope of Research

The advent of the information age has been marked by significant advancements in computational power and data storage capabilities, giving rise to the era of "big data." This proliferation of data, when harnessed through machine learning algorithms, has the potential to uncover complex patterns and make accurate predictions across various domains, including healthcare. Particularly, the integration of hardware advancements, such as graphics processing units (GPUs) with deep learning algorithms, has led to breakthroughs in fields like computer vision and natural language processing, thereby transforming the landscape of medical decision-making [LBH15, KSH12].

The healthcare industry, in its transition from paper-based to electronic medical records, has amassed a wealth of digitized health information. This pivotal shift, enriched by the integration of detailed medical imaging scans, affords a thorough representation of a patient's health status. The incorporation of longitudinal data within these EMRs, when analyzed through advanced deep learning modalities such as transformers [VSP], unlocks the potential to discern complex health patterns over time. This capability is instrumental in enhancing the precision and efficacy of clinical care, as it facilitates the early detection of health issues and the tailoring of treatment plans to individual patient needs.

In addition to EMRs, medical imaging has also been able to leverage recent advancements in computer vision technology. These developments have transformed the process of interpreting

medical images by automating what were once labor-intensive and time-consuming manual tasks. This shift in approach significantly lightens the workload related to technical image analysis, freeing clinicians to allocate more time and attention to direct patient care. Such improvements in efficiency not only raise the quality of healthcare services but also increase the capacity of healthcare providers to care for a greater number of patients. This optimized utilization of resources and time ultimately leads to enhanced patient outcomes and overall healthcare system effectiveness.

To complement this digital evolution, comes the availability of social media information, which offers both insights into patient behaviors, and environmental factors as well as valuable information regarding populations that enables public health experts to identify trends and understand health-related behaviors on a community scale with greater precision.

Together, these diverse data types, encompassing clinical, imaging, and social media-derived information, hold immense potential for impacting healthcare delivery. By enabling a more holistic understanding of health determinants, they support the development of personalized treatment plans and informed public health interventions, therefor enhancing healthcare decisions at both individual and population levels.

However, despite these advancements, the application of computational methods to healthcare decision-making faces significant challenges. A primary concern is the limited availability of annotated training data, particularly in clinical imaging, where data is scarce due to regulatory restrictions and the high cost of data collection and manual annotation [GGS16, Ker18, Taj20]. Additionally, many clinical guidelines are developed with a population-based approach, which may not fully incorporate the comprehensive data available for individual patients. This methodology often results in guidelines that, while efficient, may oversimplify complex medical scenarios thereby diminishing the potential for personalized and effective healthcare management. [CV]. Furthermore, Although social media data is extensive and holds the potential to inform and influence medical decisions at a population level [BBC15], the methodologies for effectively and safely harnessing this resource for medical decision-making

remain unclear. These gaps underscore the need for novel computational approaches that can address these challenges and support more informed and effective healthcare decisions.

## 1.2 Contributions and Overview

This dissertation proposes innovative computational methods to address the aforementioned challenges, with a focus on enhancing medical decision-making through the application of deep learning, reinforcement learning, and social network analysis. Each chapter of the dissertation contributes to bridging the knowledge gaps and advancing the field of computational healthcare.

Chapter 2 introduces SLice Integration of Volumetric features Extracted by pre-trained Residual neural networks (SLIVER-net), a protocol leveraging transfer learning to predict disease-related biomarkers from limited annotated 3-dimensional imaging data. By utilizing external datasets of 2-dimensional images [Ker18, DDS09], SLIVER-net demonstrates superior performance in identifying risk factors for retinal disease from optical coherence tomography (OCT) images. This addresses the critical issue of limited data availability and demonstrates the potential of computational methods to automate identification processes in clinical images.

Chapter 3 tackles the limitations of Clinical practice guidelines by introducing the Importance Transformer (IT), an offline reinforcement learning approach that optimizes expected outcomes through a novel loss function incorporating weighted importance sampling. This method shows promise in managing chronic diseases while outperforming recent utilizations of the transformer architecture into reinforcement learning [CLR, JLL], thus aligning with the dissertation’s aim of developing robust and accurate methods for policy formulation.

Chapter 4 explores the impact of social influence on vaccine decision-making by analyzing the use of vaccination profile frames (VPFs) on Facebook. This study provides insights into the causal factors promoting VPF adoption and assesses the potential negative social interactions resulting from expressing support for a polarizing issue like COVID-19 vaccination

[Met21, Oz18, Sch18]. It highlights the potential of social media platforms as tools for influencing public health decisions and underscores the importance of computational methods in understanding and leveraging social influence in healthcare decision-making.

Chapter 5 addresses the challenge of modeling the transmission dynamics of COVID-19 to inform social distancing strategies [KTG20]. It introduces a hierarchical Bayesian model based on the SEIR framework, fitted to current COVID-19 case data, to provide precise parameter estimates and formal confidence intervals for critical parameters like the reproduction number [LGW20, NMS20]. This approach demonstrates the practical application of computational methods in guiding healthcare strategies at a population level.

Overall, this dissertation presents a cohesive narrative that integrates computational methods with healthcare decision-making. It addresses key challenges in the field and provides novel solutions that have the potential to significantly improve patient care and health outcomes. Through a combination of deep learning, reinforcement learning, Bayesian statistics, and social network analysis, the dissertation contributes to the advancement of computational healthcare, ultimately supporting safer, more efficient, and data-driven medical decisions.

## CHAPTER 2

# Automated identification of clinical features from sparsely annotated 3-dimensional medical imaging

### 2.1 Introduction

The application of deep learning, specifically Convolutional Neural Networks (CNNs), has proven to be successful for detecting and predicting disease from medical image data[GG16, Rot16, Qi 16, Ker18, Taj20]. However, the application of deep learning to novel tasks has been hampered by the availability of appropriately annotated training data. Biomedical research questions, in particular, present an inherent challenge in terms of sample size. While large datasets have been released in collaboration with medical imaging (e.g., CheXpert[Irv19] (224,316 X-rays), ISIC[CGC18, TRK18] (25,331 dermoscopic images), ABCD-NP[Pfe18] (8500 MRI volumes), and others, e.g. <http://www.grand-challenge.org/>), current regulations (e.g., HIPAA in the United States) restrict the ability to collect sufficient data to apply deep learning to novel questions. Generally, clinical and biomedical research reports are based on small cohorts numbering in hundreds. For example, Lutkenhoff et al., 2015[Lut15] established the largest annotated cohort of patients (143) with disorders of consciousness, and the ImageCLEF initiative curated 403 CT scans for the study of tuberculosis[CDD15]. Additionally, Lei et al.[LBA17] analyzed 138 patients to determine the risk for age-related macular degeneration. In addition to the extensive clinical time required to collect cohorts, there is the added burden of manually annotating patient information to enable machine learning[GG16, Rav17, NZA16, KSB17]. All these factors present a high cost for applying



deep learning methods to new data modalities and address novel questions.

Transfer learning[GG16, Rav17, PY10] can be used to address the small number of annotated (or labeled) samples by introducing information from another domain. However, when the data consists of 3-dimensional volumes, transfer learning cannot be directly applied unless other 3-dimensional volumes are available in sufficient quantity for reference in external datasets. Unlike resources for 2-dimensional images such as ImageNet[DDS09], no such resource is available for 3-dimensional data (e.g., CT, MRI, OCT, etc.). To circumvent this problem we developed a protocol for applying deep learning to a dataset with limited annotated 3-dimensional imaging data. Our approach leverages external datasets of 2-dimensional images and uses transfer learning to predict AMD-related biomarkers in 3-dimensional volumes. We transformed 3-dimensional to 2-dimensional data to make it compatible with the external set. Converting 3-dimensional to 2-dimensional data results in loss of information, therefore, we introduced an operation (slice integration) to counter the information loss. We name this approach SLice Integration of Volumetric features Extracted by pre-trained Residual neural networks (SLIVER-net).

To illustrate the effectiveness of SLIVER-net, we tested the ability of SLIVER-net to identify risk factors for retinal disease from optical coherence tomography (OCT) images. Because of its high axial resolution and histological detail, OCT is able to assess the integrity of the retinal layers[She09, COP18, BW15] in a variety of conditions including optic nerve disorders[GT13], retinal diseases[Kea12], and systemic conditions which may have ocular manifestations[DWB11, Kah18]. OCT has been particularly transformative in the management of age-related macular degeneration (AMD), the leading cause of blindness in developed nations. Initially, AMD may manifest drusen, which are accumulations of material under the retinal pigment epithelium (RPE). Vision may be relatively good at this early or intermediate stage. Eventually, a significant number of patients develop macular neovascularization (MNV) and/or geographic atrophy (GA), which are considered late manifestations and associated with considerable loss of vision. Effective treatments (anti-

vascular endothelial growth factor, or anti-VEGF) have been developed for MNV, but thus far, there is no treatment for GA. In addition, despite the availability of treatments for MNV, many “successfully” treated patients eventually go on to develop atrophy and vision loss. The best outcomes for the treatment of active MNV are observed in patients who are treated early while the neovascular lesions are small. Therefore, identifying patients who are at high-risk for progression to late AMD is essential to identify appropriate intervals for monitoring patients with earlier stages of AMD. A number of OCT risk factors for progression to late AMD have been defined and include intraretinal hyperreflective foci (which are thought to represent migration of RPE into the retina), hyporeflective cores within drusen (shown to correspond to calcific nodules[TPF18]), subretinal drusenoid deposits, and high central drusen volume. Recently, Lei et al (2017)[LBA17] proposed a system using OCT images for integrating these factors into a simple score that could reflect a given patient’s risk for conversion to late AMD. This system was later validated by Nassisi et al (2019)[Nas19] in a post-hoc analysis of intermediate AMD fellow eyes from subjects enrolled in the HARBOR study. Despite this compelling data regarding these OCT biomarkers which could be used to risk stratify patients and define appropriate intervals for monitoring, most clinicians do not have time to assess these OCT features in the context of a busy clinical practice. Ideally, these risk factors for progression should be detected automatically from the OCT, which would allow a risk score to be immediately available to the clinician. Such a risk score could also potentially be used to identify high-risk patients for enrollment into early intervention trials or to monitor disease progression over time in a more precise or quantitative fashion. Moreover, beyond its immediate clinical impact, an automated system to assess risk on OCT could be used for research investigations to probe large datasets such as the UK Biobank or the electronic health records and image databases of large health systems. This would allow the variability of the evolution of these biomarkers to be more precisely characterized. An automated risk score could also be used as quantitative endophenotype in genetic discovery studies, particularly those aimed at identifying genetic risk factors for disease progression.

We applied SLIVER-net to automatically identify these factors, henceforth termed “biomarkers”. Recent applications to OCT images have focused on predicting glaucoma[An19, Asa19], different severities of AMD[RLO19], and other diseases[Ker18, Kuw19, Fau18]. Because the clinical and biological bases for these biomarkers are still under investigation, there are relatively few examples with which we can develop a deep learning approach. SLIVER-net specifically targets such scenarios, in which the number of annotated 3-dimensional images is small (in the hundreds). Still, SLIVER-net was able to outperform current methods and sometimes better than the retina specialists. Our results demonstrate that our method is superior to expert retinal image graders. Notably, the improvements provided by SLIVER-net are primarily driven by transfer learning and slice integration, both of which are not limited to biomarker prediction nor OCT classification, and thus applicable to other 3-dimensional imaging modalities. Our analysis was done on a few hundred annotated images and demonstrates the utility of SLIVER-net for analyzing a small dataset and generalizing the annotation for a larger database.

## 2.2 Results

### 2.2.1 The SLIVER-net model

Our model, SLIVER-net, is a novel deep neural network architecture designed to operate on 3-dimensional images despite a limited number of manually annotated examples. In order to cope with the small sample size of labeled data SLIVER-net leverages external information through transfer learning from 2-dimensional images, then fine-tuned using a small set of labeled 3-dimensional images (with medically relevant annotations). The labels of the 2-dimensional images are not required to have any medical relevance, as previous investigations have shown that models learn to represent domain-general features in the transfer learning paradigm[PY10]. Typically, the 3-dimensional volumes with desired labels can number in the hundreds, while the external dataset will consist of tens of thousands,

or ideally millions of images. After training SLIVER-net can be applied to a 3-dimensional image to predict the annotated outcomes without further need of the external dataset.

To enable transfer learning between images and volumes SLIVER-net differs from standard algorithms in two ways. First, it re-frames the 3D OCT volume as a 2D “tiling” (e.g., mosaic) of slices, allowing for the use of transfer learning with currently available 2-dimensional datasets. Second, there are additional layers to the deep neural network which enable SLIVER-net to preserve the 3-dimensional spatial structure lost by tiling. (See Methods: Table 2 for further details).

The SLIVER-net model itself consists of three steps. First, the re-framed OCT volume (tiled images) is passed through a “backbone” convolutional neural network (CNN), for which the output is a representation in an abstract feature space. Then, a slice aggregation operation is applied to compress this representation and obtain information that is shared across adjacent slices. Finally, a decision module operates on this compressed representation to determine the presence or absence of biomarkers. A more detailed description of SLIVER-net is provided in the Methods.

### **2.2.2 AMD-Related biomarker prediction**

In order to demonstrate its utility, we applied SLIVER-net to biomarker prediction from OCT, which has been the primary driver of breakthroughs in the understanding and characterization of novel biomarkers associated with AMD[Nit19]. The identification of these biomarkers in an OCT scan requires careful manual inspection and annotation of each slice (termed a B-scan) within the OCT volume, which is highly laborious and time-consuming. It is therefore desirable to develop automatic tools that will replace manual annotation. Thus, we developed SLIVER-net to automatically predict biomarkers in early and intermediate AMD.

Data were collected across three sites: University of Miami (369 patients), Case Western Reserve University (248 patients), and University of Pennsylvania (390 patients). We employed

an external validation approach, where data from two of the sites, the University of Miami and Case Western Reserve University, were used to develop and validate the model, and data from the University of Pennsylvania were reserved as an external testing set (see Methods for additional details). The separation into three different datasets ensured that there was no overlap between the patients used for model development and testing. In total, the training and testing sets included OCT volumes from 1007 patients, which is currently the largest available dataset annotated for these biomarkers[Nit19].

In order to overcome the challenge of a limited dataset, we incorporated a large publicly available dataset[Ker18], containing 84,495 2-dimensional OCT images (only horizontal B-scans passing through the foveal center) using transfer learning. These 2-dimensional fovea-centered OCT images provide only partial information since they do not contain 3-dimensional volume information and no information about macular regions beyond the foveal depression. This scenario fits the case for which SLIVER-net was designed. We trained SLIVER-net using this external information from 2-dimensional fovea-centered scans, along with the 3-dimensional information from the OCT volumes from the University of Miami and Case Western Reserve University.

SLIVER-net was successfully able to predict the four AMD-related OCT biomarkers evaluated in this study. Three of these biomarkers, intraretinal hyperreflective foci, subretinal drusenoid deposits, and hyporeflective drusen cores, were manually annotated, while the other biomarker (high central drusen volume) was determined based on information provided from another OCT device (Cirrus OCT). In addition, SLIVER-net was able to use the OCT data alone to predict another marker (reticular pseudodrusen) determined by infrared reflectance.

### **2.2.3 Comparison of SLIVER-net to state of the art deep learning approaches**

We compared SLIVER-net with two alternative models: a 3D CNN and a 2D CNN using the same image stacking approach. 3D CNNs, which are commonly used for MRI and CT analysis[JLT18, HSV17, MNA16], represent the current state of the art in volumetric image

analysis. 3D CNNs are able to consider the 3-dimensional structure in a volume instead of operating slice by slice but require very large amounts of training data due to the large number of model parameters. Specifically, 3D CNNs have a substantially larger number of parameters compared to standard 2D CNNs. In addition, we also included a 2D CNN which used the same image tiling approach as SLIVER-net, which serves as a baseline model for assessing the effectiveness of transfer learning and slice pooling. The alternative deep learning models (see Methods) were trained to predict biomarkers associated with AMD using the same training data from the University of Miami and Case Western Reserve University (see Methods for more details about the train and test sets).

Due to the strongly imbalanced nature of biomarker prevalence, the models were evaluated using area under the ROC curve (AUROC) and precision-recall curve (AUPRC) metrics. On the University of Pennsylvania test set (740 volumes), the 3D CNN predicted all biomarkers with a mean ROC area under the curve (AUC) of 0.81[95% Confidence Interval (CI): 0.75,0.86], and a mean precision-recall AUC of 0.22[CI: 0.17,0.33], and the 2D CNN performed with mean ROC AUC of 0.79[CI: 0.67,0.82] and a mean precision-recall AUC of 0.19[CI: 0.16,0.28]. SLIVER-net achieved a mean ROC AUC of 0.94[CI: 0.91,0.96], and a mean precision-recall AUC of 0.41[CI: 0.34,0.51] thus showing significant improvement over the alternative approaches in terms of ROC AUC (p-value < 0.001) and precision-recall AUC (p-value < 0.001). The performance on each individual biomarker is shown in Figure 2.1.

#### **2.2.4 Comparison of SLIVER-net with specialist clinician assessments**

Additionally, we compared SLIVER-net’s predictions against expert retinal image graders (retina specialists who had been certified for OCT image grading by the Doheny Image Reading Center) with respect to the manually annotated biomarkers. Within the test set from University of Pennsylvania, 100 patients were randomly selected and their OCT volumes were read by an additional three retina specialists.

We observed that SLIVER-net outperformed all clinician experts in identifying subretinal

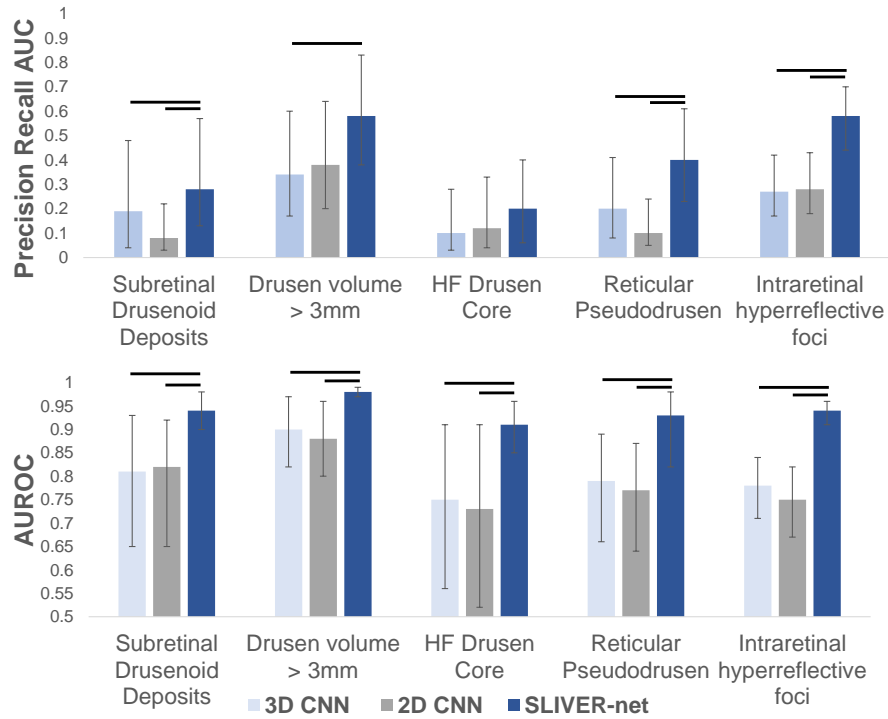


Figure 2.1: SLIVER-net performance. SLIVER-net (dark blue) was compared with a 3D CNN backbone approach (light blue) and 2D CNN (gray). SLIVER-net significantly outperformed both the 3D CNN and 2D CNN in identifying each biomarker in terms of area under the ROC (AUROC) and area under the Precision-Recall curve (Precision-Recall AUC). Top: Precision-Recall AUC for each biomarker. Bottom: ROC AUC for each biomarker. Horizontal bars indicate a significant difference in performance between the two models. Error bars represent 95% confidence interval (CI) calculated using a bootstrapping procedure.

drusenoid deposits, two out of the three clinicians in identifying intraretinal hyperreflective foci in terms of both ROC metrics and precision-recall (Figure 2.2), generally predicting fewer false positives while maintaining the same sensitivity (Figure 2.3). However, SLIVER-net was inferior in identifying hyporefective drusen cores. We also observed that SLIVER-net was successful at predicting high central drusen volume and reticular pseudodrusen which clinicians would not be able to assess without additional equipment.

Cases where SLIVER-net disagreed with specialist annotations were sent to the same clinician panel with an additional senior specialist for review (Table 1). The post-hoc review revealed that most of SLIVER-net’s errors occurred during difficult reads, in which the biomarker was small, subtle or located in close proximity to another structure which made it difficult to distinguish the feature from the background (e.g. a hyper-reflective focus close to the RPE surface). In addition, there was disagreement among the clinician panel in many of the cases where SLIVER-net produced a false positive (Figure 2.4). For subretinal drusenoid deposits, 16 of the 19 false positives (84.2%) did not have a consensus among annotators; for hyperreflective foci, 16 of the 20 false positives (80%) did not have a consensus; and for hyporefective drusen core, 10 out of 59 false positives (16.9%) did not have a consensus from the annotators (examples are visualized in Figure 2.5). After review, some of these false positives were deemed to be errors in the initial annotation, and in these cases SLIVER-net detected these biomarkers while the clinician panel did not (Figure 2.4: A, B). This further highlights the potential of SLIVER-net as an aid to clinicians in assessing for the presence of these biomarkers.

### **2.2.5 Effect of sample size on the model performance**

We found that SLIVER-net outperforms a standard 3D-CNN in the setting of a relatively small sample size. However, the necessary number of annotated samples required to achieve high performance is unclear. To address this question, we re-trained SLIVER-net with a reduced number of OCT volumes available and measured the performance on the test set



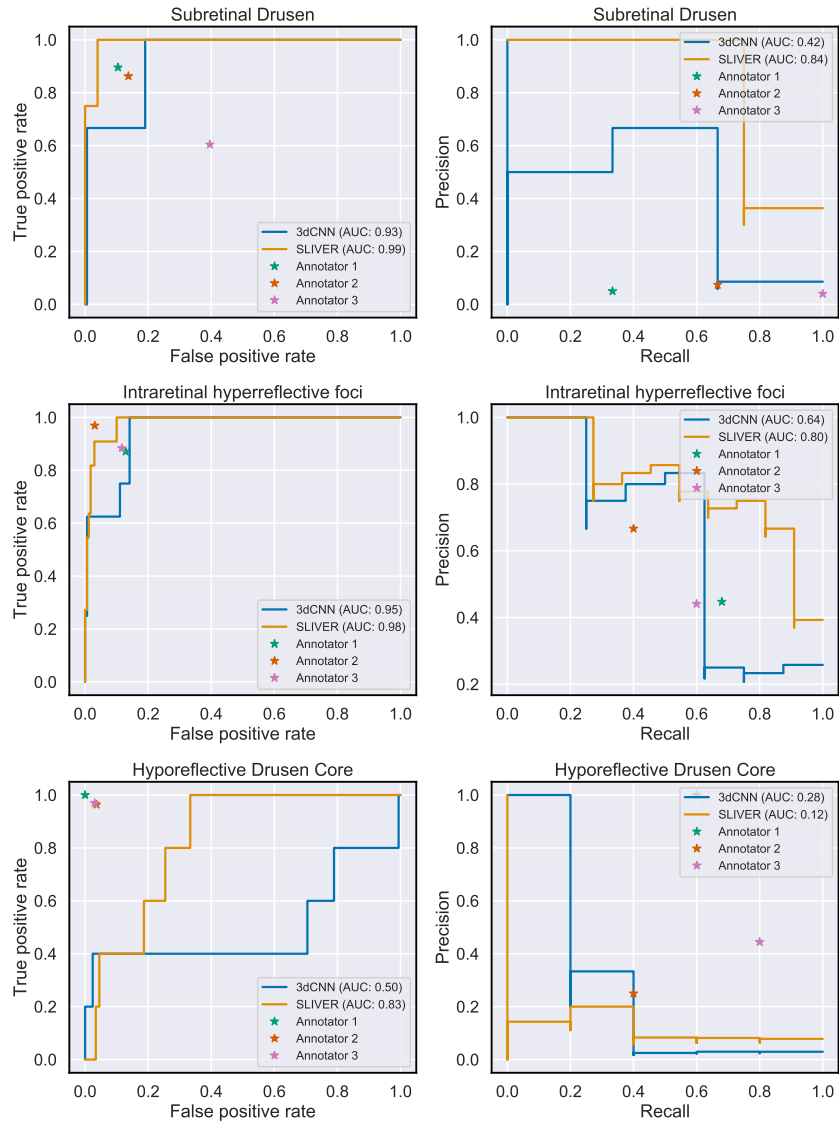


Figure 2.2: Comparison of model with clinicians. Our model identified three biomarkers that were annotated by clinicians. We present ROC (left column) and precision-recall (right column) curves for SLIVER-net and the baseline 3d CNN model along with individual annotator performance. For subretinal drusenoid deposits, SLIVER-net appears to outperform retina fellows in terms of both AUC and precision-recall, while the reverse is true for hyporeflexive drusen cores.

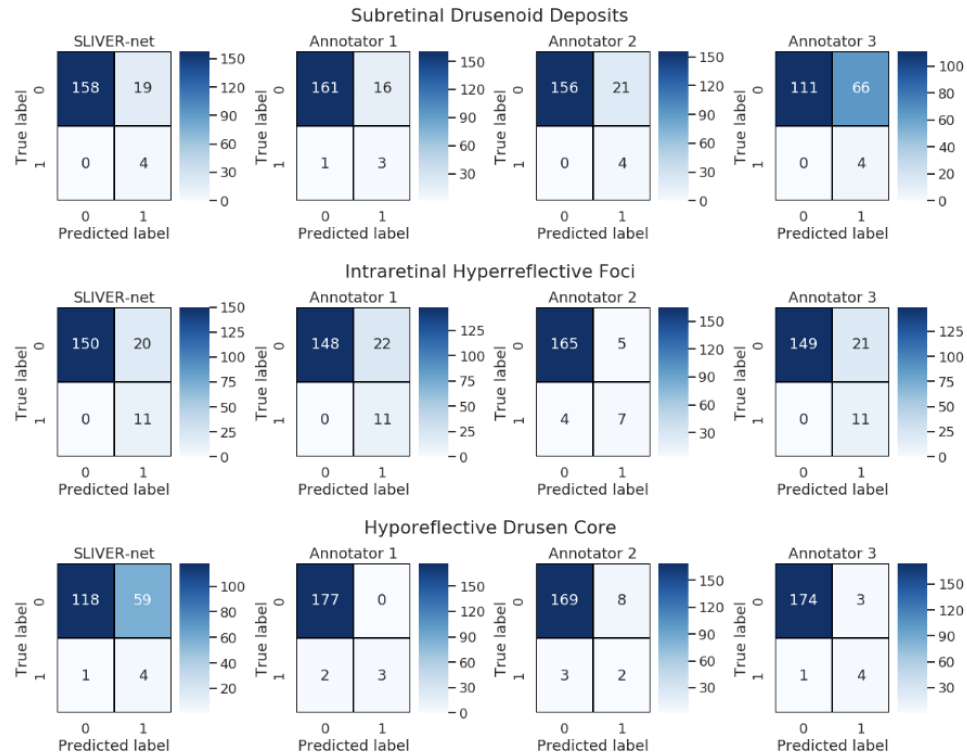


Figure 2.3: Confusion matrices for SLIVER-net and the three retinal specialist annotators. 100 of the 390 test set patients were selected for comparison with clinician performance. The remaining 290 patients were used to compute the SLIVER-net threshold, which was selected to match the mean sensitivity of the annotators. For Subretinal Drusenoid Deposits and Intraretinal HRF, SLIVER-net displays a similar sensitivity to clinicians while operating at fewer false positives.

Table 2.1: Discordant cases were reviewed by a senior retina specialist grader (SS). Upon re-review the senior retina specialist disagreed with the original ground truth grading in some cases, but in all discordant cases the findings were borderline. Observations with regards to the cause for difficulty in ground-truth assessment are provided. FP: false positive; FN: false negative; IHRF: intraretinal hyper-reflective foci; SDD subretinal drusenoid deposits; hDC hyporeflective drusen core; RPE: retinal pigment epithelium.

<b>Post hoc analysis of discordant cases between algorithm and ground truth</b>			
	<b>Discordant after review</b>	<b>Concordant after review</b>	<b>Observations from post hoc review</b>
IHRF FP (N=5)	1	4	Small IHRFs could be observed but were close to the minimum threshold size to be included
IHRF FN (N=4)	2	2	IHRF were in close proximity to the RPE band making separation from the band more difficult to discern
SDD FP (N=10)	2	8	Poor quality of B-scan images makes it more difficult to separate the SDD from the outer retinal bands (EZ, RPE)
SDD FN (N=7)	1	6	SDDs very small in size
hDC FP (N=10)	5	5	Drusen of smaller size making assessment of internal reflectivity difficult. Level of hyporeflectivity was borderline
hDC FN (N=1)	0	1	Feature missed by grader

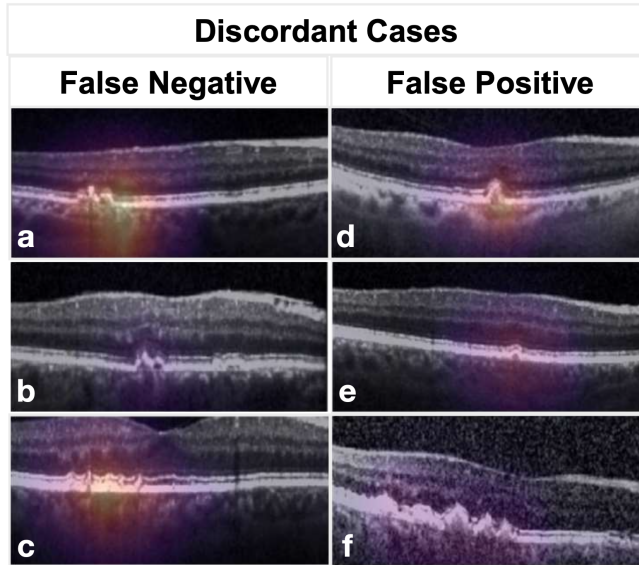


Figure 2.4: Examples of discordant cases. B-scans of example cases where SLIVER-net’s determination disagreed with the expert human graders, with heat map overlay highlighting the most informative regions of the image as determined by the algorithm. In A,B virtually no separation can be seen between retinal pigment epithelial (RPE) band and the drusen, which presumably made it difficult for the algorithm to determine that these were intraretinal hyper-reflective foci (IHRF). In fact, on post-hoc review, the senior retina specialist sided with the algorithm. In C, the heat map highlights the relevant features, but the algorithm failed to identify these tiny conical or spike-like elevations as subretinal drusenoid deposits (SDD). It should be noted that no clear distinction in reflectivity is observed between the SDD and the underlying RPE. In D, the heat map highlights a drusen but there are no apparent IHRF. However, there are occasional tiny bright dots in the Henle’s layer which are due to retinal capillaries but may have been confused as IHRF. This is a true false positive. In E, the algorithm detected a drusen with hyporeflective core, but the drusen was small  $< 40\mu\text{m}$  in height. By definition, graders do not assess the internal reflectivity in lesions this small. In F, the algorithm also determined hDC to be present, but the internal reflectivity of the drusen, while reduced, is not dark enough to be called hyporeflective. This is also a true false positive.

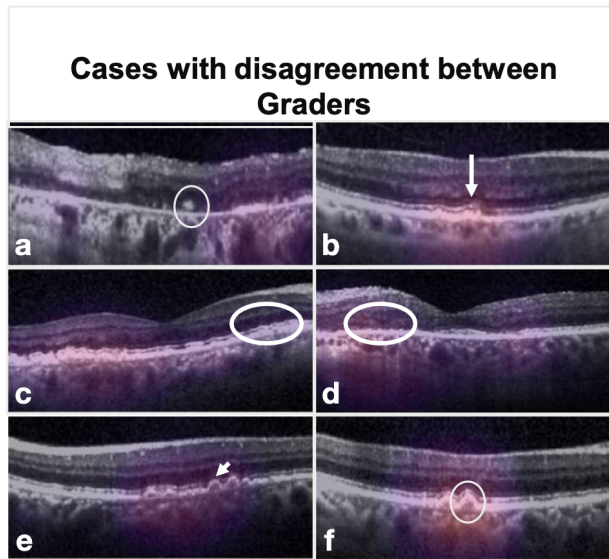


Figure 2.5: Composite of B-scan Images of Example Cases with Disagreement between Multiple Graders. Top row: IHRF, Middle row: SDD, Bottom row, hDC. In A. An IHRF is clearly visible (white circle) but is in a region of atrophy. Some graders excluded consideration of the feature as a result. This finding was correctly detected by the algorithm. B. A tiny brighter dot (arrow) is observed in the ELM band. This was interpreted by some graders as a possible IHRF. However, the feature is too small and the reflectivity is not as bright as the RPE band. This finding was correctly excluded by the algorithm. C and D. The EZ has a slightly “wavy” profile suggestive of possible underlying subretinal drusenoid deposits (within the white circles). In both these cases, the algorithm correctly identified the presence of these subtle SDD. E. The drusen (white arrow) is relatively small and its height is borderline for being  $\geq 40\mu\text{m}$ , which is the minimum threshold set by the grading protocol in order to be able to assess internal reflectivity. Graders disagreed with regards to whether the lesion met the size criterion. F. The internal reflectivity of the drusen is slightly reduced but is clearly brighter than the vitreous overlying the retina. The reflectivity is not sufficiently reduced to be confident that a hDC is present, and hence the disagreement between graders. IHRF: intraretinal hyper-reflective foci; ELM: external limiting membrane; EZ: ellipsoid zone; SDD: subretinal drusenoid deposits; hDC: hypo-reflective drusen core.

(Figure 2.6). We observed that a sample size of 200 OCT volumes was sufficient for SLIVER-net to achieve a mean ROC AUC of 0.89[CI: 0.86,0.92] and a mean precision-recall of 0.25[CI: 0.22,0.34], which is significantly better than the standard 3D CNN trained on the entire 1,202 OCT volumes available in our training cohort (mean ROC AUC 0.81[CI: 0.75,0.86]). With a sample size of 400 volumes, SLIVER-net achieved a mean ROC AUC of 0.93[CI: 0.90,0.95] and a mean precision-recall of 0.36[CI: 0.30,0.46] which is not significantly different from its top performance, which, as previously shown, was at the level of expert retina graders. In this case, SLIVER-net was able to achieve the state-of-the-art and expert-level performance with a sample size three

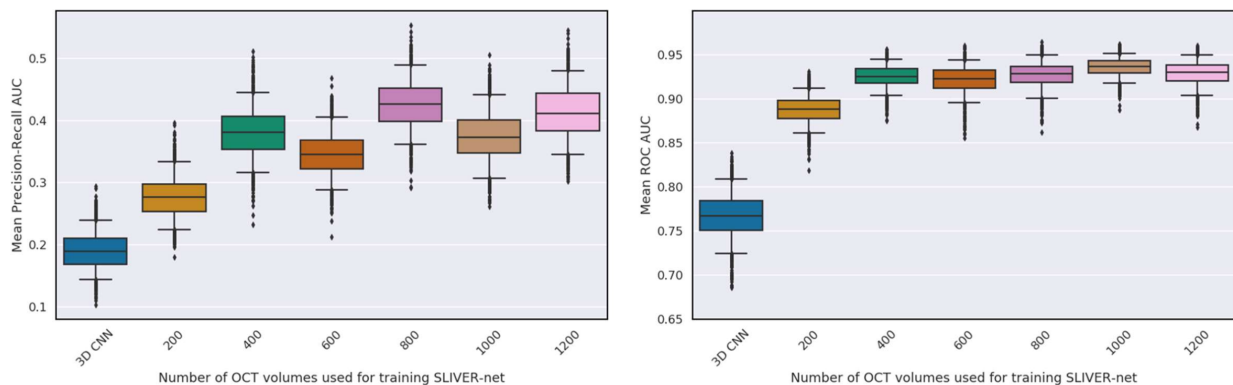


Figure 2.6: 3D CNN was trained on full data. Top: Mean precision-recall AUC across all biomarkers. Bottom: Mean ROC AUC across all biomarkers. Error bars represent 95% confidence interval (CI) calculated using a bootstrapping procedure.

### 2.2.6 Identifying traces of biomarkers outside of the macula

One advantage of deep learning is its ability to detect patterns without the usage of handcrafted features when given a sufficient amount of labeled data. In some cases, it is possible to annotate an object using one source, then train a model on a different one allowing the network to discover patterns unknown to researchers. This operation is useful when the information exists in the data but is unidentifiable by a human specialist.

In current practice, infrared reflectance (IR) imaging is commonly used to identify reticular pseudodrusen (RPD). RPD are now known to correspond to the subretinal drusenoid deposits which can be observed on OCT. Unlike IR images whose field of view is usually 30 degrees or larger, OCT volumes obtained in clinical practice are commonly limited to a 6x6mm (approximately 20 degrees) macular region centered on the fovea. RPD, however, are more frequently found in the more peripheral portions of the posterior pole outside of this 6x6mm macular region. As a result, these lesions will not be identified on review of the OCT alone, thus potentially leading to an underestimation of the risk of progression to late AMD in these individuals. To determine if this limitation could be overcome, we took advantage of companion IR images available with the OCT volumes in the Amish dataset and labeled these IR images for the presence of RPD. SLIVER-net successfully predicted the presence of RPD with an ROC AUC of 0.93[CI: 0.82,0.98] and precision-recall AUC of 0.40[CI: 0.22,0.61], significantly better than chance, using the OCT scans limited to the macula. This suggests the existence of patterns available in OCT scans that are still unknown to human specialists.

### **2.2.7 Transfer learning improves model performance**

Our training data consisted of 1202 annotated 3-dimensional volume images for biomarker prediction. Among these volumes, the prevalence of biomarkers ranged between 2 and 8 percent (see Table 1: Methods), while deep learning models generally require many more. A key component of SLIVER-net was flattening the OCT volume into an image by stacking the different slices into one long image (see Methods). This allowed us to incorporate a large publicly available dataset<sup>4</sup> using transfer learning, which is commonly used to address prediction problems when the amount of training data is small[PY10]. Under this paradigm, the model is “pre-trained” on a similar task, usually with a larger dataset. The model is then fine-tuned for the task at hand (see Methods for details).

SLIVER-net was pre-trained on the OCT dataset collected by Kermany et al., [KZG18]. This data consisted of 84,495 2D horizontal OCT B-scan images (e.g., slices) passing through

the fovea but were labeled with other ocular diseases (Choroidal neovascularization (CNV), Diabetic Macular Edema (DME) and Drusen). The pre-trained network was then fine-tuned for the biomarker prediction task.

We evaluated the performance of SLIVER-net with and without pre-training. Pre-training the model with the Kermany data (reported above) resulted in significantly ( $p < 0.001$ ) better performance when compared with training the model from scratch (mean ROC AUC 0.88[CI: 0.83,0.92]), mean precision-recall AUC 0.24[0.20,0.33], Figure 2.7).

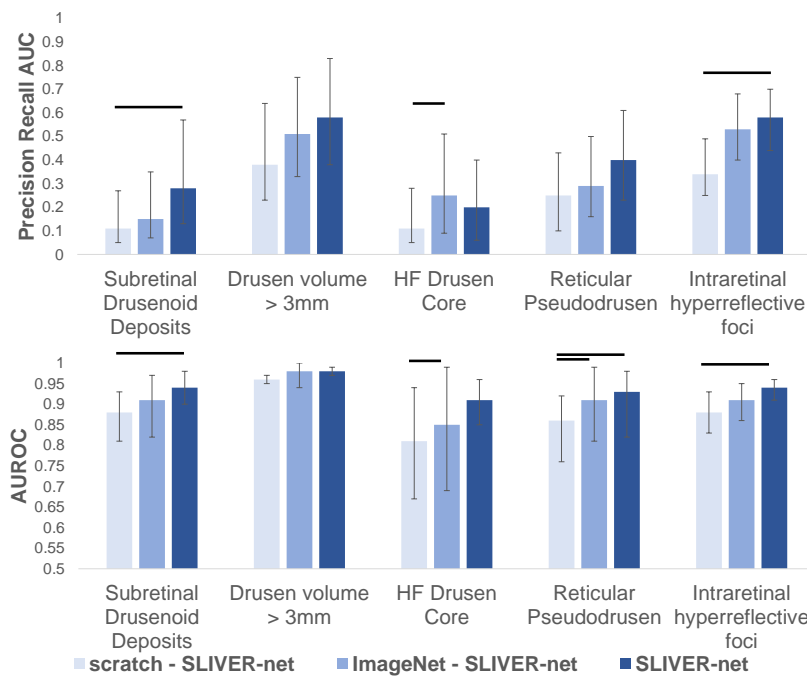


Figure 2.7: Comparison of external datasets and their effect on performance. SLIVER-net trained from scratch (light blue), pre-trained using ImageNet (blue), and pre-trained using Kermany (dark blue). Top: Precision-recall AUC for each biomarker. Bottom: ROC AUC for each biomarker. Horizontal bars indicate statistically significant differences ( $p < 0.05$ ). Error bars represent 95% confidence interval (CI) calculated using a bootstrapping procedure.



### 2.2.8 The tradeoff between quantity and quality of external data

The effectiveness of the transfer learning procedure depends on the size of the external data, as well as its similarity to the target task. While the Kermany data above contained nearly 85,000 OCT scans, there are even larger but less related datasets. Natural images from the ImageNet [DDS09] dataset (over 1 million samples with 1000 classes) may provide a good foundation for the transfer learning approach based on the sheer volume of training data. We thus compared the performance of SLIVER-net pre-trained with data from Kermany et al. 2018 (Kermany-SLIVER) against the same model pre-trained with data from ImageNet (ImageNet-SLIVER). Kermany-SLIVER outperformed ImageNet-SLIVER with a mean ROC AUC of 0.94[CI: 0.91,0.96], and a mean precision-recall AUC of 0.41[CI: 0.34,0.51] compared with 0.92[CI: 0.87,0.95] ( $p < 0.01$ ) and 0.35[CI: 0.30,0.45] respectively (see Figure 2.7) despite the difference in the number of exemplars. This is in line with recent findings[WKW16] that while training set size is essential, it is beneficial in terms of performance to pre-train networks using related data.

### 2.2.9 Robustness to the number of slices available in each volume

OCT acquisition parameters are not standardized in current ophthalmic practice. Notably, retina practitioners may determine the number of slices (B-scans) to acquire on a patient-by-patient basis, resulting in volumes with differing resolution and field of view. While the data acquired in this study were of the same resolution and field of view, we simulated scans of different field of view and resolution in order to assess SLIVER-net’s robustness to such changes.

First, we artificially varied the field of view around the macula available in each volume (see Methods) and observed that varying the field of view did not significantly affect performance for biomarker prediction (Figure 2.8). Then, we simulated different B-scan resolutions by down-sampling the number slices in each volume (see Methods), again observing that varying

volume resolution did not significantly affect the model’s performance for biomarker prediction (Figure 2.9). In both scenarios, we have observed that SLIVER-net was robust to different sizes and resolutions of OCT scans, making it useful in various clinical scenarios and under different resource constraints.

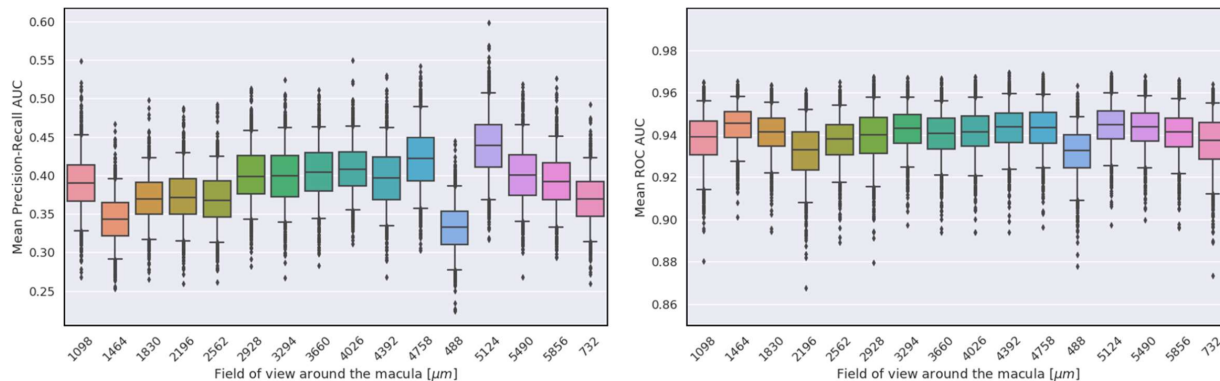


Figure 2.8: Left: Mean precision-recall AUC across all biomarkers. Right: Mean ROC AUC across all biomarkers. Error bars represent 95% confidence interval (CI) calculated using a bootstrapping procedure.

## 2.3 Discussion

The application of deep learning to new studies depends on the ability to train models with limited data. In this work, we developed a new deep learning technique, SLIVER-net, to predict clinical features from OCT volumes. Our approach provides these predictions using a relatively small number of annotated volumes (hundreds), and an even smaller number of positive training examples. SLIVER-net is based on two main ideas. First, we use transfer learning to borrow information about the structure and parameters of the network from publicly available large datasets. Unfortunately, there are no large datasets that include volumes, and we, therefore, use transfer learning using the 2D images. In order to account for this, our second idea is to model the volume as a 2-dimensional image by tiling the volume

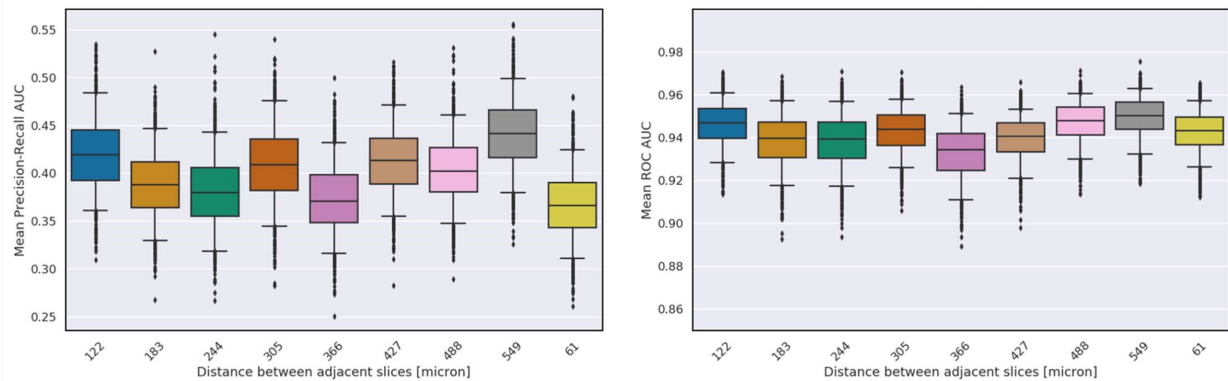


Figure 2.9: Left: Mean precision-recall AUC across all biomarkers. Right: Mean ROC AUC across all biomarkers. Error bars represent 95% confidence interval (CI) calculated using a bootstrapping procedure.

scans, and then adding to the neural networks additional layers that take into account the fact that two adjacent images in the tiled image are adjacent in the original 3D volume.

We demonstrate our approach using OCT volumes, which are widely used in current ophthalmic practice. Specifically, we used SLIVER-net to identify clinically useful OCT biomarkers which have been shown to predict the risk for progression to late AMD[Fau18]. We found that for most features, SLIVER-net was able to identify these AMD-related biomarkers in agreement with senior expert clinician graders and was superior to junior graders. In many cases, as revealed by a post-hoc review, SLIVER-net identified additional biomarkers that were missed by the initial annotation. SLIVER-net is considerably more powerful than standard deep learning techniques used for medical volumes such as 3D CNNs. Despite having very few annotated samples from our original dataset, SLIVER-net was able to outperform the current state of the art methods. Particularly, our approach significantly improved the average AUC from 0.81 achieved by 3D CNNs to 0.94 achieved by SLIVER-net. The models were compared using an external test set acquired at a separate institution, which, in contrast with single-site and single-dataset studies, provides support that SLIVER-net can be portable across institutions.

At a practical level, SLIVER-net provides a general framework for addressing prediction problems with a limited sample size of labeled data. Its success was primarily driven by transfer learning and slice integration, both of which are not limited to biomarker prediction nor OCT classification. Thus SLIVER-net presents a feasible approach to the application of deep learning to new problems involving 3-dimensional imaging modalities. While typical machine learning solutions cite requirements in the tens of thousands in terms of training samples, our investigations showed that SLIVER-net approached maximum performance with only 400 training samples (Figure 2.6), which more closely matches sample sizes required for clinical validation. Using the transfer learning framework, predictive and data-driven applications can potentially be pursued concurrent to clinical validation without devoting additional resources to annotation.

The early application of deep learning and automated image analysis to relatively new imaging modalities such as OCT can also provide a synergistic development at technical and clinical levels. We included reticular pseudodrusen (RPD) as a biomarker of interest because it is a lesion which may in some cases be present only beyond the typical macular OCT scanning field commonly used in clinical practice, and is thus instead detected using larger field of view infrared reflectance imaging. Interestingly, SLIVER-net was able to successfully detect the presence of RPD using the smaller field macular OCT information alone, which suggests that lesions which fall outside the macula may be associated with subtle alterations in the macula which remain to be understood. Future work utilizing multiple imaging modalities that are available for use in clinical practice, may reveal other novel findings which may be encoded in the OCT data.

The ability to automatically identify these high-risk biomarkers for AMD progression has important clinical implications. Lei et al[LBA17] have already shown that the presence of these biomarkers can be translated to a simple score that can risk stratify patients presenting to the clinic. Automated biomarker detection could lead to a more precise quantification of not only the presence but the extent or severity of the biomarker or feature of interest, which

could further improve the predictive accuracy of the biomarker[Nas18]. Such a risk score could be used to prognosticate disease and to define appropriate intervals for follow-up and monitoring. This is particularly relevant as home OCT devices are now becoming available for telescreening. In addition, such a scoring system could also be used to identify high-risk patients for enrollment in clinical trials for early intervention therapies. Automated biomarker detection could also prove to be invaluable in a number of research applications such as the study of the appearance and evolution/progression of these biomarkers in large AMD datasets. Investigations such as this may provide new insights into the pathogenesis of AMD.

## **2.4 Methods**

### **2.4.1 Data**

#### **2.4.1.1 Biomarker prediction data**

OCT scans were acquired from 1,007 patients as part of a longitudinal study on AMD progression in an elderly Amish population. These scans were acquired from three different sites: University of Pennsylvania (390 patients), Case Western Reserve University(248 patients), and University of Miami (369 patients) using the Spectralis system (Heidelberg Engineering). The research was approved by the institutional review boards (IRBs) of the respective institutions and all subjects signed written informed consent. All research was conducted in accordance with the tenets set forth in the declaration of Helsinki. All imaging data were transferred to the Doheny Image Reading Center (DIRC) in a de-identified fashion. The image analysis research was approved by the UCLA IRB. Two volumes (97 B-scans, with an in-plane resolution of 496 x 512 and dimension of 6x6mm on the retina – roughly a 20-degree field of view) were acquired from each patient. Only scans that were determined to be good quality, as assessed by a senior retina image grader at the Doheny Image Reading Center, were used for model development and validation. Under this criterion, we excluded

72 volumes, resulting in 1942 OCT volumes in total. Data from the University of Miami and Case Western Reserve University (1202 volumes) were used for model training, and data from the University of Pennsylvania (740 volumes) were withheld for testing.

Four biomarkers (hyperreflective foci, hyporefective cores within drusen, subretinal drusenoid deposits, and high central drusen volume), and reticular pseudodrusen as identified using IR imaging, were selected for this study. A single retina specialist reviewed each Spectralis OCT volume, manually recording the presence of hyperreflective foci, hyporefective drusen cores, and subretinal drusenoid deposits. The remaining two biomarkers were identified using different devices. The Cirrus OCT system (Zeiss) was used to quantify central drusen volume, and reticular pseudodrusen were identified using an infrared reflectance image. In accordance with previous publications[LBA17, Nit19], a high central drusen volume was determined to be a value of  $\geq 0.03\text{mm}^3$  within the central 3mm zone centered on the fovea. It is important to emphasize that the Spectralis system cannot produce a drusen volume measurement, though the drusen are visible on the OCT. In addition, while subretinal drusenoid deposits (SDD) evident on OCT appear to correspond to reticular pseudodrusen (RPD), RPD are commonly present only outside the macula, and thus RPD may be present on an IR image (which covers a 30-degree field of view) without evidence of visible SDD on the OCT. Table 2 summarizes the prevalence of these biomarkers within the dataset.

The OCT volumes of 91 patients randomly selected from University of Pennsylvania were annotated by an additional three junior reading center clinician graders. These labels were used to assess inter-rater reliability as well as model comparison.

#### **2.4.1.2 Transfer learning data**

We compiled two external datasets to pre-train our models. One dataset was ImageNet [DDS09], which consists of millions of training images comprised of a thousand object categories. ImageNet has been commonly used in transfer learning applications for natural images, and it has been shown that models pre-trained on ImageNet perform well on other do-

Table 2.2: The biomarkers used for this study and their prevalence throughout the three datasets.

	Training set	Testing set	Total
Number of patients	617	390	1007
Number of OCT volumes	1202	740	1942
Hyperreflective foci (IHRF)	89 (7.4%)	49 (6.6%)	138 (7.1%)
Hyporefective drusen core (hDC)	33 (2.7%)	13 (1.8%)	46 (2.4%)
Subretinal drusenoid deposits (SDD)	23 (1.9%)	13 (1.8%)	36 (1.9%)
High central drusen volume	40 (3.3%)	19 (2.6%)	59 (3.0%)
Reticular pseudo drusen (RPD)	41 (3.4%)	20 (2.7%)	61 (3.1%)

mains [OBL14, Shi16].

We also acquired a large collection of publicly available OCT images collected by Kermany et al., 2018, which we simply refer to as “Kermany”. In this dataset, 84,495 horizontal B-scans passing through the foveal center (i.e., typically the middle slice of an OCT volume) were annotated for one of four conditions: normal, choroidal neovascularization (CNV), diabetic macular edema (DME) and drusen. While there were less than 100,000 samples in this dataset, they were more similar to our biomarker prediction data.

### 2.4.1.3 Data Preprocessing

Each slice of the volume was resampled from 496 x 512 pixels to 224 x 224 pixels[HZR16]. Then, image contrast was enhanced by clipping pixel intensities to the 2nd and 98th percentile, and resulting values were rescaled between 0 and 255.

## 2.4.2 3D CNNs

Convolutional neural networks (CNNs) comprise out of many kernels that receive an image as input and produce a representation that is most meaningful for a given task using an operation called convolution. 3D CNNs extend this approach to three-dimensional objects and are commonly applied to volume analysis. They have gained popularity in biomedical imaging (e.g. CT[HSV17, JSB19], MRI[MNA16, Kam17, DDA18, Val17]) due to increasingly capable hardware. We used a 3D version of Resnet18[HZR16] to compare against the 2D approach. The input to the network was a 3D volume of size 224x 224 x 97 and the output was a prediction score range 0 to 1 for each biomarker representing the probability the respective biomarker is present.

## 2.4.3 SLIVER-net Architecture

Our proposed approach, termed SLIVER-net, was comprised of three steps. First, the preprocessed OCT volume was passed through a “backbone” convolutional neural network (CNN), which represented the scan in an abstract feature space. Then, a slice aggregation operation was applied in order to compress this representation and capture information that is shared across adjacent slices. Finally, a decision module operated on this compressed representation to determine the presence or absence of biomarkers.

### 2.4.3.1 Step 1: Backbone networks

CNN models contain several convolutional layers stacked together (i.e., each layer’s output serves as the next layer’s input) to extract a feature map from a single image. Previous work[YCN15, EBC09] has shown that the first CNN layers (lower layers) of a deep learning model generally identify abstract features (lines, edges, corners) and the upper layers identify features that are more task-specific. In our experiments, all tested models were based on the same CNN architecture, Resnet18[HZR16]. 2D backbones (SLIVER-net) used 2D kernels



(size  $3 \times 3$ ,  $7 \times 7$ ) while the 3D-CNN backbone used 3D kernels (size  $3 \times 3 \times 3$ ,  $7 \times 7 \times 7$ ). Resnet18 was chosen since it has shown to perform well in the natural image setting [HZR16]. This model represents each  $224 \times 224$  OCT slice as an  $8 \times 8$  image.

Feature extraction on all 2D slices was computed in one forward pass. To do this, each of the 97 slices was concatenated vertically, forming a “tiled” image of  $(97 \times 224) \times 224$  (see Figure 2.10) that was passed to the model. The output of the backbone model was a  $(97 \times 8) \times 8$  image with 1024 features for each of the 97 slices.

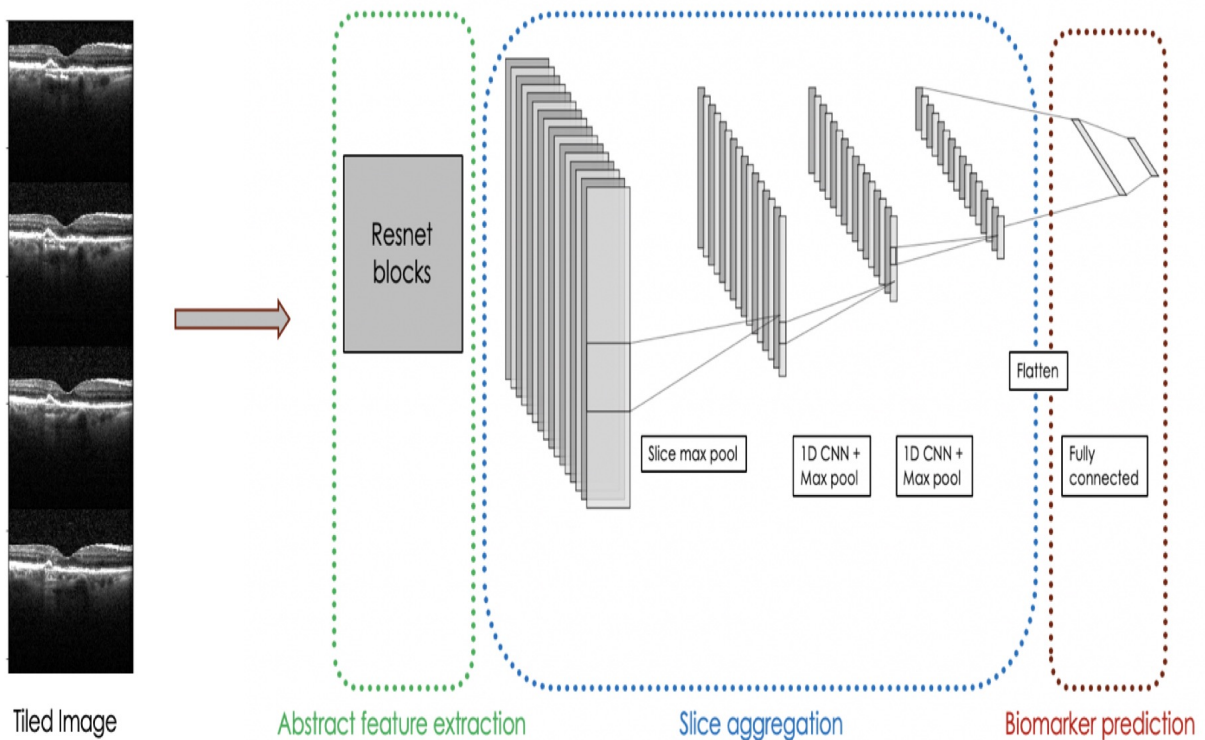


Figure 2.10: Our model operated on a 2d tiling of the OCT volume. Resnet18 served as the abstract feature extractor, and the representations for each slice were aggregated using slice integration and a 1D CNN. Finally, biomarkers were predicted using fully connected layers.

### 2.4.3.2 Step 2: Slice Integration

In a deep learning model, the final feature map produced by the CNN layers is collapsed into a feature vector, usually by taking the average across all spatial dimensions in an operation referred to as global average pooling[LLG15]. This “flattens” the feature map such that it can be passed to a decision module. We extended this operation by taking both the maximum (“max pooling”) and average (“average pooling”).

However, we observed that applying this operation globally would remove the model’s access to the local 3D structure of the OCT volume. In order to preserve correspondence among neighboring slices, we performed average and max pooling within each of the 8 x 8 backbone outputs, producing a 97 x 1024 representation of the volume. Then, a small 1D CNN was added to aggregate these slices before they were passed to the decision layer. This Slice Integration procedure was a primary driver of the success of SLIVER-net.

### 2.4.3.3 Step 3: Decision module

Biomarkers were predicted in a multi-task approach, in which the single network simultaneously predicted the presence of all targets. Our prediction “head” consisted of only one hidden layer with 1024 hidden units, feeding to an output layer of 5 units with a sigmoid activation function, corresponding to the biomarkers. By simultaneously optimizing for separate tasks, the multi-task paradigm provides an implicit regularization, improving generalizability[CQY16, AEP07].

### 2.4.4 Training

Data acquired from the University of Miami and Case Western Reserve University were used to develop the models. These data were randomly split into training (80%) and validation (20%) sets. Models were implemented using PyTorch[PGM19] and optimized using the Adam optimizer with default parameters[KB14] and a weight decay of 0.01. For each model, the

learning rate was chosen from values between 1.0 and  $10e-7$  using the learning rate finder implemented in the Fastai library[HG20]. Models were trained with a batch size of 32, and training continued until validation loss stopped decreasing for 20 consecutive epochs (i.e., passes through the training dataset). The model weights that achieved the lowest loss on the validation set during training were chosen for evaluation on the test set.

### **2.4.5 Transfer learning**

One limitation of the Resnet and other CNN feature extractors is that they require a large amount of data to train. A typical solution to this is to apply transfer learning[PY10], in which the network is first trained on an existing but similar dataset, and then “fine-tuned” on the dataset of study.

#### **2.4.5.1 Model pretraining**

We evaluated the ImageNet and Kermany datasets for their suitability for transfer learning. While ImageNet is a much larger dataset, the Kermany set. consisted of OCT images similar to our data.

The original labels for the candidate datasets (image classification for ImageNet, and disease diagnosis for Kermany), were not aligned with our biomarker prediction task. However, it has been observed[YCN15] that some convolutional neural networks extract general features applicable to most visual tasks. We used the following approach to apply transfer learning to biomarker prediction: (1) We trained a network for the original task of the auxiliary dataset. For both datasets, a Resnet18 feature extractor was trained for its respective task (object classification or disease classification) for up to 50 epochs (with early stopping) and a learning rate of  $1e-3$ . (2) We discarded the decision layers, which were specialized for the auxiliary task, and (3) replaced the decision layer with a randomly initialized one appropriate for the target task. (4) Only the new decision layer for biomarker prediction was then trained with

our training set without updating any of the parameters in the feature extractor. 5. Finally, the whole network was updated using a reduced learning rate of  $1e-5$ .

#### **2.4.6 Model evaluation**

Model performance on the test set was quantified in terms of the area under the receiver operating characteristic (ROC) curve, as well as the Precision-Recall (PR) curve. A 95% confidence interval was estimated for model performance using a bootstrapping procedure. For each bootstrap iteration, we randomly resampled from the test set with replacement and calculated performance metrics. We repeated this 5000 times and selected the 125th and 4,875th values of the sorted list to define the 95% confidence interval. Performance metrics were compared using a Wilcoxon signed-rank test [Wil45] (i.e., nonparametric t-test).

#### **2.4.7 Model explainability**

In clinical settings it is of high importance for statistical models to communicate some rationale behind decision making in order to build trust between the machine learning algorithm and the clinical user. To address this issue, we provide explainability maps along with its predictions to show important regions as inferred by the algorithm (Figures 2.4,2.5). The explainability maps are produced by visualizing the backbone representation of each OCT volume. The representation of each slice (an  $8 \times 8$  feature map with 512 channels) is averaged across channels to create an  $8 \times 8$  feature image for each scan (97 total scans in each volume), which shows the average local importance across all channels. Then, the feature image is interpolated to match the sizes of the original input. The feature image and the original input are shown together to produce the explainability map.

#### 2.4.8 Simulating model performance with different acquisition parameters

We assessed the robustness of SLIVER-net to different acquisition parameters by artificially varying the OCT volumes. In each case, we trained SLIVER-net on the transformed data and observed performance on the test set with the same transformation.

To manipulate field of view, we used various numbers of slices taken around the macula. We evaluated performance when 9 central slices (488 microns) were available up to 97 slices (5856 microns). Then, to evaluate the SLIVER-net’s performance on the resolution of each volume along the Z-axis, i.e., the distance between two nearby slices, we used different sampling rates to down-sample the number of slices in each volume, thus simulating lower-resolution OCT volumes. We varied the distances between two nearby slices from 61 microns (97 slices total, the standard resolution of this study) up to a range of 549 microns between each (11 slices per volume).

## CHAPTER 3

# Augmenting Clinical Practice Guidelines Using Reinforcement Learning and Causal Transformers

### 3.1 Introduction

Clinical care routinely involves planning treatment for patients which includes carefully considering potential risks and benefits of the treatment options. Clinical practice guidelines (CPGs) published by medical associations are based on the best available population-level evidence and are intended to assist healthcare professionals in making clinical decisions. However, CPGs may be ambiguous or sub-optimal when considering polychronic patients, that suffer from multiple intersecting chronic conditions. These complexities pose challenges because CPGs are oriented to single conditions, and it is left to clinician judgement to adjudicate between conflicting recommendations from multiple guidelines. For example consider an aging population that exhibits increasing clinical complexities and care demands, resulting in patterns of super-additive costs when diseases interact [CV]. Application of disease-specific CPGs to patients with multiple diseases can lead to competing recommendations and the potential for adverse drug-drug or drug-disease interactions. For example, medications indicated for heart failure could compromise kidney function in those with kidney disease, or, nonsteroidal anti-inflammatory drugs (NSAIDs) may be suggested to treat osteoarthritis pain, but they have relative contraindication in patients with a history of peptic ulcers disease. To account for the patient's unique circumstances, such as demographics, family and disease history, or individual physician practice patterns, doctors may deviate from applicable

guidelines partially or fully. While these deviations may be appropriate in certain cases, they can also lead to unwarranted variation and poorer health outcomes. In contrast to deviations that manually personalize clinical care, deviations may also result from professional uncertainty, such as lack of specialized domain expertise or uncertainty about treatment options [Wen]. This presents a unique opportunity for AI-based solutions to learn from CPGs and their deviations in order to design better clinical practice guidelines that can direct choices supporting better health outcomes.

Applying Reinforcement Learning (RL) approaches to clinical decision support systems is an active area of research [LSN]. Our aim is to improve the fundamental problem of using RL to enhance clinical practice guidelines. Specifically, offline Reinforcement Learning is an AI approach suited to ingest clinical histories and learn treatments whose patient outcomes can only be observed in the future. Several studies have demonstrated that offline RL offers a promising framework for utilizing longitudinal medical data to generate accurate medical recommendations that can include both chronic and acute conditions. In cancer treatment, researchers have used offline RL to optimize radiation therapy scheduling [YS, TLC], to determine optimal chemotherapy dosages [ZZS, Hum], and to personalize cancer treatment policies [LSP]. In addition, offline RL has been proposed for treatment planning in Parkinson’s disease [WKV] and Type 2 diabetes treatment [OPL]. In critical care, offline RL has been applied to medication recommendations for sepsis [PDW, RKA] and management of mechanical ventilation in intensive care units [PCC].

These approaches have been effectively applied in data rich environments such as the ICU (Intensive Care Unit), or operating room. However, outside these settings, there are several limitations. First, they assume that a patient’s health state can be measured in regular intervals, which is often not met in practice. Specifically, this assumption of regularly observed time series data is critical for the applicability of Markov Decision Processes. Second, although clinical guidelines contain many treatment options, such as lab requests, medications, and medical procedures, previous research only included a limited number of possible treatments

such as one or two drug dosages. Thus, these approaches are currently not suitable for many common clinical scenarios, such as the management of chronic disease for polychronic patients, where patient data is observed very irregularly, and the number of combinations of treatment and other clinical actions is large.

To address these limitations, our approach relies on recent work [CLR, JLL, HSZ] that utilizes the sequential nature of decision-making to treat the RL problem as a conditional sequence modeling problem and leverage the transformer’s ability to effectively capture long-range dependencies while also synthesizing disparate well-performing actions across many medical histories. However, both the decision transformer (DT) and the trajectory transformer (TT) only optimize indirectly for the expected outcome. In its learning process, the decision transformer optimizes for the actions observed in the data, in effect doing behavioral cloning and at test time, trajectory optimization is done by conditioning a high target return. The trajectory transformer also does not optimize directly for the expected outcome. Instead, it learns to predict the reward and actions. Given the likely actions it then uses beam search to identify the trajectory that produces the highest reward. We propose the Importance Transformer (IT) that directly optimizes the expected outcome by defining a novel loss function that incorporates weighted importance sampling.

We demonstrate the effectiveness of IT using a cohort of patients with either type 2 diabetes (T2D) or heart failure or both. Using weighted importance sampling, a common statistical method used to evaluate offline policies, we estimate the potential outcome of diabetes’ severity level, heart failure (HF) severity level, and a combination of both. We compared IT with current state-of-the-art offline RL transformer-based methods and observed superior performance on T2D, and T2D + HF, and on par performance for HF.

As an assistive clinical decision support system, it is important to establish a trustworthy grounding of our modelling approach and our evaluation. Therefore, we conducted several diagnostic tests to challenge our modelling assumptions. These include a sensitivity analysis of weighted importance sampling for behavioral policy estimation and a qualitative analysis to



examine whether sequential information is in fact a significant factor in the policy’s decision process. Both tests highlight the trustworthiness of our modelling and evaluation approach. In summary, we showed that the approach presented in this paper can support clinical decision making in cases where CPGs are not fully defined or integrated in polychronic settings. This is accomplished by learning personalized treatment pathways and thus absorbing warranted deviations and potentially reducing uncertainty in treatment options.

## 3.2 Results

In this study, we introduced an offline reinforcement learning approach called the Importance Transformer (IT). This GPT-based architecture utilizes both static information (sex, age) and sequential information, as captured by the medical state, and performs actions at each medical encounter to produce action recommendations that results in optimal clinical outcomes. To do so, IT uses a composite loss function that encourages it to recommend actions to enhance CPGs that are associated with high rewards in the future while preventing it from deviating too much from actions that were observed in the data.

For training and evaluation of our approach, we analyzed electronic health record (EHR) data obtained from multiple provider groups and hospital systems in the United States. For more details on our cohort see the method section.

We first compared IT with two state-of-the-art offline reinforcement learning methods, the Decision Transformer (DT) and the Trajectory Transformer (TT). To do so, we used estimated expected reward as the evaluation criteria where four different outcomes were examined as rewards: disease severity level for both T2D, HF, and a combination of HF and T2D. Second, to solidify our findings, we performed a sensitivity analysis for our weighted importance sampling evaluation criteria. We analyzed the method’s sensitivity to the quality of the estimated behavioral policy which is the policy used to generate the observed actions. Finally, we present a permutation analysis to evaluate the importance of temporal information

	Mean	Median	Min	Max
Encounters per patient	24.4	23	48	10

Table 3.1: Encounter Statistics

	Sex		Age			
	Female	Male	49 -	49 - 64	64 - 75	75 +
Prevalence	55	45	8.0	31.0	29.4	31.6

Table 3.2: Demographics Statistics

for our IT approach.

### 3.2.1 Data Summary

In our data we consider 803,746 patients with a total number of health care encounters of 19,627,746, all of whom had some level heart failure (HF) or diabetes (T2D), as determined by internal mapping based on procedures labs and medications. In addition, each patient had at least 10 encounters between 2014 and 2018. A statistical summary of the encounter statistics is given in table 3.1.

Table 3.2 provides a summary of the demographics represented in the data in terms of gender and age.

And finally, table 3.3 gives a summary of the prevalences for each diabetes and heart failure severity level.

Patients selected were of ages between 35 and 89 years during the data collection window and have at least one clinical encounter in at least two successive years. Patients with Diabetes and Heart Failure were identified using diagnosis (ICD10, Symmetry Episode Treatment Group maps), procedure (CPT4), drug (PCC and DCC from Symmetry Drug code hierarchies) and lab values (LOINCs (Logical Observation Identifiers Names and Codes)).

	0	1	2	3	4
Diabetes Prevalence (%)	55.8	27.0	13.1	3.8	0.3
HF Prevalence (%)	64.3	16.9	12.2	5.5	1.1

Table 3.3: Prevalence of Disease Severity Levels

All diagnosis codes were converted to ICD10 codes using the General Equivalence Mapping (GEM) published by CMS. The clinical timelines of the patients showed different severity levels for both these conditions labeled using four discrete levels 1-4 with 1 being the lowest. Given the variation in coding and to prevent temporal spikes in severity overly influencing the severity labeling, a simple smoothing approach was used wherein the severity of the patient is the max severity over the last 8 weeks with at least two occurrences.

### 3.2.2 Evaluating expected outcomes

A key metric of our model is that of expected outcomes (see equation 3.1). Comparing expected outcomes across state-of-the-art transformer-based models, namely DT and TT, demonstrate superior performance along several reward dimensions. Specifically, as shown in table 3.4, IT outperforms DT and TT with the rewards for diabetes and the combined reward of diabetes and heart failure. For the reward related to heart failure alone, IT is found to be on par with DT and TT. Moreover, IT achieves a higher expected reward than the average reward in our cohort for diabetes, and the combination of diabetes and heart failure.

To further increase the trust in our finding, we have also conducted a diagnostic check by assessing the performance of a random decision model. Intuitively, we expect a random decision policy to exhibit a lower-level performance compared to all AI policies (IT, DT, TT) as well as the average reward in our cohort. Interestingly, this is not the case in all scenarios (see table 3.4). In all scenarios, it is indeed the case that the random decision policy performs worse than the average reward in our cohort. Further, for the scenarios involving disease

Model	T2D	HF	HF+T2D
IT	<b>2.65</b> [2.31, 2.67]	<b>0.38</b> [0.31, 0.45]	<b>2.95</b> [2.58, 2.99]
DT	0.36 [0.36, 0.37]	0.28 [0.27, 0.3]	0.64 [0.63, 0.65]
TT	0.87 [0.38, 1.09]	0.35 [0.21, 0.66]	-0.64 [-1.11, 1.43]
Random	0.09 [0.04, 0.13]	0.09 [-0.09, 0.15]	0.22 [0.04, 0.3]
Observed	0.35	0.28	0.63

Table 3.4: Summary of estimated expected outcomes for different models across different rewards with 95% confidence intervals (best values highlighted).

severity, all AI policies out-perform the random one with the notable exception of TT in the combined diabetes and heart-failure reward model. With respect to expected outcomes, TT exhibits unstable behaviors with a large variance and the median performance matching the performance of a random policy.

### 3.2.3 Sensitivity of WIS for behavioral policy estimation

So far in this study, we have used the WIS method to estimate the expected outcome of a suggested policy. To make such estimates, the WIS method relies on the ratio between the suggested policy,  $\pi(\alpha)$ , and the observed policy,  $\pi_{\text{obs}}(\alpha)$ . However, since the observed policy is only estimated, WIS might be sensitive to the quality of this estimator, and we are presented with the challenge to identify what part of the estimated expected reward is due to the decision policy and what part is due to the quality of the estimation of the behavioral policy. To address this, we employ another diagnostic check to substantiate the trust in the results presented in this paper by studying the quality of the observed policy estimator. Specifically, we isolate the quality of the estimator of the behavioral policy  $\pi_{\text{obs}}(\alpha)$  and employ as before a random decision policy.

Intuitively, employing a random policy we expect to be worse off than the observed average

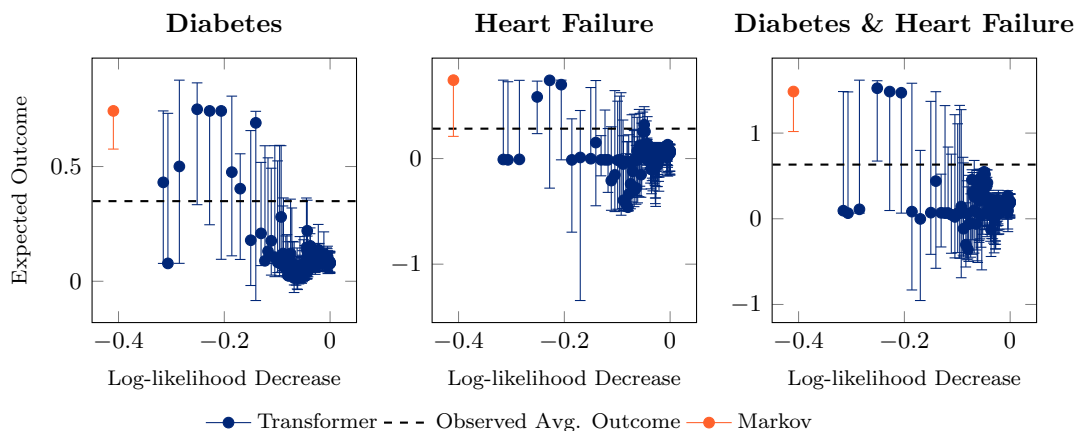


Figure 3.1: Sensitivity Analysis of expected outcomes under different modelling assumptions for the observed policy.

outcome of the cohort. Any deviation from this expectation is due to the quality of the behavioral policy’s estimation. A good estimation of the behavioral policy would therefore result in a significantly lower expected outcome than the observed average outcome of the cohort. Contrary, a bad estimation of the behavioral policy should exhibit high variance and potentially even higher expected outcomes than the observed average outcome of the cohort.

With these considerations in mind, we sought to measure the WIS expected outcome estimate for a random policy using models of varying quality as the estimated behavioral policy. To generate these models, we trained a transformer-based model and saved its weights throughout the training process, resulting in multiple models with a gradually increasing performance, as measured by the data log-likelihood. The observed policies are modelled using 1-step Markov and Transformer models. We then evaluated the WIS expected outcome estimates for a random policy using these models.

Our results, presented in figure 3.1, demonstrate that the WIS estimated expected outcome in the convergence area (0%-5% on the x-axis) where the model achieved highest performance, is significantly below human performance, with a narrow confidence interval. However, as we move away from this area on the x-axis, and the model’s quality decrease, the WIS

estimation of a random policy produces outcomes that are no longer significantly better than CPG’s performance with increasingly wider confidence interval. We also notice that a basic Markovian model that has a log-likelihood that is 40% smaller than that of the converged transformer, produces a WIS estimate far above the observed average outcome of the cohort. This confirms that a first intuitive approach to model the behavioral policy using a Markovian model is too simplistic as it is not able to capture the dynamics in the data well enough to produce reliable and trustworthy decision policies.

### **3.2.4 Permutation testing of temporal importance**

So far, we have based our modelling assumptions on being able to ingest and handle sequential data. However, we have not surfaced any insights on whether our decision policy relies on sequential data or only certain events in medical history. One strategy may be to analyze the attention weights, yet a simpler and more quantitative approach to reveal whether temporal information affects the decision policy is to conduct a permutation test on the sequences fed into our model. Concretely, this involves randomly shuffling the order of patient encounters and using IT to generate recommended actions based on the shuffled data. Importantly, this also involves retraining the behavioral policy to account for the distribution shift of the actions given the health state. In our case, we retrained the behavioral policy on the original data without any positional embedding, which removes time information. To justify this decision to retrain the behavioral policy, we compared the log-likelihood of the original and retrained behavioral policies on the shuffled data. Our results show that the log-likelihood of the retrained behavioral policy was significantly higher. Consequently, the results presented for the shuffled data will be based on the retrained policy.

Naturally, we expect the expected reward of the decision policy to be lower than the expected reward of the original data. Thus, it would confirm that the sequential information inherent in medical records can be successfully exploited by IT.

We repeated the permutation test 100 times (see table 3.5) and observed that the expected

Reward	P-value (100 permutations)
T2D Severity Level	0.01
HF Severity Level	0.06
T2D + HF Severity Level	0.01

Table 3.5: P-values indicating the significance of sequential data in modelling next best actions.

outcome of the algorithm on the permuted data was never found to exceed its performance on the original data for T2D, and T2D+HF, resulting in a test significance level of  $p < 0.01$ . For HF, the expected outcome of the algorithm exceeded its performance on the original data six times resulting in a test significance level of  $p < 0.06$ . These findings suggest that temporal factors have a strong impact on the performance of NBA (Next Best Action) for T2D, and T2D+HF and a weaker impact for HF.

### 3.3 Methods

#### 3.3.1 Data

This study used de-identified Electronic Health Records (EHR) data between 2014 and 2018 from the OptumLabs Data Warehouse (OLDW) [WSD]. The database contains longitudinal health information on enrollees and patients, representing a mixture of ages, ethnicities, and geographical regions across the United States. The data in OLDW include diagnosis, procedure, drugs (prescribed and administered) and laboratory results for over 60 million patients. Because data was de-identified or a Limited Data Set in compliance with the Health Insurance Portability and Accountability Act and customer requirements, the UnitedHealth Group Office of Human Research Affairs deemed that the study was not subject to ongoing Institutional Review Board oversight.

**Time aggregation:** We combine all encounters that took place within a single month into a singular encounter. In terms of actions, this implies that all distinct actions that occurred within the same month will be included as part of the same encounter. As for severity levels, we select the highest severity level recorded that month as the severity level for that encounter. This consolidation of actions and severity levels provides a concise representation of the data, allowing for statistical signal processing through machine learning while retaining enough temporal information.

**Rewards:** In our analysis, we have employed a variety of reward structures. Specifically, for hospitalization or emergency room encounters, we assigned a reward of -1 if such an event occurred during the encounter and 0 otherwise. In the case of diabetic or heart failure severity level, we set the reward to be  $+k$  if the severity level decreased by  $k$  levels, and  $-k$  if it increased by  $k$  levels. These diverse reward mechanisms enabled us to comprehensively evaluate the performance of our system across a range of clinical scenarios.

### 3.3.2 Modelling

IT is designed in a comparable way to that of the Decision Transformer [CLR] but with several modifications. Figure 3.2 shows an overview of the IT architecture.

### 3.3.3 Model Input

**Action representation:** We assigned each possible combination of actions from the available action space with a unique action tokens (AT). This allows the model to create a joint distribution where each combination of actions can be assigned with a probability score during model training and inference. To incorporate some measurement of safety, following [KCB] we excluded action combinations which were rarely observed in the data using the following exclusion criteria: 1) we generated a conditional distribution for each action token,  $\pi_0(a_t, s_t)$ , conditioned on severity levels. 2) Each AT that did not pass the probability



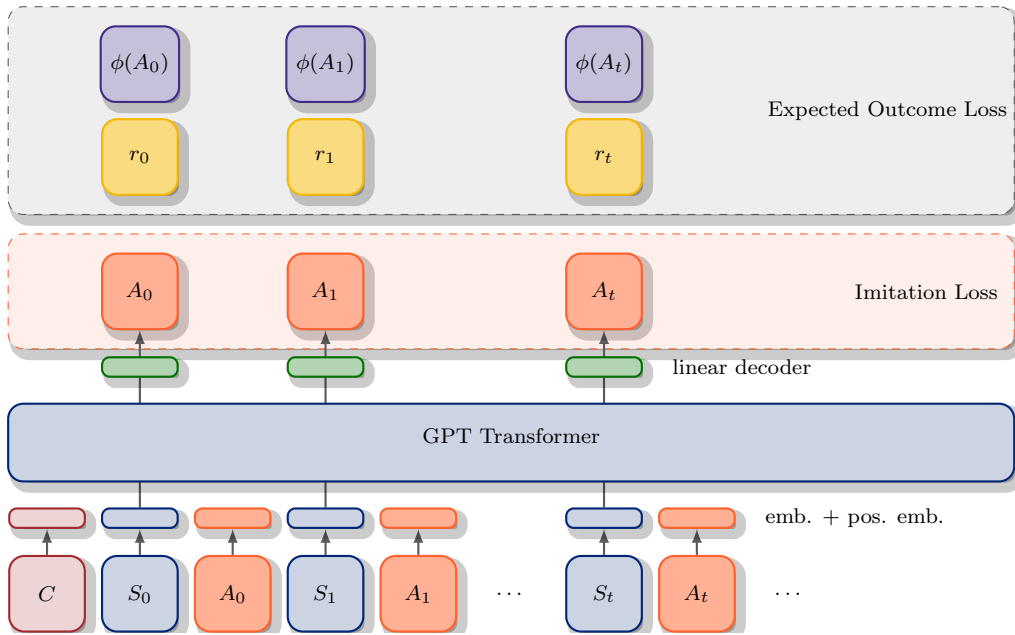


Figure 3.2: The bottom part shows the sequential input to the model. The middle part shows the DL architecture, which includes an embedding layer, a GPT transformer, and a linear decoder. The top part shows the parts of the models that are related to the loss functions.

threshold for any of the severity levels was excluded. 3) The probability threshold was set such that at least 95% of encounters will not be affected. This reduced the total number of unique action combinations from a theoretically possible  $2^{15}$  to 411.

To prevent throwing away information, each of the excluded ATs was replaced by a similar AT, from the remaining ATs pool. Such similarity was determined by choosing the common AT that contained the largest number of intersecting actions (greater than zero) with the excluded AT. In the case of a tie between several common ATs, the one that is more frequent in the data was chosen.

**Sequential representation:** IT uses sequential data where patients’ encounters are converted into the tuple  $(s_t, a_t, r_t)$  which records the patients’ current physiological state ( $s_t \in S$ ), the action token representing the actions that were taken during the encounter ( $a_t \in A$ ), and a medical outcome ( $r_t \in R$ ) often referred to in the reinforcement learning literature as the reward. To these, an additional token  $c \in C$  is added representing static information such as a patient’s age and sex. Together, all such tuples generate a medical trajectory ( $\tau = \{c, \{s_t, a_t, r_t\}_{t=0}^T\}$ ) per patient that can be used as input into a deep learning transformer architecture [JLL].

### 3.3.4 Architecture and Training

**Architecture:** IT is designed in a comparable way to that of the Decision Transformer [CLR] but with several modifications. Like DT: 1) IT has a transformer-based architecture that can generate actions autoregressively based on sequential contextual information. 2) IT uses a linear layer for each modality that projects the raw inputs to the embedding dimension. 3) IT learns an embedding for each timestep and adds it to the embeddings of the tokens in the timestep. 4) The tokens are then processed using a GPT (Generative Pre-trained Transformer) model that predicts future action. Unlike DT, IT does not receive rewards as part of its input. Instead, rewards are used to generate the expected outcome loss, which is described in the next section. Figure 3.2 shows an overview of the IT architecture.

**Training using expected outcome reward:** We started by stratifying a patient randomly into either one of the three datasets for training, validation, and testing. Following that, IT was trained by sampling mini batches of sequences covering up to two years of information (I.e., maximum 24 time steps – one for each month where we have an observation) from the training data. The probability score corresponding to the input token  $s_t$  was trained to predict the actions token  $a_t$  using the cross-entropy loss. These losses are then averaged across all time steps for all patients to generate the Imitation Loss. In addition to that, we used importance sampling theory to add an Expected Outcome Loss described by the following:

$$L_{\text{expected-outcome}} = \frac{\sum_{i=0}^N w_i R^{(i)}}{\sum_{i=0}^N w_i} \quad (3.1)$$

$$w_i = \prod_{t=0}^{T_i} \frac{\phi_{\theta}(a_{t,i} \mid s_{\leq t,i}, a_{< t,i}, c)}{\pi_o(a_{t,i}, s_{\leq t,i}, a_{< t,i})} \quad (3.2)$$

$$R^{(i)} = \sum_{t=0}^{T_i} \gamma^t r_{t,i} \quad (3.3)$$

Where  $\phi_{\theta}(a_{t,i} \mid s_{t-1,i}, a_{t-1,i})$  is the transformer probability score for the  $a_t$  token for patient  $i$  given additional context tokens, e.g., demographic variables, and  $\pi_o(a_t, s_t \mid s_{t-1,i}, a_{t-1,i})$  is the observed conditional probability for the  $a_t$  token for patient  $i$  – both conditioned on the state  $s_{t-1}$  and  $a_{t-1}$ .  $w_i$  is the importance weight for patient  $i$  calculated as the product of ratios between the model’s outcomes and observed probabilities and  $R^{(i)}$  is the total discounted reward for patient computed as the discounted sum of patient  $i$  rewards with discount factor  $\gamma$ . This loss encourages the transformer to recommend actions that yield high rewards. The observed conditional probability  $\pi_o(a_{t,i}, s_{t,i} \mid s_{t-1,i}, a_{t-1,i})$  can also be replaced with a more general term  $\psi(\theta')(a_{t,i})$  generated by a transformer that was trained similarly but using only the imitation loss to represent a more general probability score that is conditioned on the entire history up to time  $t$  than the conditional probability. The final loss used to train the IT model was a combination of the above losses:

$$L = L_{\text{imitation}} + \lambda L_{\text{expected-outcome}} \quad (3.4)$$

Where  $\lambda$  is a hyper-parameter used to determine how much the model is encouraged to optimize for future outcomes while deviating from observed actions that are captured by the imitation loss. To find a robust value for  $\lambda$ , we used bootstrapping to choose the  $\lambda$  values that produced the highest low confidence interval (2.5%) on the validation set. For other hyperparameters, transformer hyperparameters were chosen to have the same values as described in TT [JLL] and learning rate was chosen using the Learning Rate Finder implemented by the Pytorch Lightning package.

### 3.3.5 Evaluation

Off-policy estimation using weighted importance sampling: Importance sampling is a technique frequently used to estimate expectations when the data generating distribution is not directly accessible. It involves utilizing a surrogate distribution to approximate the desired expectation. It is often used for evaluating offline RL models [KCB] and a generalization of the inverse propensity weighting approach [RR] often employed in causal analysis [IR, Imb, IR, Rob]. In the case of off-policy estimation (OPE), we are interested in estimating a potential outcome across an entire population caused by following the suggested policy  $\pi^*$ . However, since we cannot apply policy  $\pi^*$  without a clinical trial, we can only use samples drawn from the behavioral policy  $\pi_b$  to estimate the effect of applying  $\pi^*$  to make clinical recommendations. In such a scenario, we can use importance sampling to estimate the average effect of following policy  $\pi^*$  using the formula:

$$\hat{V}_{\text{WIS}} = \mathbb{E}_{\pi^*} [R] \approx \frac{\sum_{i=1}^N W_i R^{(i)}}{\sum_{i=1}^N W_i} \quad (3.5)$$

$$W_i = \frac{\pi^*(\tau^{(i)})}{\pi_b(\tau^{(i)})} \quad (3.6)$$

### **3.3.5.1 Quality levels of the behavioural policy**

To estimate the expected reward, WIS (Weighted Importance Sampling) relies on the ratio between the suggested policy and the behavioral policy. However, given that the observed policy is merely an estimation, the sensitivity of WIS to the quality of the estimator becomes a pertinent concern. Therefore, to investigate the influence of the quality estimate of the behavioral policy on the WIS estimation, it is beneficial to estimate the behavioral policy at various levels of quality.

To achieve this, we first established the log-likelihood of actions as the metric for defining estimation quality where higher log-likelihood values for observed actions indicated higher-quality models. Next, to generate models of varying quality levels, we employed a transformer model trained to imitate the observed actions in our dataset. Using gradient descent, the model underwent iterative improvement, starting from random initialization and ending at a fully converged model where at each gradient step, the model became increasingly adept at imitating the clinical actions we had observed. Throughout the optimization process, we saved the weights of the trained model, yielding different models representing distinct stages of training and possessing diverse levels of quality. By estimating the behavioral policy at different quality levels, we sought to comprehend the influence of the quality estimate on WIS estimation.

### **3.3.6 Permutation Testing**

To gain deeper insights into the nature of our decision policy and its reliance on specific events in medical history versus the sequential nature of the data, we used a permutation test methodology applied to the input sequences. The following steps were undertaken to execute this methodology: 1) Random Shuffling: Initially, we randomly shuffled the order of patient encounters within the data. This step was crucial to disrupt any inherent temporal relationships present in the original sequence. 2) Suggested Policy Probability Score: After

shuffling the data, we applied our decision policy model, the Importance Transformer (IT), to the shuffled data. By doing so, we obtained the suggested policy probability score for each of the actions on the shuffled data. 3) Behavioral Policy Probability Score: Additionally, we utilized a transformer-based behavioral model to obtain the behavioral policy probability score for each action within the shuffled data. It is important to note that during this permutation analysis, the model used to compute the behavioral policy was trained without time embedding. This adjustment was made to account for the absence of temporal information in the shuffled data. Notably, this model yielded a superior likelihood score compared to the one trained on the original data. 4) Expected Reward Calculation: Using the suggested policy scores from step 2 and the behavioral policy scores from step 3, we computed the expected reward using Equation 1, which captures the essence of our reward estimation mechanism. 5) Repetition and Comparison: To ensure robustness and statistical significance, we repeated this process 100 times, each time performing a different shuffle of the data. Subsequently, we measured the number of times the expected reward obtained by using IT on the shuffled data surpassed that of using IT on the original data. This comparison gave us valuable insights into the impact of the permutation and temporal disruption on our decision policy’s effectiveness.

### **3.4 Discussion**

We presented promising results along several key clinical outcomes in utilizing offline reinforcement learning methods to augment clinical practice guidelines. Our principal hypothesis was to capture justified clinical variations from clinical practice guidelines to customize treatment for the unique circumstances of the patient and enhance CPGs such that more complex disease pathways are captured and are able deliver superior clinical outcomes. The AI employed showed empirically superior performance compared with average CPG outcomes seen in the data and current transformer-based state-of-the-art methods. Introducing this approach into practice may justify variations in care choices while optimizing outcomes. One

contributing factor to outcome variability is the time constraint on reviewing lengthy medical histories from multiple encounters, which can lead to decisions made from partial information containing recent encounters only. In addition, clinicians may rely on recent states of patients' health, whereas disease progression is often complex and spans multiple encounters. In this context, finding ways to effectively analyze and interpret long-term medical records is crucial to improving CPGs to capture complexities arising from poly-chronic disease pathways. We used several technical innovations to deal with the nature of clinical records. These include relative positional embeddings that encode the time between encounters, the sparsity of the encounters and the clinical outcome, and a loss function using weighted importance sampling.

However, there are several limitations. Importantly, our data source is based on claims, which are primarily intended for administrative purposes and can lag by a number of months due to the claims adjudication process. However, the codified nature of claims provides very concrete action spaces, which is important for reinforcement learning approaches. In the approach presented in this paper, we used a carefully curated dataset that had been annotated by clinical domain experts to provide us with details about disease severity. Further research will explore the ingestion of claims and electronic medical records into our transformer model. This can be further extended into pre-training the transformer architecture on large datasets and finetune on learning personalized clinical practice guidelines.

The sparsity of clinical outcomes also poses a modelling challenge. It is often not straightforward to make a qualitative assessment about a health state, so we expect innovations that deal with this in more principled ways. One recent approach called decision stacks separates the state, action spaces prediction with two models, which is able to account for different semantics in state space and action space [ZG]. When training decision stacks end-to-end, it can learn implicit rewards per step that may help in alleviating the sparse reward problem. In future work, we also would like to address the explicitly engineered action spaces accounting for combinations of treatments. One solution to this is to factor action spaces as proposed in [TMS].

With these future directions, we hope that the high-capacity nature of transformer models will be able to tease out even more nuance and hence provide improved clinical recommendations. Also, our solution does not necessarily generalize to other clinical outcomes. For example, we found that use-cases involving hospitalizations (not shown in this article due to brevity) as a clinical target produced results that showed a deterioration in performance.



## CHAPTER 4

# The heterogeneous effects of social support on the adoption of Facebook’s vaccine profile frames feature

### 4.1 Introduction

Widespread acceptance of COVID-19 vaccines is essential for achieving the coverage required for herd immunity, but in many countries, a sufficiently large proportion of people are still hesitant about receiving available vaccines [Laz21, Sol21]. Reaching sufficient vaccine coverage has been challenging due to barriers at multiple levels. One commonly used classification system describing these barriers is the 4Cs model which segments people based on the main driver for hesitancy: confidence (lack of trust in health institutions and pharmaceutical interventions), convenience (structural barriers preventing vaccination conversion despite intent), complacency (low perception of disease risk), and calculation (significant information searching) [BBC15]. Furthermore, as people make vaccine decisions, Social Contagion Theory suggests that social influence also plays a role as these considerations are influenced by belief in the decisions of others [CF13]. Recent studies have shown that such social influence can have a substantial effect on eventual vaccine decision-making, with positive associations found between acceptance and beliefs about the intentions of others to vaccinate [Bru13, Bre17, AEO21, KGK21, MCG21]. These associations are amplified when the others in question are close, trusted ties from a person’s social network [GGM20, Lau22, Rab22] In order to activate this social influence channel, people need to have accurate information about the beliefs of their social network, yet it’s currently unclear to what extent people are aware of

the vaccine decisions of others in their social network and may misestimate the degree of acceptance/resistance based on the amplification of a relatively small number of voices. For example, being exposed online to amplified messages of concerns regarding vaccine safety could decrease confidence and move calculations toward hesitancy [Loo21]. On the other hand, positive indications that trusted ties have chosen to vaccinate can combat this phenomenon and result in increased confidence and adjust factors such as calculation and complacency towards intent [KGK21]. In general, we do know that people underestimate others' adherence to a range of COVID-19 preventative behaviors [GAL21], biasing their perception of social norms towards non-compliance.

In order to make people more aware of the vaccine perceptions of their network connections, Facebook, in partnership with public health agencies, recently launched vaccination profile frames (VPFs) to enable users to surround their profile picture with a supportive message with respect to vaccination [Met21]. This form of advertising one's support for vaccination is the raw material that may allow social influence to make progress on the 4Cs. Previous work has established the impact of social proof-driven behavior change on Facebook, in non-health-related areas such as voting [Bon12], friending [ST20], and activism for social issues [SA15]. However, little is known about the factors that drive social signaling of vaccination support on social media, their relative importance, their overlap with factors that drive vaccine decision-making more broadly, and whether there are any negative downstream effects of sharing one's support. In this study, we explore these issues in the context of VPF usage on Facebook.

Our first research objective (RQ1) seeks the factors that promote VPF adoption, with a particular focus on determinants related to exposure to the adoption decisions of a user's friend network. The VPF feature rollout, coupled with our knowledge of the overall Facebook social graph and user demographics, provides the variation and controls which enables us to address RQ1 quantitatively along a number of key dimensions. Specifically, (1) promotions for VPFs for the frames took on several forms, including those with/without the social context of

friend adoptions, allowing us to observe the effects of social proof, (2) among promotions with social context, friends were selected at random to produce a mix of relationships, enabling a study of tie strength effect, and (3) VPF promotions were held back from a set of random users, giving us an interventional setting to validate our main findings. Together, these factors allow us to determine the effects and heterogeneity of social context on VPF adoption using both observational and experimental data.

Our second research question (RQ2) addresses the potential that while promoting vaccine beliefs may lead to improved public health outcomes in aggregate via social influence, expressing support for a polarizing issue such as COVID-19 vaccines may also come with social risks and unwanted negative social interactions [Oz18, Sch18]. To address RQ2, we searched for a detectable backfire effect against those who adopted VPFs, where we operationalized this effect to be negative actions received on Facebook that limit social ties with a VPF adopter (unfriending, unfollowing, blocking).

To our knowledge, this is the first large-scale quantitative look at how social context applies to people’s decision to socially demonstrate their choice to vaccinate, a distinct and less studied behavior compared to vaccine acceptance about which much more is known, giving attention to both positive (RQ1) and negative (RQ2) outcomes. Our results have implications for understanding the determinants of vaccine-related social signaling which is crucial for maximizing the impact of the social influence channel and for the design of messaging campaigns aiming to drive health-related behavior change via social media.

## **4.2 Results**

### **4.2.1 VPF adoption exhibits a pattern of complex diffusion**

Signaling one’s support for a polarizing issue such as vaccination comes with social risks, and may require the psychological support of first seeing friends adopt this behavior for many users to do so themselves. If so, there are downstream questions as to whether this form

of influence follows a pattern of complex (increasing dose-like effect of multiple exposures) versus simple (no increasing dose-like effect) diffusion. To determine this association, we used bivariate analysis and measured the empirical adoption probability of users in our sampled data conditioned on the number of (1) their friends that had previously adopted a VPF and (2) their friend's VPFs adoption posts they saw in their News Feed (a post alerting a user's adoption is automatically generated and shared in their friend's Feeds).

Figure 1a shows that the probability of adopting a VPF rapidly increases as the number of friends that have previously adopted grows, starting at a baseline of very close to 0% when no friends have adopted and saturating at about 2% when around 40 friends have done so. However, users do not necessarily see the adoption posts for all of their friends, as they may not scroll far enough on their Feed, scroll past these posts without viewing them, or simply not be on Facebook at that time. Therefore, we also looked at this same adoption probability conditioned on adoption post impressions in Feed that were reliably seen by the user. These results in Fig. 1b show a quicker saturation effect, at approximately 7 exposures leading to about a 4% adoption rate. Together, these results are supportive of social proof playing a role in VCP adoption, via a complex diffusion with saturation after 7 exposures on average, and requiring upwards of 40 friend adoptions for exposures to reach this level.

#### **4.2.2 Pre-existing openness to vaccines requires significantly less social proof for adoption**

Not all recipients of social proof are alike, and while Fig. 4.1 established a population-level association between social proof and VPF adoption, further segmentation of the data reveals significant adoption heterogeneity related to existing vaccination attitudes. Specifically, we divided users based on different noisy proxies for overall vaccination attitudes to evaluate the differential impact of social exposure as a function of rising openness towards vaccination. The two proxies we utilized were (1) profile county location thresholded to divide users between high/low COVID-19 vaccination rate counties (top and bottom 25 percent quartiles) and

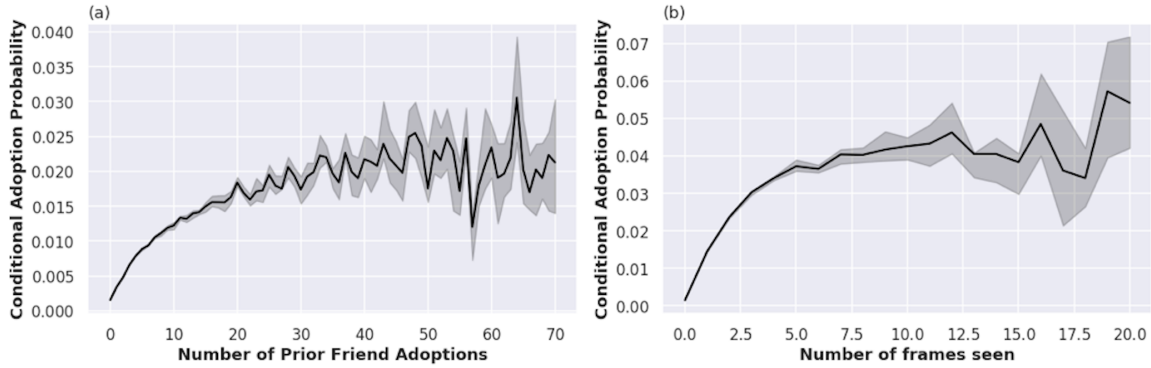


Figure 4.1: The probability of adopting a VPF is conditioned on a number of friends who have adopted and b the number of friend’s adoption posts seen. A pattern of complex diffusion is evident, in which as the number of social proof exposures increases, so does the likelihood of the user adopting the frame.

(2) the binned number of high-quality health pages followed by the user (see the “Methods” section for details on these pages) [Fre86, YF14]. In both cases, we observed that substantially less social proof is required to reach comparable adoption rates as we move up in the levels of these proxy variables.

For example, Fig. 4.2 (a) shows that as users follow more high-quality health pages from trusted health authorities, the effect of six exposures to VPF social proof increases the adoption rates by 26% (95% confidence intervals of  $\pm 19\%$ ) when comparing users who follow 10+ high-quality health pages versus those who follow none. Figure 4.2 (b) shows a 47% increase (95% confidence interval of  $\pm 56\%$ ) at 3 exposures when we utilize a location-based attitudinal proxy based on the specified home county of the user (note that at higher exposures the statistical significance of the difference disappears). Overall, these different vaccine attitude proxies highlight that substantially more social proof is required to drive comparable VPF adoption when there is existing resistance toward vaccination.

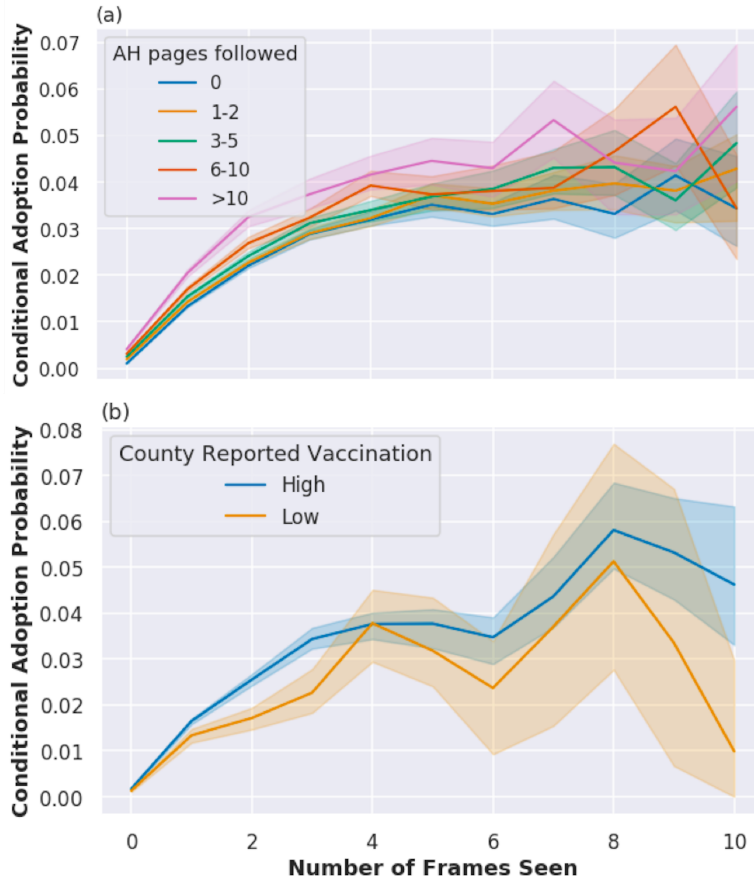


Figure 4.2: The probability of adopting a VPF, conditioned on the number of friend’s adoption posts seen, and segmented by a Authoritative health (AH) pages followed by the users, and b the COVID-19 vaccination rate in the user’s home county (binned by top and bottom 25 percent quartiles). These cuts provide proxies for pre-existing vaccine attitudes and show that significantly less social proof is required to reach comparable adoption rates as we move up in the levels (representing more openness to vaccination, in aggregate).

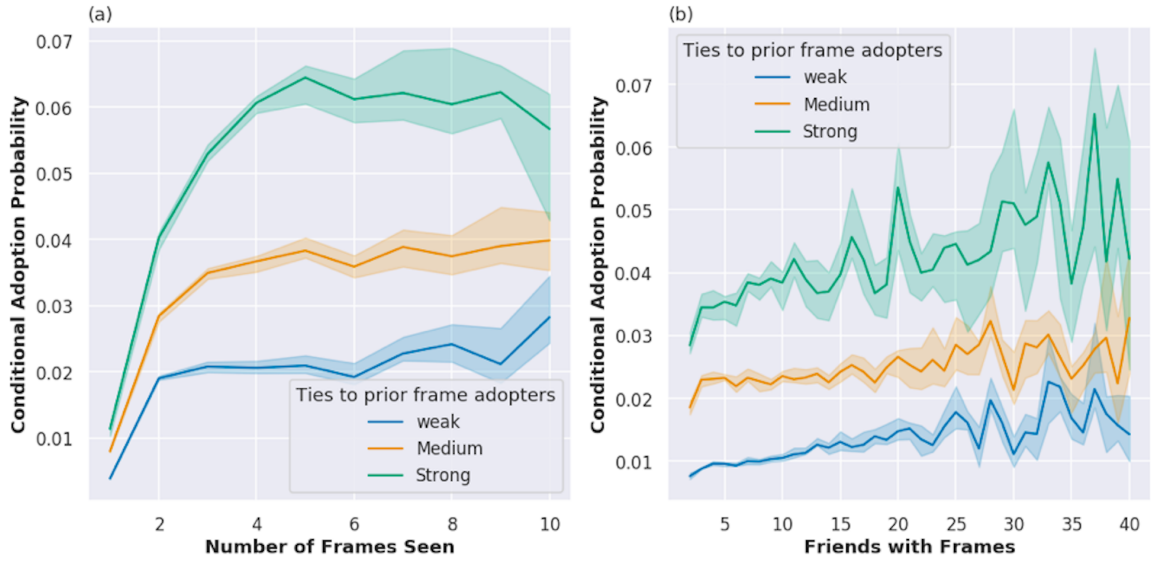


Figure 4.3: The probability of adopting a VPF segmented by the levels of tie strength with prior adopters and conditioned on the (a) number of friends who have adopted and (b) the number of friend’s adoption posts seen. Users with strong ties to prior adopters seem to be more likely to adopt the VPF when social proof exposure increases compared with users that have weaker ties with prior adopters.

#### 4.2.3 Social proof from stronger ties has a greater effect on adoption

Having shown heterogeneity in adoption response from the point of view of the recipient of the social proof, we next examined differential response when the friend providing the proof are close/far ties. Facebook users with frequent interactions on the platform, such as close friends and family, generally have a higher interpersonal influence on each other than user pairs who rarely interact or where the interactions are only one way (e.g. following a celebrity) [AW14]. Figure 4.3 shows that social proof from strong ties indeed leads to a stronger likelihood to adopt the VPF compared to weak ties, saturating at 5% versus 2% when approximately 40 friends have adopted, and at 6% versus 3% after being exposed to about 7 VPF posts.

#### 4.2.4 Influencers showed limited effect on adoption

While social proof from close ties was overall more influential, we also looked at a subset of weak ties that are of particular interest as campaign messengers, the social influencer. To do so, we examined the most followed 105 influencers in our data set and compared the adoption rates of VPFs among their followers before and after the influencers adopted the frame themselves. For a comparison control value of this difference, we matched each VPF-adopting influencer to 10 similar non-adopting ones, based on follower count and graph embeddings (see the “Methods” section for details), and looked at the changes in their follower’s adoption patterns across the same time period. Figure 4.4 (a) shows the difference in differences (DID) set up for our analysis for a particular adopting influencer and a matched control from our data. Figure 4.4 (b) shows the average DID between the influencers and their matches, as well as the permutation test-based P-values (see the “Methods” section for details). Only 5% of the 105 influencers produced a significant DID effect on the VPF adoption rates of their followers, suggesting that feature adoption was not primarily driven by social influencers, and sharpening the importance of strong ties.

#### 4.2.5 Modeling the effects of social proof on adoption

Having demonstrated significant effects and heterogeneity of social proof in isolated bivariate comparisons, we next moved to model VPF adoption as a function of these and other confounding variables in order to estimate the contributions of the different promotional formats. To do so, we implemented a logistic regression where VPF adoption is the dependent variable, exposures to the different promotional formats and whether one of these included a strong tie are the independent variables of interest, and the set of confounders included prior beliefs, general Facebook activity, number of friend VPF adoptions, and user age/gender/location.

The performance of the model as measured by the area under the receiver operator characteristic curve (AUROC) using cross-validation was 0.87 [95% CI 0.87,0.87], indicating that



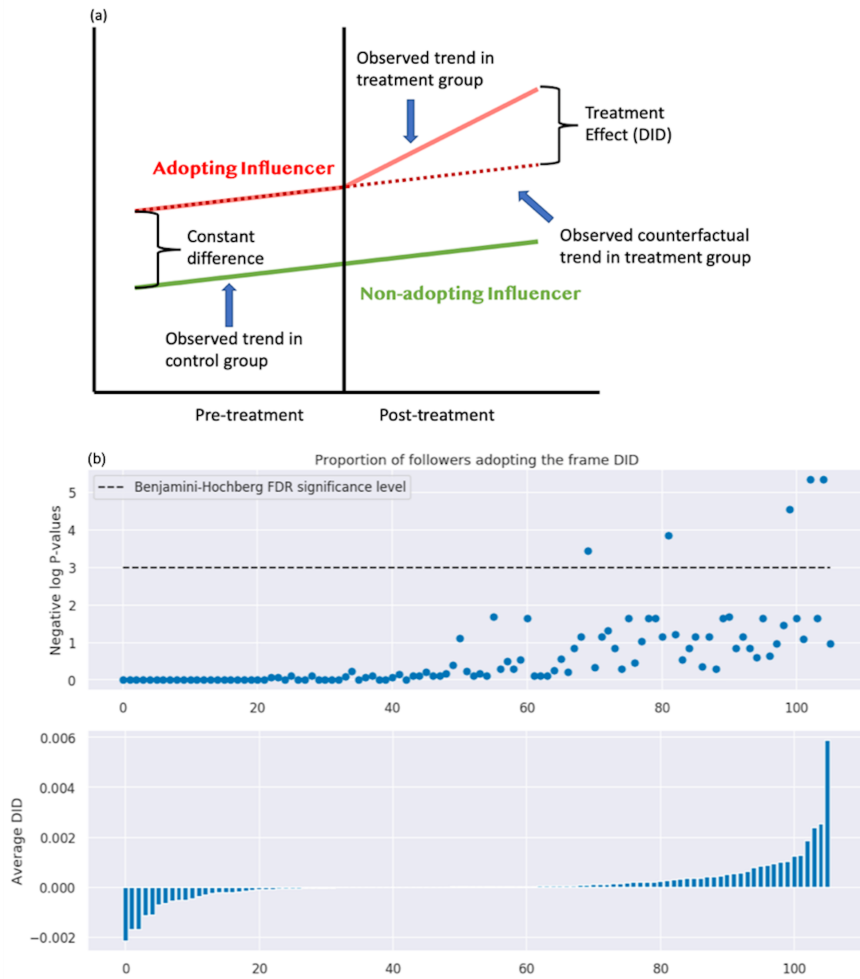


Figure 4.4: **(a)** A difference in differences (DID) approach was taken to estimate the effect of an influencer’s adoption decision on the decisions of their followers. In this illustrative example, we show an adopting influencer and a matched non-adopting control (in practice, we use 10 matched controls per influencer). We estimate the effect of the adopting influencer’s decision on her followers by looking at the departure from the counterfactual provided by the non-adopting influencer’s followers’ behavior. **(b)** A permutation method enables deriving an empirical null distribution of DID values per influencer, allowing determination of P-values (adjusted for multiple hypotheses testing), and revealing that only about 5% of influencers show a significant effect at  $\alpha=0.05$ .

Table 4.1: Regression Models Performance

Model	AUROC [95% confidence interval]			
Model with a State Vaccination Rate feature	0.866 [0.866,0.867]			
Model with a State indicator	0.870 [0.870,0.870]			
Model with Demographics only	0.760 [0.759,0.762]			
<b>Regression coefficients</b>	<b>coef</b>	<b>std err</b>	<b>z</b>	<b><math>P &gt;  z </math></b>
High friend count	-0.15	0.05	-2.92	3.50E-03
Influencer’s post seen	-0.22	0.34	-0.65	5.16E-01
Strong ties post seen	0.65	0.06	11.5	1.25E-30
friend agg only	1.47	0.08	18.23	2.94E-74
Frame post seen only	2.51	0.05	46.27	7.14E-293
Frame post and friend agg	2.56	0.06	40.68	7.14E-293
Profile Prompt QP	0.69	0.03	21.77	4.53E-105
Newsfeed QP	0.16	0.03	4.96	6.92E-07
Intercept	-5.17	0.13	-40.22	7.14E-293

the logistic model is a good choice to describe the relationship between the outcome and the dependent variables (Table 4.1). Figure 4.5 shows that discovery by social means has a significantly stronger effect when compared with the non-social promotion that appears on a user’s profile page (OR=6.18 and 95% CI of [5.46,6.88] for profile frame post; OR=2.18[1.83,2.54] for friend aggregation post). The effect is even greater when comparing social discovery to non-social promotions appearing on a user’s Newsfeed (OR=10.50[9.10,11.77] for profile frame post; OR=3.71[3.07,4.32] for friend aggregation post). The significant coefficient scores (log odds scores) for social discovery and the high OR values compared to non-social discovery highlight and quantify the value of providing social proof to drive VPF adoption.

In addition, when social proof from a strong tie (close friend or family) is provided in the social discovery, we estimate a further increase in the adoption odds score of 1.9 (holding

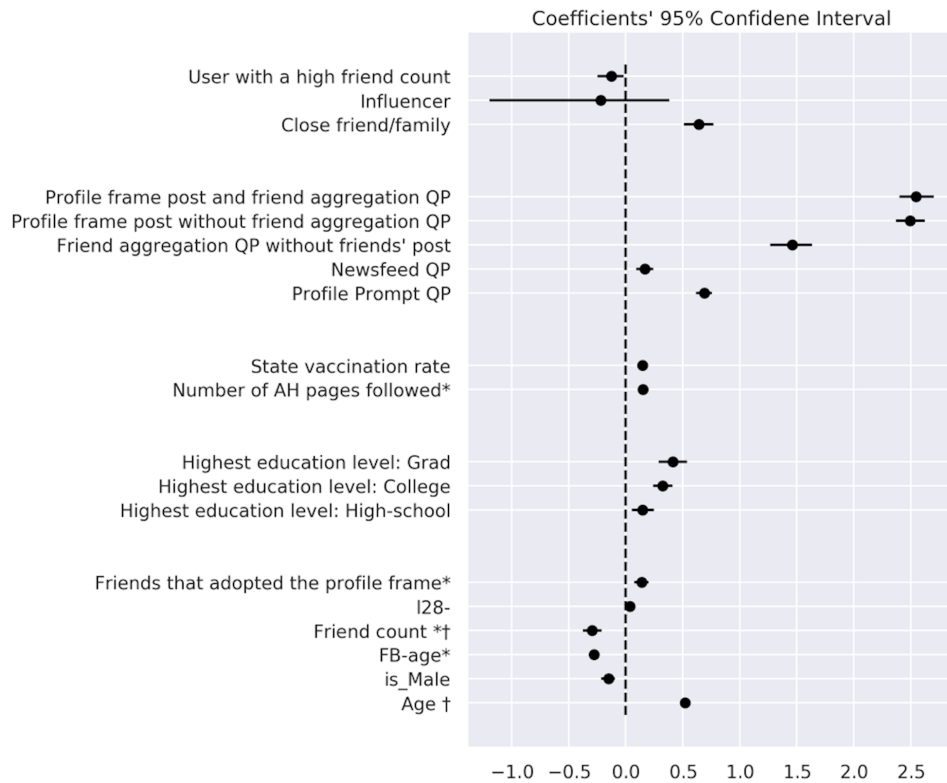


Figure 4.5: Logistic regression coefficients for adoption conditioned on the discovery channel (different QPs=“quick promotions”), tie strength of any included social proof, and a broad set of confounders, including those where heterogeneity in adoption response was observed. The coefficients (log odds scores) imply that VPF adoption is strongly affected by social aspects such as seeing a vaccination post from a close friend or seeing a promotion informing users that their friends have adopted the frame. FB-age is defined as the number of days since the user has signed up to the Facebook platform. “l28-” is an inactivity variable defined as the number of days within the last four weeks at which the user has not been active on the FB platform. (\*) marks a log transformation. (†) marks a standardization transformation for zero mean and unit variance.

everything else constant). Consistent with our findings of influencer effects in the previous section, we also found no significant effect when an influencer was included in the social proof. These results confirm the value of close ties in driving VPF adoption in a controlled model setting.

#### **4.2.6 A randomized field experiment provides causal support for social proof and tie-strength effects on the adoption**

While our modeling results point to strong social proof effects that are amplified when the source is a close tie, these estimates are observational without causal interpretation given that we cannot rule out the existence of uncontrolled/unobserved confounders. To provide some causal support for these main conclusions, we utilized an experiment that held out a random set of eligible US Facebook users from the social aggregation promotion. This control group did not receive the Newsfeed promotion shown in Fig. S2B, while the test group did so. Since eligible users (active 18+ Facebook users from the US with at least three friends that had previously adopted the VPF) were assigned to the test and control group completely at random, the conditional ignorability assumption holds and the causal effect of the VPF on eligible users can be estimated [HR]. The covariate balance across control and test groups is shown in Table 2, and the average treatment effect (ATE) on the treated for this promotion was a 0.15% (95% CI=[0.12%,0.18%]) increase in VPF adoption rates (Fig. 4.6, left column) which amount to a relative increase of 75% (95% CI=[56%,79%]) in VPF adoption. In cases where the adoption rate among the control is very small, the absolute effect size can also appear quite small, and so it is important to also consider the relative effect size which indicates the proportional increase caused by the treatment. Here, it is estimated that among eligible users the VPF would increase adoption of the frame by 75%. In addition, when considering the millions of users active on the Facebook platform in the US, an absolute effect size of 0.15% would result in a large increase in people who adopt the VPF.

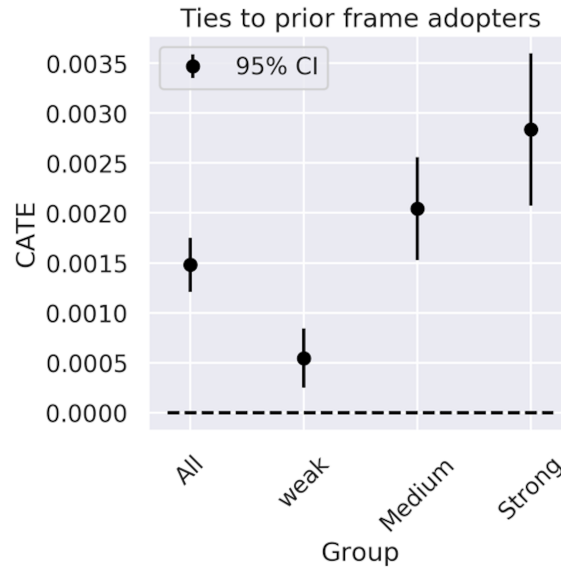


Figure 4.6: Users in the treatment arm received the friend aggregation post as a means of social proof for VPFs. The three friends for this format were selected at random, enabling estimation of conditional average treatment effects conditioned on approximated tie strengths to the friends in the aggregation. The findings show an increasing trend in CATE correlated with increasing levels of tie strength.

In addition to ATE, we also attempted to estimate the conditional ATE (CATE) on the treated given an idiosyncrasy of this experiment: friends who adopted were selected for the social aggregation post at random, but the friend’s identities were not retained in our logging. While the random selection created the variation to estimate a CATE of tie strength on adoption, the lack of friend identities led us to use each exposed users’ maximum tie strength across all their friends who previously adopted for conditioning levels of tie strength (weak, moderate, strong). The rationale for choosing this conditioning scheme is that users with higher tie strengths on average will tend to see more social promotions from closer friends (in the aggregate). Figure 4.6 shows an increasing trend in CATE correlated with increasing levels of tie strength, with the strongest level significantly higher than the weakest one (CATE strong=0.28%, CATE weak=0.05%, difference=0.23% (two-sided P-value=8.4e-9).

#### 4.2.7 Causal machine learning reveals additional heterogeneous treatment effects

Estimated exposure to different tie strength levels of social proof showed significant CATE differences in our pre-planned experiment. To search for other potential heterogeneous treatment effects (HTE), we applied the causal forest algorithm [ATW19] to the experimental data and our full set of covariates from the modeling section, allowing us to rank covariates by their contribution to heterogeneity in adoption upon treatment with the social aggregation promotion (Fig. 4.7). This analysis confirmed that tie strength is a major HTE contributor, showing up second in the ordered feature importance scores. Lower in the list, we also observe adopter friend count, health pages followed, and state vaccination level as drivers of heterogeneity, showing that the data-driven HTE discovery approach confirms the observations in Figs. 4.1 and 4.2, where we see that these features drive heterogeneity, but to a lesser extent than tie strength.

The strongest HTE-driving feature was the age of the user, with older users ( $> 50$  years old) showing higher CATE (Fig. 4.8(a)). One factor that may be contributing to this age effect is tie strength. When looking at the proportion of friends who have adopted the VPF, segmented by levels of tie strength, we see that older users tend to have proportionally stronger tie friends (Fig. 4.8(b)). In particular, we see that once an age cohort has a proportion of close ties around 0.25 (Fig. 4.8 (c)), the CATE effect becomes significantly and consistently different from 0. This is approximately the point at which it becomes probabilistically more likely than not to select at least one strong tie when choosing three friends at random for the aggregation post. There may be other explanations for the association between age and adoption, including differential risk perception and incentives to promote vaccination among older users, but these are outside the scope of our study.

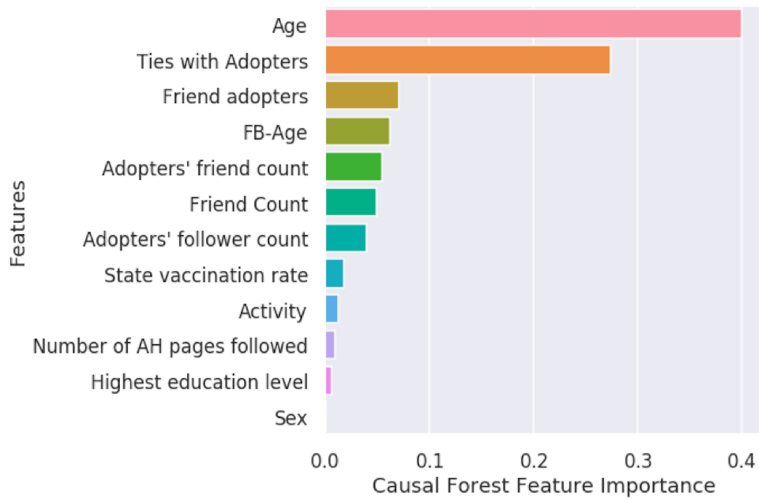


Figure 4.7: Each bar represents the importance of the associated feature in maximizing the heterogeneous treatment effect.

#### 4.2.8 Backfire effects of VPF adoption

Having established a positive relationship between VPF adoption and exposure to a friend’s VPF, we next moved to examine the potential negative side of exposing one’s vaccine beliefs openly. Specifically, we examined whether there are significant differences in targeted negative actions (unfollowing, unfriending, and blocking) upon VPF adoption for individual adopters in the 2 weeks before/after adoption.

Our findings revealed extremely low effects of VPF adoption on the number of negative actions received by adopters, suggesting that adoption did not have any significant backfire effects on individual adopters. The average treatment effect (ATE) on the treated was a 0.86 units increase in the relative difference in negative actions (Fig. 4.9 (a)). This finding was robust to stratification of users by a propensity to adopt (see the “Methods” section for stratification details), with an average effect size across the strata having a Cohen’s  $d=0.06$  ( $< 0.2$  suggests small differences between the two distributions; Fig. 4.9 (b-f)).

In addition, across all strata, we do not observe any temporal variation in negative action trends after VPF adoption. Specifically, we observe a statistically insignificant peak (95%

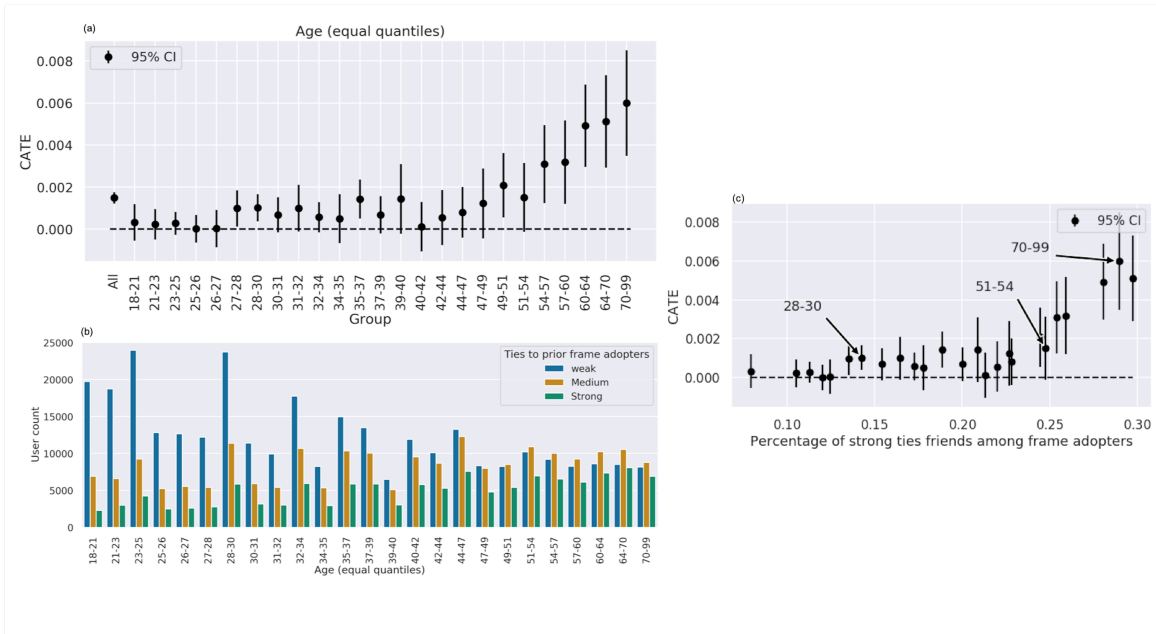


Figure 4.8: Heterogeneous treatment effects grow with user age (a). This pattern is likely driven by strong ties, as older cohorts tend to have proportionally more stronger tie friends who have adopted (b). As the age cohort's strong tie friend proportion exceeds 0.25, we see increasing CATE, significantly different from 0 (c) text annotation shows select cohort age ranges).



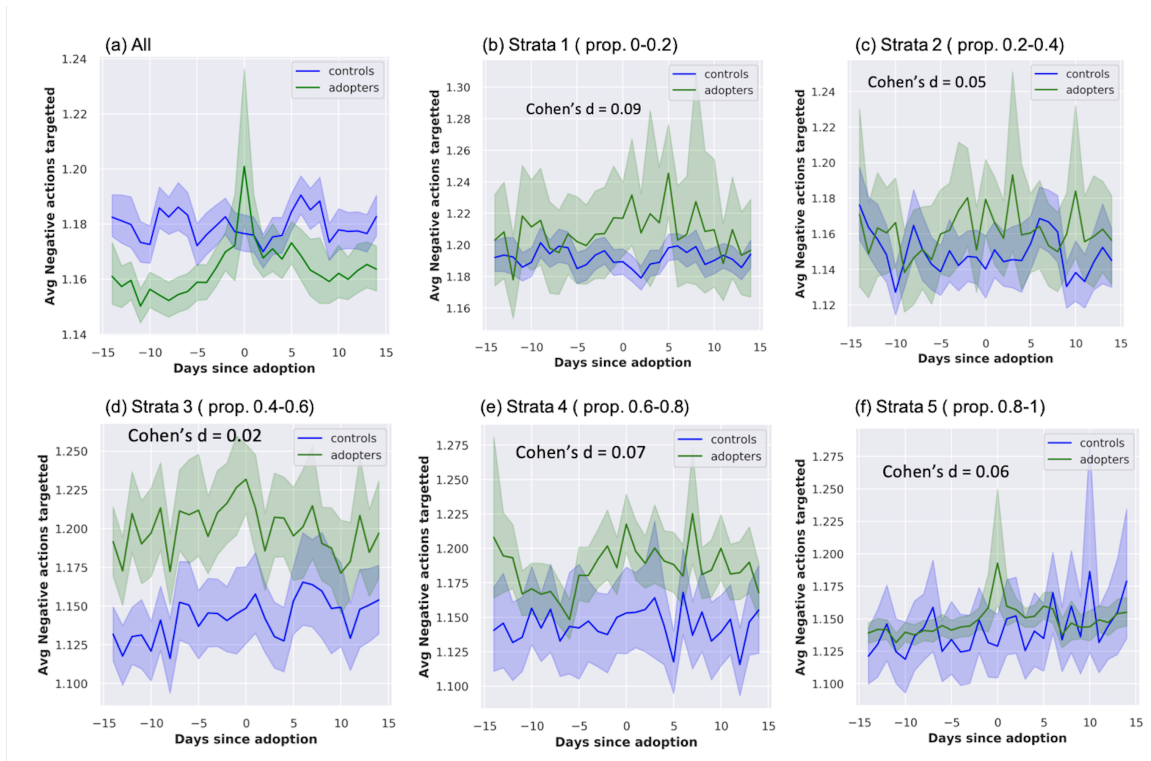


Figure 4.9: (a) Shows the average negative actions per day across all adopters and controls in the sample. (b-f) Show the average negative actions per day, across adopters and controls in each stratum matched based on the propensity score to receive treatment (VPF adoption).

CI [1.176, 1.233]) in negative actions on the day of adoption, however, this flattens to the baseline value prior to adoption immediately.

## 4.3 Methods

### 4.3.1 VPF adoption and promotional exposure data (RQ1)

This study was conducted using de-identified data logged by Facebook in the normal usage and launch of VPFs in accordance with Facebook’s data use policy. The full dataset contained  $\sim 1$  million users in the US who adopted a VPF (among a set of  $\sim 40$  official VPFs available for the initial feature launch) within the analysis window starting 2021-04-25 and ending

2021-05-08. We also selected a large randomly sampled set of  $\sim 10.5$  million US non-adopters who were active on FB and at least 18 years of age. All analyses to measure diffusion effects and modeling to uncover adoption drivers were conducted using samples drawn from these parent sets of users while preserving the adoption rate we observed in purely random samples from 2021-04-25 (0.64% adopters). A flowchart illustrating user selection is shown in the supplementary section.

To serve all downstream analyses with these user samples, we also collected standard demographic and Facebook activity controls, exposure counts for different VPF promotions (2 with social context, 2 without), VPF adoption dates, the number of accredited health pages they followed at the start of the analysis window (a proxy for prior vaccine beliefs), how many of their friends adopted VPFs prior to the user’s adoption, and available tie strengths for any social context displayed in promotions (see subsection “Additional data section” in the “Methods” section for details on health pages and tie strengths). Table 4.2 lists descriptive statistics for the data.

Among the promotions considered in this study are 2 non-social variants (a message in the user’s feed or profile page to adopt VPFs), and 2 social versions (messages in the feed that show that a single friend or a set of 3 friends have adopted VPFs). Examples of these promotional messages are shown in the Supplementary.

### **4.3.2 Influencers’ matching and difference in differences (RQ1)**

To examine the effects of influencer adoption on the adoption changes of followers, we chose the most followed 105 Facebook pages in the US that had adopted the VPF between 04/01/2021 and 06/24/2021. For each influencer, we measured the percentage of followers that adopted the frame one week before and after the influencer’s adoption date and then calculated the difference between the two measurements. As a comparison counterfactual value of this difference, we matched each VPF-adopting influencer to 10 similar non-adopting ones, based on follower count and pre-trained Facebook graph embeddings (minimum cosine distance to

Table 4.2: Descriptive statistics of features used in the regression analysis.

<b>N=10,822,480</b>					
<b>Adopted the VPF:</b>					
Yes=69,508					
No=10,752,972					
<b>Numerical variables</b>	<b>Mean</b>	<b>Std</b>	<b>Min</b>	<b>Median</b>	<b>Max</b>
Age (years)	43.24	16.69	18	41	99
FB age (days)	2988.96	1700.62	1	3582	6333
User's friend count	467.23	693.47	0	243	4987
Number of friends that adopted the VPFs	4.82	9.69	0	2	797
VPF posts seen	0.5	1.23	0	0	104
User's "State vaccination rate" (%)	53.04	8.01	35.72	52.85	73.34
Number of days that contained Facebook activity within the last four weeks	21.51	9.49	1	28	28
<b>Binary variables</b>	<b>Prevalence (%)</b>				
Newsfeed promotion seen	19.85				
Profile promotion seen	14.22				
Friend aggregation promotion seen	17.14				
VPF post seen	25.66				
VPF post from close friend seen	0.74				
VPF post from influencer seen	0.06				
VPF adoptions	0.64				
VPF post from user with high friend count seen	2.23				

<b>Categorical variables</b>	<b>Prevalence (%)</b>	
<b>Highest education level</b>		
Graduate school	4.79	
College	38.76	
High school	20.69	
Unknown	35.76	
<b>Sex</b>		
Male	53.73	
Female	46.27	

neighboring pages with similar follower count; see subsection “Additional data section” in the “Methods” section for details on embeddings), and looked at the difference in differences (DID) of their follower adoption patterns. For each VPF-adopting influencer, we measured the average DID with their 10 matches as a measurement of influence on the adoption rate of their followers.

To measure the statistical significance of the DID values, we generated P-values using a permutation method to approximate the null distribution. Specifically, we randomly permuted each influencer’s 2-week follower adoption rate vector, breaking up any temporal effect that was driven by the influencer adopting a VPF. Therefore, DID calculations based on 10,000 iterations of such permuted vectors captured the null distribution, which we used to assign P-values to our observed DID values.

### 4.3.3 Logistic regression model (RQ1)

We implemented logistic regression where VPF adoption was the dependent variable and exposures to the different promotional formats and whether one of these included a strong tie

were the independent variables of interest. We also utilized a set of confounders as described in the text. A full list of variables can be found in Table 4.2.

To estimate model parameters in a robust manner, we used a bootstrapping procedure where we: (1) randomly sampled 1 million users from among VPF adopters and non-adopters, maintaining the adopter ratio observed in purely random samples from 2021-04-25 (0.65% adopters; Supplementary Fig. A.1); (2) fitted a logistic regression model using the statsmodel package in Python [SP10] to produce maximum likelihood estimates of model coefficients; (3) measured the model’s performance using the area under the receiver operating characteristic curve (AUROC) on a randomly held-out sample; (4) repeated steps 1, 2, and 3 for 1000 iterations to produce a mean estimate along with a 95% confidence interval of model parameters and AUROCs.

#### **4.3.4 Randomized field experiment (RQ1)**

To determine the causal effects of social promotions, interventional data between 06/18/2021 to 07/18/2021 were collected. The experimental design was a simple A/B test where treatment was defined to be delivery of the friend aggregation post (Supplementary Fig. A.3), and the control condition was not receiving this promotion. Eligibility for inclusion in the experiment was based on being a non-adopter at the start of the experiment, age ( $\geq 18$  years old), location (US-based user), not having received a friend aggregation promotion within 2 weeks of the start date, and having at least three friends who had already adopted the VPF. Approximately 645K users met this eligibility condition, with roughly 323K randomly chosen for treatment and 321K as controls. The experiment outcome was the adoption of a VPF before the experiment’s end date.

Table 4.3 lists descriptive statistics for the experimental data, showing the covariate balance between treatment conditions.

Table 4.3: Descriptive statistics of features used in the randomized field experiment analysis.

<b>Experimental data</b>	<b>Cases</b>		<b>Controls</b>	
<b>N=644,231</b>	<b>N=321,438</b>		<b>N=322,793</b>	
<b>Numerical variables</b>	<b>Mean</b>	<b>Std</b>	<b>Mean</b>	<b>Std</b>
Age (years)	40.72	14.74	40.73	14.71
FB age (days)	3642.04	1425.84	3640.85	1428.11
User's friend count	810.91	956.64	810.32	955.72
Number of friends that adopted the VPFs	9.51	13.69	9.47	13.68
AH pages followed	2.69	8.82	2.68	9.52
User's "State vaccination rate" (%)	52.37	8.03	52.41	8.03
Number of days that contained Facebook activity within the last four weeks	26.81	3.55	26.82	3.54
<b>Binary variables</b>	<b>Prevalence (%)</b>			
Non-social promotion seen	41.4		42.2	
Friend aggregation promotion seen	31.2		0	
VPF post seen	24.9		24.8	
VPF adoptions	0.37		0.22	
<b>Categorical variables</b>	<b>Prevalence (%)</b>			
<b><i>Highest education level</i></b>				
Graduate school	5.3		5.2	
College	50.2		50.1	
High school	24.3		24.5	
Unknown	20.2		20.2	

Categorical variables	Prevalence (%)	
<i>Sex</i>		
Male	59.81	59.68
Female	40.19	40.32
<i>Ties with prior friend adopters</i>		
Weak	48.1	48.3
Medium	32.6	32.5
Strong	19.3	19.2

#### 4.3.5 Causal Forest for heterogeneous treatment effects (RQ1)

To search for heterogeneous treatment effects in the field experiment, we used causal forests which leverage the random forest algorithm to find sub-groups on which the conditional average treatment effect is maximized. We utilized the causalforest package in R for model fitting and feature importance metrics from the trained model to score the contribution to an effect size of each included covariate, which encompassed the same set as the predictive model described above. In general, covariates that were used more often and in earlier stages of tree building are provided with higher importance scores.

#### 4.3.6 Backfire effects (RQ2)

To examine the association between VPF adoption and backfire effects, we collected data on select negative actions taken against adopters by other users in the two weeks preceding and succeeding VPF adoption. Specifically, we selected actions that sever the relationship or inhibit information flow between a user pair: unfriending, blocking, and unfollowing. We included 475K users in our study who adopted a vaccine profile frame between 2021-04-14 and 2021-04-18 and received at least one negative action against them in the 2 weeks surrounding

VPF adoption. As controls, we sampled 507K non-adopters (from the same time period) identified by the criteria described above for RQ1.

We conducted a stratified propensity score analysis between adopters and control users who did not adopt the VPF. As covariates, we included demographic variables (age, gender), friend count, account tenure, number of vaccine profile frames seen, and accredited health (AH) pages followed by the user. To identify treatment (VPF adopters) and control users who are statistically similar to one another along the covariates, we match individuals with similar propensity scores into strata. Each stratum, then, consists of matched treatment and control users and lets us estimate the effect of VPF adoption on backfire effects within each stratum. To compute the propensity scores, we built a logistic regression model (accuracy=0.87, F1 score=0.84) with the above covariates to predict one’s likelihood to receive the treatment (VPF adoption). Then, based on the empirical distribution of propensity scores, our stratified matching approach groups treatment and control users with similar propensity scores into 5 strata. Table 4.4 shows the balance in covariates per strata. Lastly, we compute the average treatment effect per stratum with the outcome as a relative difference in negative actions targeted before and after adoption. Weighing the average treatment effect per stratum with the number of treated users in that strata gives the final average treatment effect on the treated.

### **4.3.7 Additional data**

#### **4.3.7.1 Tie strength**

We used an internal scoring of edges in the Facebook friend graph which gives higher weights to pairs of users who interact more often and more directly. These scores were clustered into five non-overlapping intervals representing tie strength buckets. Scores in the highest bucket were annotated as strong ties, those in the next two were medium, and the ties in the lowest two buckets were considered weak. In general, this annotation scheme tends to place close



Table 4.4: Descriptive statistics of features used in the backfire effect analysis.

Covariates	Adopters		Controls	
	<b>Strata 1: 29,307</b> <b>Strata 2: 37,854</b> <b>Strata 3: 37,191</b> <b>Strata 4: 65,034</b> <b>Strata 5: 305,264</b>		<b>Strata 1: 344,205</b> <b>Strata 2: 93,479</b> <b>Strata 3: 33,172</b> <b>Strata 4: 18,761</b> <b>Strata 5: 18,015</b>	
	Mean (per strata)	Std (per strata)	Mean (per strata)	Std (per strata)
<b>Age (years)</b>	35.66	8.16	30.76	7.9
	53.97	9.41	54.12	8.02
	38.09	17.4	48.49	21.36
	45.67	8.06	47.37	12.29
	53.52	13.86	51.17	14.26
<b>FB-age (log)</b>	3.29	0.55	3.29	0.51
	3.27	0.59	3.25	0.61
	3.34	0.51	3.33	0.54
	3.33	0.55	3.39	0.48
	3.43	0.45	3.45	0.42
<b>User's friend count</b>	1422.46	1406.76	1618.5	1408.17
	1169.54	1296.8	1052	1228.65
	1471.66	1427.36	1201.73	1286.96
	1438.88	1417.13	1312.49	1268.47
	1229.19	1312.69	1429.05	1341.87

Table 4.4: Descriptive statistics of features used in the backfire effect analysis.

Covariates	Adopters		Controls	
	<b>Strata 1: 29,307</b>		<b>Strata 1: 344,205</b>	
	<b>Strata 2: 37,854</b>		<b>Strata 2: 93,479</b>	
	<b>Strata 3: 37,191</b>		<b>Strata 3: 33,172</b>	
	<b>Strata 4: 65,034</b>		<b>Strata 4: 18,761</b>	
	<b>Strata 5: 305,264</b>		<b>Strata 5: 18,015</b>	
	Mean	Std	Mean	Std
	(per strata)	(per strata)	(per strata)	(per strata)
<b>Number of VPFs seen</b>	0	0	0	0
	0.03	0.18	0.01	0.09
	0.77	0.42	0.52	0.5
	0.98	0.14	0.9	0.29
	3.25	2.96	1.92	1.3
<b>AH pages followed</b>	4.68	6.99	2.64	4.79
	7.91	12.02	6.26	11.45
	6.52	12.6	7.01	16.04
	7.95	13.34	9.38	19.97
	13.47	28.94	16.18	50.37

friends and family into the strong tie bucket.

#### 4.3.7.2 Page embeddings

We utilized internal pre-trained Facebook page-dense embeddings which embed nodes in the page-page graph where edges are determined by followers/fans, links posted, topics discussed,

and other features (Slide 37 from [Fac]). These vectors place nodes such that a low cosine distance from a query node gives pages that are most similar. In terms of our application to influencers, they return similar celebrities.

#### **4.3.7.3 Accredited Health pages**

Facebook has annotated pages from global and US health organizations such as the WHO, UNICEF, and CDC, as well as local/regional sources of trusted health information to be disseminated on various platform surfaces [Met20]. We leveraged this list of pages and utilized their follows as a proxy for pre-existing vaccine beliefs, the assumption being that users with strong negative views towards vaccination will likely not follow such pages.

#### **4.3.7.4 Vaccine data**

We downloaded COVID-19 vaccination rate data at county granularity from the CDC website (<https://data.cdc.gov/Vaccinations/COVID-19-Vaccinations-in-the-United-States-County/8xkx-amqh>) on 06-01-21, which represents aggregate county vaccination coverage up until that date. For Fig. 4.2, we used the county-level metrics for the percent covered with 1+ dose among the 18+ population.

## **4.4 Discussion**

In this study, we show that social influence plays a significant role in increasing VPF adoption, an example of a health behavior change where users choose to advertise their support for vaccination to their social network (RQ1). We found evidence that adoption follows a complex diffusion process, where multiple instances of social proof increase the probability of adoption (Fig. 4.1), and that there is significant heterogeneity in this response associated with factors such as prior vaccine beliefs (Fig. 4.2), whether users become aware of the feature in a social

context (Fig. 4.6), and tie strength when social context is provided (Fig. 4.3). In short, significantly more exposures are needed to achieve comparable adoption levels when the user holds more resistant views to vaccination, or when the exposures lack social context from strong ties. We also jointly modeled adoption using these factors, controlling for a variety of confounders, to arrive at estimates of relative contribution among the factors, revealing that social support from strong ties is the most influential factor in driving adoption (Fig. 4.5).

This observational result was validated using a field experiment where a promotional message presenting multiple friends who had already adopted was held back from a control group, confirming the value of strong ties in an interventional setting, and therefore giving this relationship a stronger interpretation than simple association (Figs. 4.6 and 4.7). This randomization also provides evidence that this effect is not due exclusively to homophily, which we expect to be comparable between treatment and control groups, but representative of social influence.

These results, for the distinct and much less understood behavior of advertising one’s vaccine support, are consistent with the literature on factors that contribute to vaccine acceptance and add to the research in a number of ways. As in previous studies [Bru13, Bre17, Bru19, AEO21, KGK21, MCG21], we find that social influence is a strong determinant that drives a vaccine-related decision. While some studies have shown demographics to be more important [Bru19], we found that social influence is the strongest determinant. Also consistent with previous studies [GGM20, LL21, Rab22], we found strong ties to be the most influential form of social influence and age to be the strongest demographic determinant.

With respect to weak ties, many studies have reinforced the notion of the “strength of weak ties” [Gra] in diverse areas such as job searches [Raj22], scientific publications [FMF22], novel information propagation [BRM12], and many others. For vaccine acceptance, the results are mixed, with some studies showing that weak ties do matter [MCG21, Rab22], and others showing that they do not [SA23]. Our results support the latter conclusion, with very little effect found from the exposure to VPF adoption by weak ties. In addition, the value

of influencers, which are generally weaker ties with high follower count, has not yet been established although many vaccine messaging campaigns utilize such celebrities [Ive21, Lor21]. Our results show that these users are not an influential choice for providing VPF social proof (Fig. 4.4). One reason for these findings could be that for socially sensitive topics such as vaccines, the deeper affinity that a user has with their strong ties is a necessary precondition for being influenced to publicly disclose one’s views. While weak ties and influencers may still hold value in providing novel social capital to influence downstream decision-making [KSE21], they by themselves do not seem to trigger the advertising of beliefs publicly for this socially sensitive issue.

Finally, our study is also novel in our ability to estimate the “dose effect” of social influence on a vaccine-related decision in our findings of complex diffusion dynamics. These dynamics have implications for public health messaging campaigns, motivating designs that plan for multiple exposures per user to saturate conversion rates. When the campaign aims to make in-roads with those lacking vaccine confidence, our results showing the heterogeneity that comes with prior beliefs suggest that far more exposures will be needed to reach comparable conversion rates. If the campaign is budget constrained and cannot reach such high levels of exposure, it may instead be a better use of resources to go after one of the other factors from the 4C model.

On the opposite side of positive behavior change (influencing frame adoption), it was possible that exposing one’s vaccination beliefs via a VPF could also lead to unsolicited, negative reactions on Facebook. Despite vaccinations being a polarizing issue in the United States, we found no evidence of a backfire effect in which users exposed to their friend’s adoption responded by limiting social ties (RQ2, Fig. 4.9). While there can still exist other forms of such an effect, the fact that we did not observe increases in aggregate unfriending, unfollowing, and blocking suggests that campaigns in which identifiable social proof of vaccination is provided may not need to be overly concerned about large scale observable social contraction as an unintended downstream effect. Given the fact that VPF adoption

awareness came largely from social promotions or friend posts on Facebook’s News Feed, the absence of backfire effects cannot be attributed simply to the natural consequence of homophily or echo chambers as the Facebook friend graph has previously been shown to be cross-cutting with respect to user beliefs and interests [GMW10, LL21].

This study has a number of limitations. First, we do not know to what extent these findings generalize to other health behaviors beyond VPF adoption. While it’s unlikely that our findings idiosyncratically only apply to profile frames, and more likely that the learnings transfer to a variety of other health communications aimed at behavior change, it’s unclear where the boundaries of generalization are. To define these, further work which varies the behavior, exposures, and tie strengths is needed.

Another limitation of our work is that we are not generally able to distinguish to what extent the mechanism of adoption is driven by homophily versus social influence (outside of having the field experiment which allows us to control for homophily via randomization in one promotional format). This differentiation has significant implications where, if homophily is the dominant driver, such campaigns are largely converting people who already hold open views towards vaccination and have come together to form ties on social media, but not necessarily making gains to influence those lacking confidence and in areas of the network which may be in most need for behavior change to drive public health objectives such as minimum thresholds for herd immunity.

The data available for this study also presented some additional limitations. We only looked at controlled exposures from a fixed set of promotional formats for which data was cleanly logged and available for analysis. Of course, users may see that a friend has adopted a VPF outside of these opportunities, such as when they are visiting a friend’s profile page or via an organic friend post in their feed. This presents opportunities for “treatment” exposure that we were not able to detect, although we don’t believe there is any systematic under-estimation that would skew our conclusions. We also did not control for all confounders in our observational data analysis, either because they were unknown or because we did

not have proxies for all known ones, and therefore the results from the regression model cannot be interpreted causally. Building a complete dependency graph for variables (and their operational proxies) which may be influencing both exposure and adoption would allow us to more completely control for possible confounders, and bring the interpretation of the model coefficients closer to causality.

While the use of the field experiment did bridge this correlation/causality gap to some degree, we also note that the experiment design could have been improved to include multiple factors representing additional promotional formats and cohort properties, ties could have been chosen in the treatment arms more systematically to introduce controlled variation, and we could have selected additional endpoints to collect from users via pre/post treatment surveys to segment treatment effects by key pre-treatment variables and to estimate intent changes.

With respect to our backfire analysis, we looked for increases in specific events that limited direct connections on the Facebook social graph (unfriending, unfollowing, and blocking) upon adoption of an official VPF. We did not study other forms of negative social interactions, such as counter-speech in comments, negative reactions, adoption of anti-vaccination frames, or negative actions against non-adopters. Therefore, we cannot rule out these and other forms of backfire effects.

Finally, we note that VPF adoption is not the final endpoint of interest for public health purposes, and this study did not look at how increased adoption led to increases in intent or uptake.

Despite these opportunities for improving the study design in future work, our present results strengthen previous findings (based largely on small-scale surveys) that there is heightened value in positive vaccine messages containing social proof from close friends and family and that online delivery of such messages can help drive health-related behavior change at scale. We believe this result can help inform design choices made by policymakers and campaign designers to optimize public health communications. Overall, when there is the

opportunity to deliver messages containing social support, and there is a choice in which ties to select, our results argue for including the social proof from the strongest ties possible to most effectively leverage the social influence causal channel (RQ1), and that providing this social proof does not result in social contraction as an unintended side effect (RQ2).



## CHAPTER 5

# A Statistical Model for Quantifying the Needed Duration of Social Distancing for the COVID-19 Pandemic

### 5.1 Introduction

Understanding the near-future implications of the COVID-19 pandemic is one of the most fundamental questions the scientific community is trying to answer in the past few months. While strategies such as social distancing have been widely employed to mitigate the impact of the pandemic on healthcare resources, the necessary timing, frequency, intensity, and effectiveness of these interventions is largely unknown. One of the key unknowns in these strategies is the duration of time for which social distancing needs to be imposed to flatten the pandemic curve. Answering this question requires an accurate model of the transmission trajectory of SARS-CoV-2.

Several recent studies [PLR20, LPC20, KTG20, GBB20] have attempted to understand the transmission trajectory of SARS-CoV-2 using variants of compartmental models, such as the SEIR model [KM27, Het00]. Prem et al. [PLR20] fit an age-structured variant of the SEIR model to case data from Wuhan. They used this model to investigate the effect of lifting restrictions on returning to work and concluded that a premature and sudden lifting of interventions could lead to an early secondary peak. Li et al. [LPC20] fit an SEIR model to SARS-CoV-2 case data across 375 cities in China during 10-23 January 2020. The

model separately considered documented and undocumented infections. Further, they also integrated mobility data across cities. Using their model, they concluded that  $\approx 86\%$  of the cases were undocumented and these undocumented infections were the source of  $\approx 80\%$  of the documented cases. Kissler et al. [KTG20] provide an elegant analysis using a variant of the SEIR model that takes into account various factors that modulate the transmission, including the effects of social distancing, seasonality, immunity, and cross-immunity, resulting in a highly detailed model that can predict, among other things, the time until social distancing is no longer required to flatten the curve. Using their model, they conclude that even under the assumption of full immunity as a response to infection, the time required for social distancing is at least 2022 assuming no vaccine or medication is found by then. Giordano et al. [GGB20] consider a model with eight stages of infection: susceptible (S), infected (I), diagnosed (D), ailing (A), recognized (R), threatened (T), healed (H) and extinct (E). They apply their model to data from Italy to conclude that social distancing will need to be combined with testing and contact tracing to control the pandemic.

Unfortunately, even though the SEIR model is an established model, it is unclear to what extent the accuracy of the prediction of the time of social distancing is affected by the choice of the parameters. Further, the choice of parameters in these models in the context of SARS-CoV-2, has been a subject of debate within the scientific community. One of the key parameters that determine the transmission trajectory is the reproduction number,  $R_0$ . Published values of  $R_0$  range from 1.4 to 7.23 [LGW20, NMS20]. Kissler et al. chose to set peak  $R_0$  as ranging from 2.2 to 2.6, based on the fit of their model to historical data on related coronavirus (HCoV-OC43 and HCoV-HKU1) cases. This choice was made since at the time of publication, there was not enough SARS-CoV-2 data to establish these parameters. However, currently, there is an opportunity to adjust the predictions based on SARS-CoV-2 data as opposed to previous related viruses.

In this work, we fit a statistical model of transmission dynamics building upon the SEIR model. However, instead of fitting this model to previous strains of the SARS virus, we fit

the model to data from current COVID-19 cases. A challenge with our approach arises from the limited case data available in a given location. Particularly, we demonstrate that the key epidemiological parameters that determine the end of social distancing (the reproduction number  $R_0$  and the average time spent in the infectious state  $\tau$ ) have large uncertainties associated with them which, in turn, lead to substantial uncertainties in estimates of the end of social distancing.

To obtain more precise parameter estimates, we formulate a hierarchical Bayesian model that allows the sharing of statistical strength across the location-specific models. Specifically, while each location is allowed to have its own values of the two parameters, these location-specific parameters are assumed to be drawn from a distribution centered around global parameter values. We estimate these global parameters using a marginal likelihood maximization framework. We then use these global parameter estimates, integrating over their uncertainty, to estimate the range of times till the end of social distancing in a new location. The resulting approach not only gives us point estimates (for parameters such as  $R_0$  and for the time to end social distancing) but also provides formal confidence intervals.

We apply our framework to COVID-19 cases from six locations (New York, Spain, Germany, France, Denmark, and the UK) to estimate global and location-specific parameter estimates. We show that these parameters provide a good fit to the data from each of the locations. Finally, we use the global parameter estimates to estimate that the time to end social distancing will be in October 2020 (assuming permanent immunity, no seasonality, and that social distancing reduces the effectiveness of transmission by 60%). We provide open-source software that can be applied to diverse locations to estimate transmission parameters and predict the required duration of social distancing. Although our analysis and motivation stems from the current COVID-19 pandemic, our method is general, and can be applied to other future pandemics.

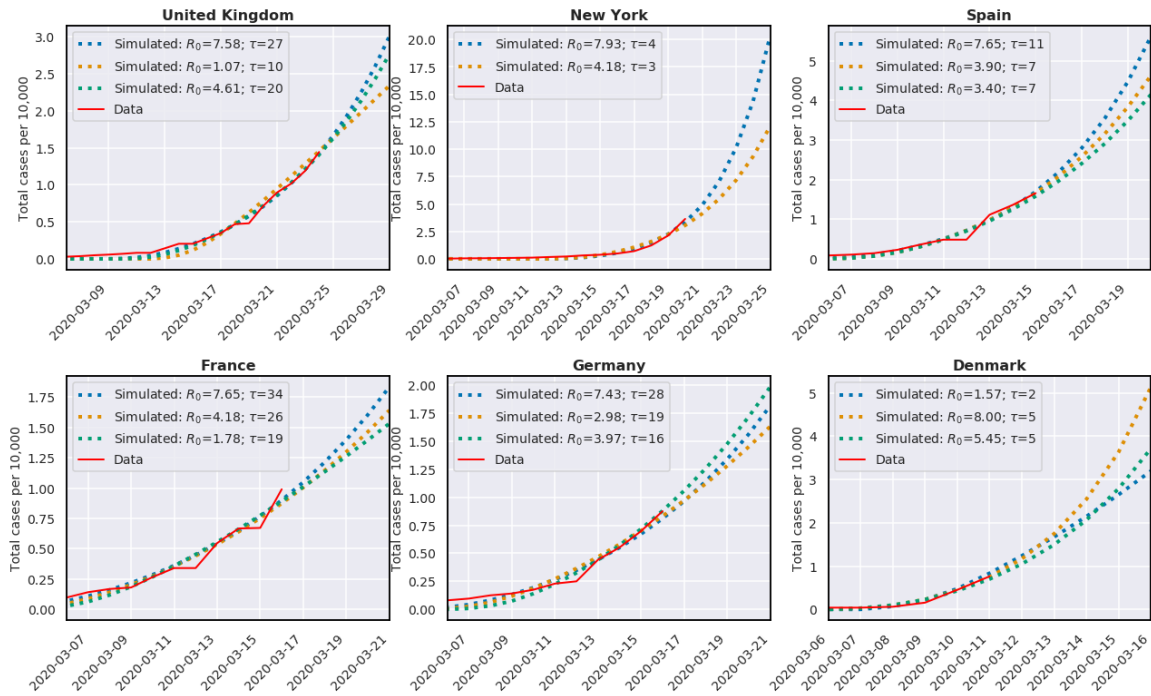


Figure 5.1: Comparison of the observed trajectory of the number of cases in United kingdom, New York, Spain, France, Germany, and Denmark (prior to the date where social distancing was imposed). We provide fits based on region-specific parameters (we choose sets of parameters that all lie within the 95% confidence set). The different sets of parameters diverge significantly in the subsequent dates showing the under-determination of this model.

## 5.2 Results

### 5.2.1 Estimates of $t_{\text{end}}$ from region-specific parameter estimates

We consider COVID-19 data[DDG20, The20] from six locations: UK, Spain, Germany, France, Denmark, and New York. Since our goal is to estimate the parameters ( $R_0$  and  $\tau$ ) in the period when no social distancing was imposed, we restricted our analysis to the dates prior to when social distancing was imposed in each of these regions.

Figure 5.1 shows the parameter estimates when we fit a SEIR model to each of the six regions. While each of the models appears to fit the data in each of the regions, there is

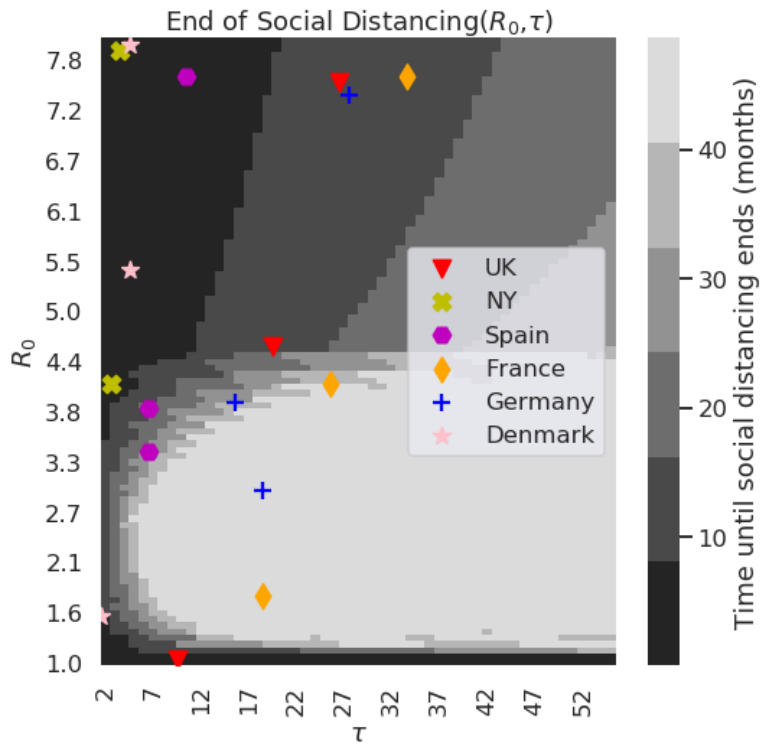


Figure 5.2: The time until social distancing ends (in months) based on the SEIR model, using different  $R_0$  and  $\tau$  values. For each of the regions (Spain, United Kingdom, New York, France, Germany, and Denmark) we also marked the parameters that provided a good fit as shown in Figure 5.1.

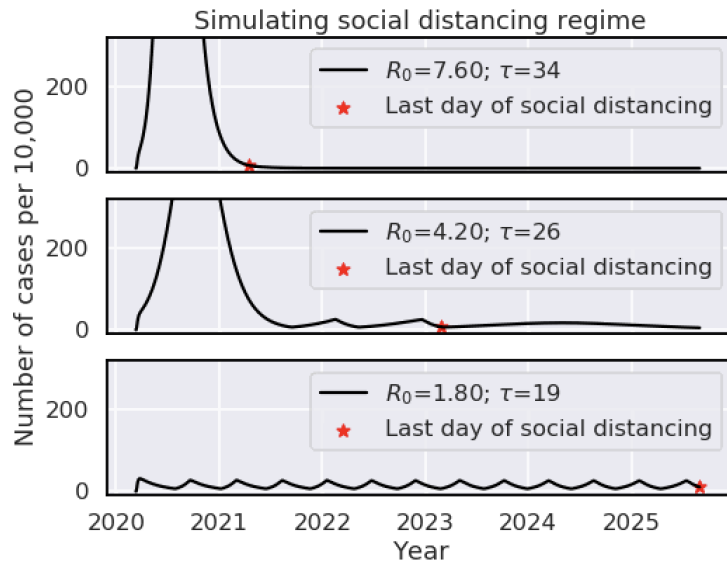


Figure 5.3: Simulating the number of cases under the social distancing regime where social distancing is turned on when the number of cases exceeds 35 per 10,000 and is turned off when it drops below 5 per 10,000. We show 3 different sets of parameters matching data taken from France as seen in Table 2.

Table 5.1: The maximum-likelihood estimates for  $R_0$  and  $\tau$  for every region

Region	$R_0$ [95% confidence interval]	$\tau$ [95% confidence interval]
UK	7.6 [6.6,8.0]	27 [26,31]
New York	7.9 [4.5,8.0]	4 [2,4]
Spain	7.6 [5.6,8.0]	11 [10,11]
France	8.0 [7.1,8.0]	35 [29,35]
Germany	8.0 [6.9,8.0]	28 [23,28]
Denmark	8.0 [1.0,8.0]	5 [2,5]

considerable uncertainty in the parameter estimates (see Table 5.1 for 95% CI). We note that the uncertainty in the key epidemiological parameters that determine the end of social distancing (the reproduction number  $R_0$  and the average time spent in the infectious state  $\tau$ ) leads to substantial uncertainties in estimates of  $t_{end}$ : the time till the end of social distancing (Figure 5.2 and Figure 5.3).

### 5.2.2 Estimates of $t_{end}$ using a Bayesian framework

Due to the large uncertainty in the parameters estimated in each of the locations separately, we fit our model jointly in all locations using a Bayesian framework (see Methods). The Bayesian framework assumes a prior distribution (normal) on the parameters  $R_0^{(0)}$  and  $\tau^{(0)}$ , and it estimates the posterior probability based on the data obtained in each of the countries. The estimated global parameters of the model are  $R_0^{(0)} = 7.5(6.6, 8.0)$ ,  $\sigma_R^2 = 1(1.0, 2.3)$ ,  $\tau^{(0)} = 17(7, 28)$ ,  $\sigma_\tau^2 = 121(49, 529)$ . We observe that our parameter estimates provide an adequate fit to the data in each of the locations (Figure 5.4). We then sample the parameters from the most likely distribution of the parameters  $(R_0, \tau)$  and for each set of parameters we simulate the pandemic scenario, while taking into account that social distancing reduces  $R_0$

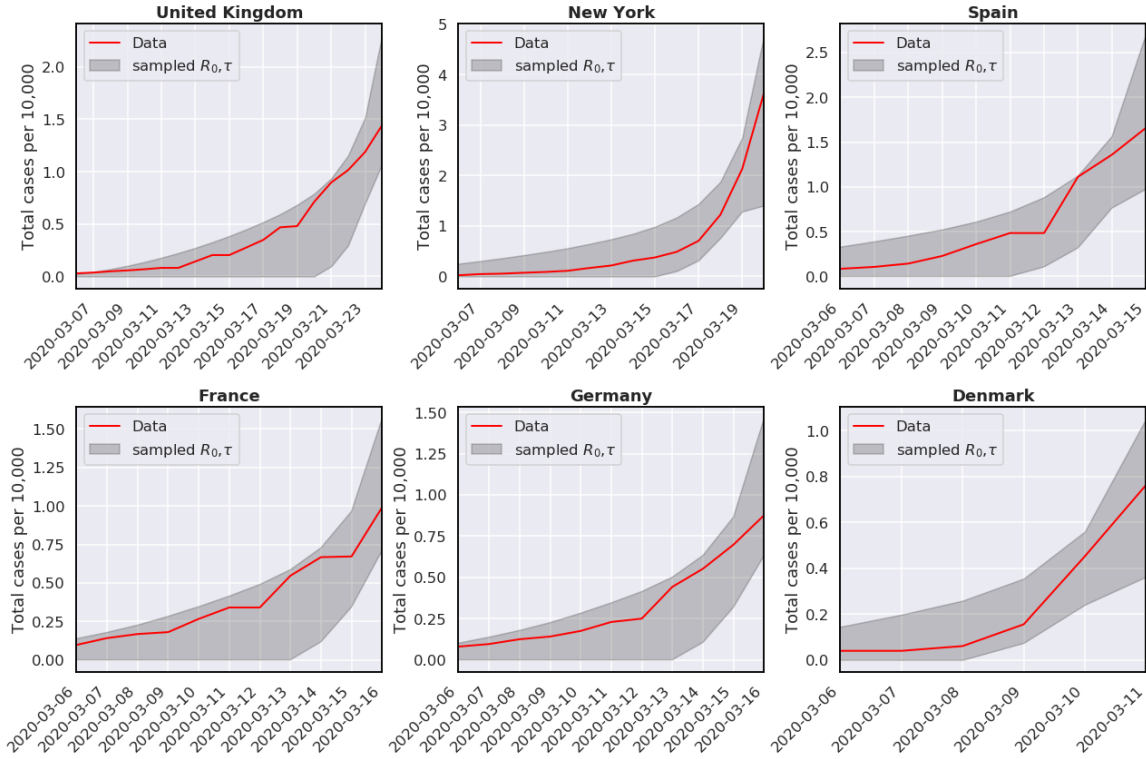


Figure 5.4: The range of trajectories for the number of cases predicted using samples from distribution implied by the global parameters estimated on all regions.

by 60%, and under the assumption that immunity is fixed for life once exposed. The latter assumption is a best-case scenario, i.e., if this assumption is relaxed then the time to social distancing is expected to increase. Furthermore, we assume no seasonality, and again, this results in a lower bound on the time for social distancing. However, since we do not have any strong evidence for specific effects of seasonality, or specific information about the duration of immunity, we chose to focus on this lower bound scenario. Under this scenario, our analysis provides a distribution of possible values for  $t_{end}$  (Figure 5.5). The mode of the distribution is in September 2020, the median is in October 2020 and the variance is 16 months. Based on these results, we obtain a more optimistic view of the time for the end of social distancing [KTG20].



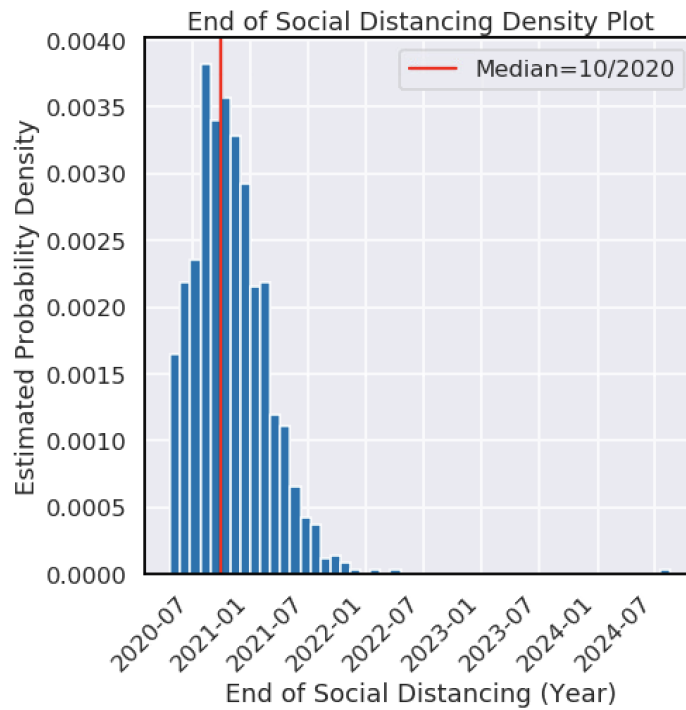


Figure 5.5: The distribution for the time until social distancing will end implied by the global parameters  $R_0^{(0)}, \tau^{(0)}, \sigma_R^2, \sigma_\tau^2$ . The median is October 2020, the mode is September 2020 and the variance is 16 months.

### 5.2.3 Sensitivity analysis.

We first wanted to check how our estimates of  $t_{end}$  are affected by the choice of the specific regions. Out of the six regions (United Kingdom, New York, Spain, France, Germany, and Denmark) we iteratively chose four regions and estimated the global parameters  $(R_0^{(0)}, \sigma_R^2, \tau^{(0)}, \sigma_\tau^2)$ . For each such set of parameters, we estimated the median of the time for social distancing by sampling 1000 samples from the distribution implied by these parameters, resulting in 1000 estimates of the time in which the social distancing will end. We observe that the median  $t_{end}$  is not greatly affected by the choice of the regions, and particularly the medians typically range from September 2020 to April 2021 (Figure 5.6 (a)).

We next wanted to examine the effect of the decrease in  $R_0$  as a result of social distancing on our estimates. We therefore fixed the values of  $R_0$  and  $\tau$  to the maximum marginal likelihood estimates ( $R_0 = 7.5$  and  $\tau = 17$ ), and varied the effect of social distancing on  $R_0$ . Interestingly, this results in a phase transition behavior where the time for social distancing will end within the next year if social distancing has a moderate effect (*i.e.*, it reduces  $R_0$  by less than 60%), or it will end within many years if social distancing has a large effect (Figure 5.6 (b)).

## 5.3 Methods

### 5.3.1 The SEIR Model

We consider the extended SEIR model that have formed the basis of a number of recent studies of SARS-CoV-2 transmission dynamics[KTG20]. This model partitions the population into susceptible, exposed but not yet infectious, infectious (mild), infectious (but not yet hospitalized), infectious (but not yet critical), hospitalized, critical (in the ICU), and removed. Given the state of the population at time  $t$ , *i.e.*, the number of individuals in each of the partitions, the model describes the state of the population at the next time point by a set

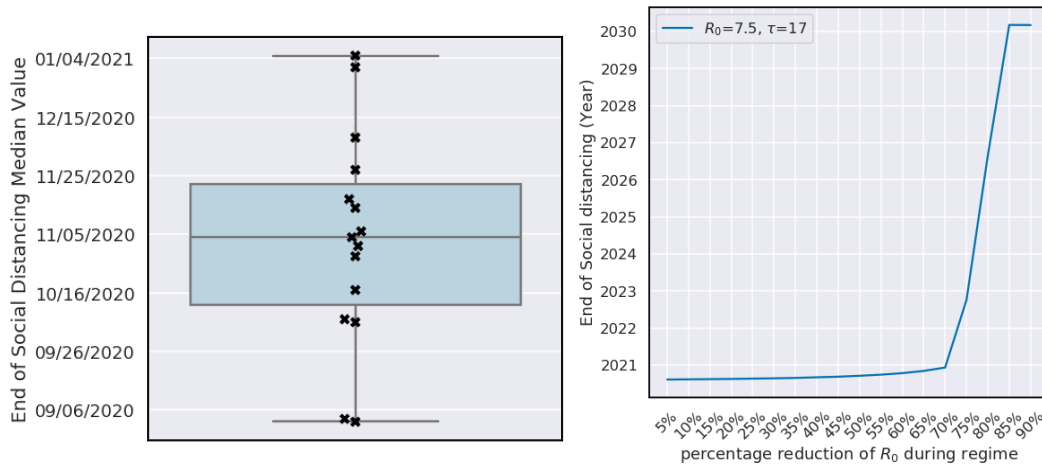


Figure 5.6: (a) Different times until social distancing will end based on different choices of the regions. Each sample was generated by choosing four regions out of the six (UK, France, Spain, Germany, New York, and Denmark), estimating their global parameters, and then measuring the median for the time social distancing will end implied by these parameters. (b) End of social distancing regime as a function of the percentage of  $R_0$  during the regime.

of ordinary differential equations which are governed by a number of parameters, such as the rate at which a susceptible individual is infected and rates at which an individual who is exposed becomes infectious, an infectious individual goes to the hospital, and so on (see Figure 5.7). Given the parameters and the state of the population at some initial time  $t_0$ , this model allows us to compute the state of the population at subsequent times which, in turn, provides a trajectory of cases in the population.

Given the trajectory of SARS-CoV-2 cases from this model, a possible social distancing strategy involves imposing social distancing when the number of critical or hospitalized cases reaches the capacity of the health system and then relaxing social distancing when these numbers are sufficiently small. Depending on the transmission trajectory of SARS-CoV-2, social distancing may need to be imposed multiple times till a sufficiently large number of individuals in the population are immune (assuming that immunity to the virus is permanent). Social distancing is assumed to affect the transmission trajectory by changing the reproduction

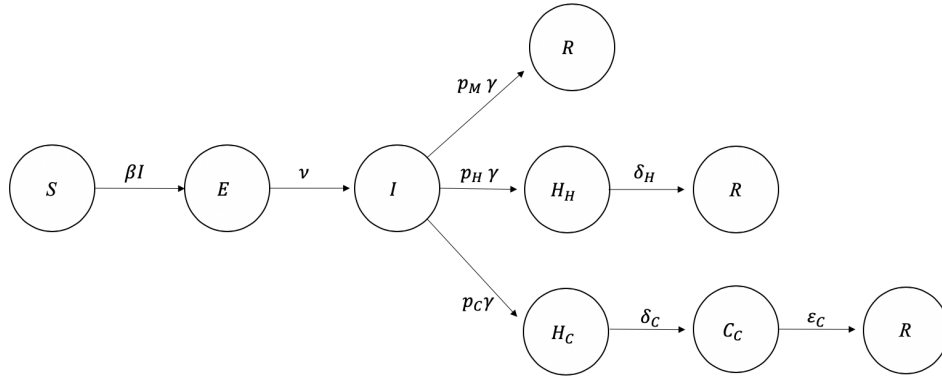


Figure 5.7: SEIR model schema. Each individual in the population begins at susceptible state  $S$ , and will enter into the exposure state  $E$  with transition rate  $\beta I$  in each time unit, where  $\beta = R_0 \gamma$ . In exposure state  $E$ , an individual will go into infectious state  $I$  with transition rate  $\nu$ . Of all the people who arrive at state  $I$ ,  $p_M$  of them will recover (state  $R$ ),  $p_H$  will be hospitalized but will never reach critical care (state  $H_H$ ), and  $p_C$  will be hospitalized to later be in critical care (state  $H_C$ ). All transitions from the  $I$  state will occur with transition rate  $\gamma$ . People in  $H_H$  will enter into  $R$  with a transition rate  $\delta_H$ ; people in  $H_C$  state will enter into critical state  $C_C$  with a transition rate  $\delta_C$ , and then enter into  $R$  state with a transition rate  $\epsilon_C$ . We set parameters  $p_M = 0.956$ ,  $p_H = 0.0308$ ,  $p_C = 0.0132$ ,  $\nu = 1/4.6$ ,  $\delta_C = 1/6$ ,  $\delta_H = 1/8$ ,  $\epsilon_C = 1/10$  as were estimated by Kissler *et al.* All states are normalized with respect to population size  $N$ .

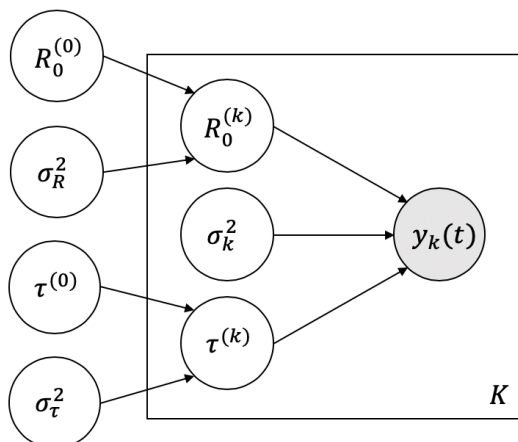


Figure 5.8: Parameter estimation diagram: We assume that the parameters  $R_0^{(k)}$ ,  $\sigma_k^2$ ,  $\tau^{(k)}$  are drawn from a distribution which is defined by the parameters  $R_0^{(0)}$ ,  $\sigma_R^2$ ,  $\tau^{(0)}$ , and  $\sigma_\tau^2$ . We then assume that the cumulative case number curve  $y_k(t)$  is generated by the process defined by these parameters. We estimate the most likely values of  $R_0^{(0)}$ ,  $\sigma_R^2$ ,  $\tau^{(0)}$ , and  $\sigma_\tau^2$  using maximum marginal likelihood approach.

number  $R_0$ .

The key parameters in this model that determine the time till the end of social distancing ( $t_{end}$ ) are the reproduction number ( $R_0$ ) and the average time during which an individual is infectious ( $\tau$ ). The parameter  $\tau$  is related to the rate at which an individual transitions out of the infectious state typically used in the SEIR model ( $\gamma$ ) as  $\tau = \frac{1}{\gamma}$ .

### 5.3.2 A Bayesian hierarchical model for parameter estimation across multiple locations

Since it is unclear whether the parameters that fit HCoV-OC43 and HCoV-HKU1 are also applicable to SARS-CoV-2, we propose an alternate approach, in which we estimate the key parameter values by fitting the SEIR model to contemporary COVID-19 cases from specific locations. The challenge in such an approach is that the limited data in a given location

leads to large uncertainty in the parameter estimates and is very sensitive to outliers.

Our approach to improving the precision of parameter estimates involves fitting SEIR models to all the locations jointly. One possible approach to do so involves setting the parameters to the same value across each location. However, this assumption is unlikely to be realistic. Instead, we endow each location-specific model with its own parameters but assume that the parameters are drawn from a distribution with global parameter values. The SEIR model has a number of parameters that control the transmission trajectory. Our model can jointly estimate all of these parameters. In our analysis, we fix all the parameters to values used in [KTG20] but estimate the values of  $R_0$  and  $\tau$ .

We assume that we have data on the observed number of COVID-19 cases from  $K$  locations:  $\{y_k(t)\}, t \in \{1, \dots, T_k\}, k \in \{1, \dots, K\}$ . Let  $f(t; (R_0, \tau))$  denote the number of infections at time  $t$  predicted by the SEIR model with parameters  $(R_0, \tau)$ . The parameters for each region are denoted  $(R_0^{(k)}, \tau^{(k)})$  and the global parameters:  $(R_0^{(0)}, \tau^{(0)})$ .

$$\begin{aligned} R_0^{(k)} &\sim \mathcal{N}(R_0^{(0)}, \sigma_R^2) \\ \tau^{(k)} &\sim \mathcal{N}(\tau^{(0)}, \sigma_\tau^2) \\ y_k(t) | (R_0^{(k)}, \tau^{(k)}, \sigma_k^2) &\sim \mathcal{N}(f(t; (R_0^{(k)}, \tau^{(k)})), \sigma_k^2) \end{aligned}$$

Each of the region-specific parameters is drawn from a normal distribution with a mean given by the global parameters. The observed cases in region  $k$  at time  $t$  are drawn from a normal distribution with mean given by the prediction from the SEIR model  $f(t; (R_0^{(k)}, \tau^{(k)}))$  with a region-specific noise variance  $\sigma_k^2$ . Further, we impose an uninformative prior on the noise variance:  $P(\sigma_k^2) \propto \frac{1}{\sigma_k^2}$ . The parameter selection schema is shown as Figure 5.8

We then have:

$$\begin{aligned} P(y_k(1 : T_k) | R_0^{(k)}, \tau^{(k)}, \sigma_k^2) &= \prod_{i=1}^{n_k} \left( \frac{1}{2\pi\sigma_k^2} \right)^{\frac{1}{2}} e^{-\frac{(y_k(t_i) - f(t_i; R_0^{(k)}, \tau^{(k)}))^2}{2\sigma_k^2}} \\ &= \left( \frac{1}{2\pi\sigma_k^2} \right)^{\frac{n_k}{2}} e^{-\frac{\sum (y_k(t_i) - f(t_i; R_0^{(k)}, \tau^{(k)}))^2}{2\sigma_k^2}} \end{aligned} \tag{5.1}$$

$$\begin{aligned}
P(y_k(1 : T_k) | R_0^{(k)}, \tau^{(k)}) &= P(y_k(1 : T_k) | f(t; R_0^{(k)}, \tau^{(k)})) \\
&= \int_0^\infty P(y_k(1 : T_k) | f(t; R_0^{(k)}, \tau^{(k)}), \sigma_k^2) P(\sigma_k^2) d\sigma_k^2 \\
&= \int_0^\infty \left( \frac{1}{2\pi\sigma_k^2} \right)^{\frac{n_k}{2}} e^{-\frac{\sum(y_k(t_i) - f(t_i; R_0^{(k)}, \tau^{(k)}))^2}{2\sigma_k^2}} \frac{1}{\sigma_k^2} d\sigma_k^2 \\
&= \left( \frac{1}{2\pi} \right)^{\frac{n_k}{2}} \int_0^\infty \left( \frac{1}{\sigma_k^2} \right)^{\frac{n_k}{2} + 1} e^{-\frac{\sum(y_k(t_i) - f(t_i; R_0^{(k)}, \tau^{(k)}))^2}{2\sigma_k^2}} d\sigma_k^2
\end{aligned} \tag{5.2}$$

Using the fact that the integrand in Equation 5.2 is a Gamma function, we have:

$$P(y_k(1 : T_k) | R_0^{(k)}, \tau^{(k)}) \propto \Gamma\left(\frac{n_k}{2}\right) \left( \frac{2}{\sum_{i=1}^{n_k} (y_k(t_i) - f(t_i; R_0^{(k)}, \tau^{(k)}))^2} \right)^{\frac{n_k}{2}} \tag{5.3}$$

We then compute the maximum marginal likelihood estimates of the global parameters using a grid search:

$$(\hat{R}_0^{(0)}, \hat{\sigma}_R^2, \hat{\tau}^{(0)}, \hat{\sigma}_\tau^2) = \arg \max_{(R_0^{(0)}, \tau^{(0)}, \sigma_R^2, \sigma_\tau^2)} l(R_0^{(0)}, \sigma_R^2, \tau^{(0)}, \sigma_\tau^2)$$

$$\begin{aligned}
l(R_0^{(0)}, \sigma_R^2, \tau^{(0)}, \sigma_\tau^2) &= \log P(\{y_k(1 : T_k)\}_{k=1, \dots, K} | (R_0^{(0)}, \sigma_R^2, \tau^{(0)}, \sigma_\tau^2)) \\
&= \log \prod_{k=1}^K P(y_k(1 : T_k) | (R_0^{(0)}, \sigma_R^2, \tau^{(0)}, \sigma_\tau^2)) \\
&= \sum_{k=1}^K \log P(y_k(1 : T_k) | (R_0^{(0)}, \sigma_R^2, \tau^{(0)}, \sigma_\tau^2))
\end{aligned}$$

We evaluate each term in the log likelihood as:

$$\begin{aligned}
&P(y_k(1 : T_k) | (R_0^{(0)}, \sigma_R^2, \tau^{(0)}, \sigma_\tau^2)) \\
&= \int P(y_k(1 : T_k) | R_0^{(k)}, \tau^{(k)}) P(R_0^{(k)} | R_0^{(0)}, \sigma_R^2) P(\tau^{(k)} | \tau^{(0)}, \sigma_\tau^2) dR_0^{(k)} d\tau^{(k)}
\end{aligned} \tag{5.4}$$

The integral in Equation 5.4 does not have an analytical solution so we evaluate the integral numerically over a grid of values for  $(R_0^{(k)}, \tau^{(k)})$ .

The grid search of the parameters in the likelihood searches for values of  $R_0^{(0)}$  between 1 and 8,  $\tau^{(0)}$  between 2 to 55,  $\sigma_R$  from 1 to 8, and  $\sigma_\tau$  from 1 to 30.

### 5.3.3 Application to predict the end of social distancing

We estimate  $t_{end}$ , the time when social distancing can be ended, in the following way. First, using a maximum marginal likelihood approach, we find the most likely parameters  $(R_0^{(0)}, \sigma_R^2, \tau^{(0)}, \sigma_\tau^2)$ . Then, we sample  $R_0, \tau$  from the distribution  $R_0 \sim \mathcal{N}(R_0^{(0)}, \sigma_R^2), \tau \sim \mathcal{N}(\tau^{(0)}, \sigma_\tau^2)$  and for each such sample we compute the estimated value of  $t_{end}$  as follows. We follow the parameter choices used in [KTG20]: assuming that immunity to SARS-CoV-2 is permanent (which provides the minimum time of social distancing), that social distancing is imposed when the number of cases exceeds 35 per 10,000 individuals and is relaxed when the number of cases drops below 5 per 10,000 individuals (these thresholds were chosen so that the number of hospital cases is below the capacity in the United States), and that each period of social distancing reduces  $R_0$  by 60%. We then simulate the SEIR scenario based on the above parameters, including  $R_0$  and  $\tau$ . This results in a distribution of values  $t_{end}$ . Additionally, we performed a sensitivity analysis where we demonstrate the effect of the choice of each of the above parameters.

## 5.4 Discussion

In this work, we fit a statistical model of transmission dynamics based on the SEIR model to data from COVID-19 cases from multiple locations. Our approach uses a Bayesian framework, resulting in a distribution of end dates for social distancing, as opposed to a specific end time, incorporating the uncertainty in the parameter choices of the model. This uncertainty is inherent to the SARS-CoV-2 pandemic, as can be viewed by the fact that  $R_0$  has been ranging in the literature from 1.4 to 7.23 [LGW20, NMS20]. We show that our approach provides a good fit for the COVID-19 cases in these locations. Our approach demonstrates that the end of social distancing will be around October 2020, under mild assumptions.

It is important to note that the assumptions made by our analysis provide a lower bound on the time for social distancing. Particularly, we assume no seasonality; if COVID-19 is



seasonal, we expect greater spread to appear in winter relative to summer (as has been observed for influenza). However, it is not clear whether COVID-19 is seasonal, and if so, to what extent, and we therefore leave this aspect for future analysis, once more data will be available. Similarly, it is currently unclear whether one acquires permanent immunity or for a short duration after getting exposed to the disease. Thus, given the lack of information about immunity, we chose to make the best-case scenario assumption in which immunity is acquired for life. Other assumptions may prolong the effects of social distancing.

Critically, other interventions such as the introduction of a vaccine, the introduction of effective medications, or the introduction of a larger number of clinical care resources such as ventilators, will change the scenarios provided in this analysis. Specifically, the introduction of a vaccine or effective medications will likely alter the parameters  $R_0$  and  $\tau$ , and would therefore result in a shorter time for social distancing. Additional resources such as ventilators will result in different thresholds set for the application of social distancing, as social distancing will only be required when there is a risk for a surge that collapses the health systems. In that case, again, the end of social distancing is expected to arrive sooner.

The above limitations need to be taken into account when interpreting our analysis. However, we note that as new data on immunity, seasonality, medications, and vaccines becomes available, these can easily be incorporated into our framework, and a revised analysis can be performed. We provide freely available code that allows for such an analysis by researchers in the community (see appendix).

Finally, we would like to point out that the issue of sensitivity of the model to the parameter choices is not specific to the SEIR model. Specifically, in our hands we have observed a similar phenomenon for other models as well (data not shown). The limited availability of data limits the certainty with which parameters of the model can be identified. Thus, it is critical that estimates from the application of statistical models to such data be accompanied by formal measures of uncertainty. We believe that the statements resulting from our analysis should also be taken in the context of the specific locations we analyzed

and the specific model that we used. Possibly, other models or other locations may provide different estimates.

Table 5.2: The date of the end of social distancing for different sets of parameters that fit the data (as shown in Figures 5.1 and 5.2)

Region used for parameter estimation	$R_0$	$\tau$	Estimated end of social distancing
United Kingdom	1.1	10	Apr, 2020
New York	7.9	4	May, 2020
Denmark	8.0	5	Jun, 2020
Denmark	5.4	5	Jul, 2020
Spain	7.6	11	Aug, 2020
New York	4.2	3	Sep, 2020
Denmark	1.6	2	Nov, 2020
United Kingdom	7.6	27	Jan, 2021
Germany	7.4	28	Feb, 2021
United Kingdom	4.6	20	Apr, 2021
France	7.6	34	Apr, 2021
Spain	3.9	7	Aug, 2021
Germany	4.0	16	Aug, 2022
Spain	3.4	7	Sep, 2022
France	4.2	26	Feb, 2023
Germany	3.0	19	Feb, 2024
France	1.8	19	Mar, 2024

## APPENDIX A

Supplementary Material - The heterogeneous effects of social support on the adoption of Facebook's vaccine profile frames feature

RQ1

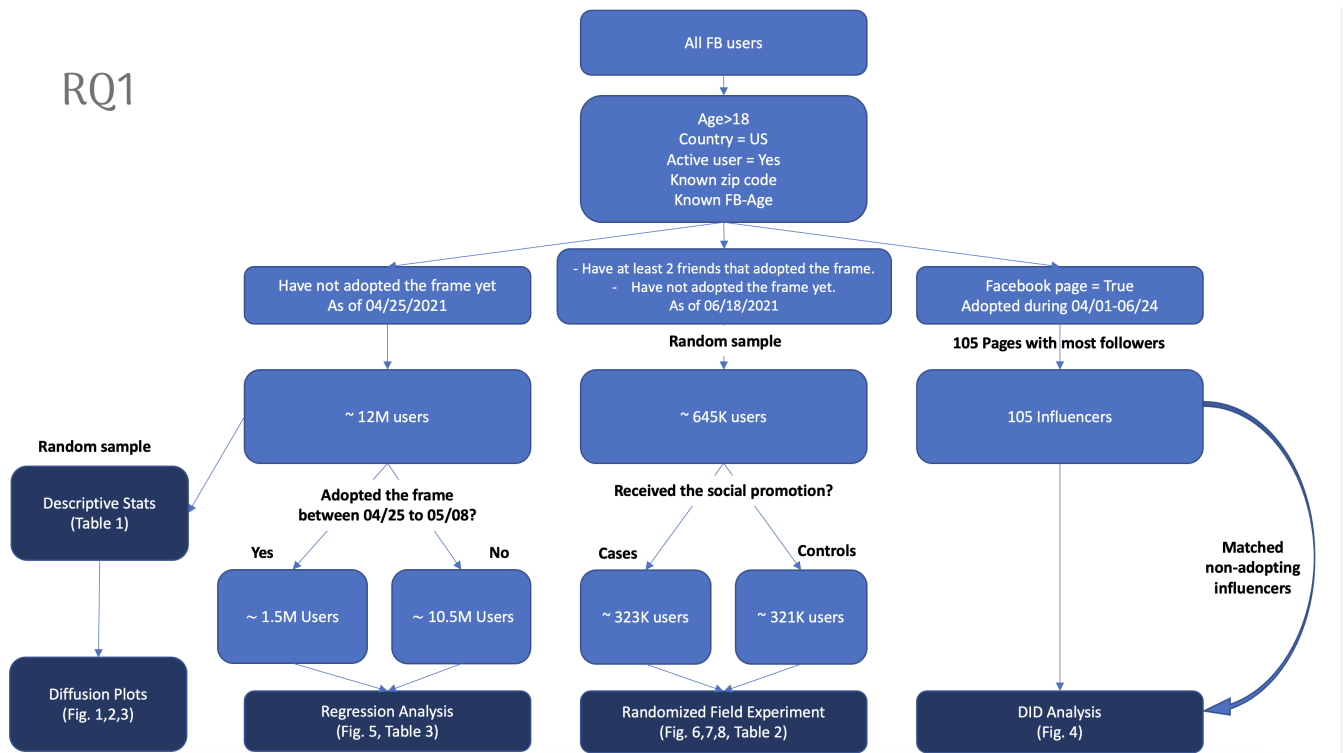


Figure A.1: Flowchart illustrating user and data selection for RQ1.

RQ2

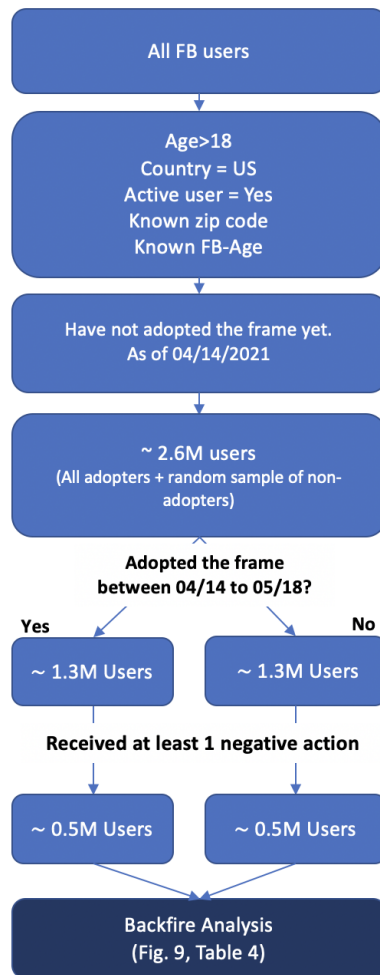


Figure A.2: Flowchart illustrating user and data selection for RQ2.

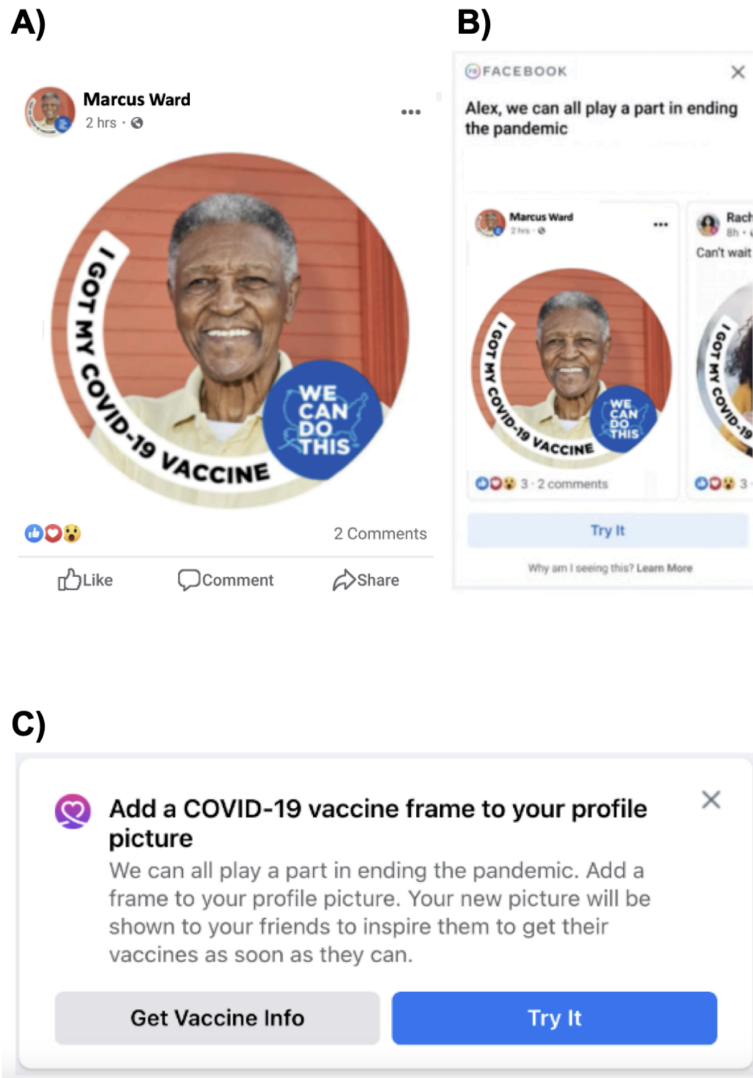


Figure A.3: Different messages promoting adoption of a Vaccine Profile Frame (VPF) **A)** VPF Post - A post that is automatically generated upon adoption and displayed to friends within their newsfeed. **B)** Friend Aggregation Post - A newsfeed post informing users that three of their friends have adopted a VPF. **C)** Profile/Newsfeed Notification - A non-social notification presented on either the user's profile page or in their newsfeed encouraging them to adopt the frame.

## APPENDIX B

### Supplementary Material - A Statistical Model for Quantifying the Needed Duration of Social Distancing for the COVID-19 Pandemic

#### B.1 Ordinary differential equations

The SEIR model described by Kissler *et al.* is defined by the following set of ordinary differential equations:

$$\begin{aligned}\frac{dS}{dt} &= -\beta SI; & \frac{dE}{dt} &= \beta SI - \nu E \\ \frac{dI}{dt} &= \nu E - \gamma I; & \frac{dH_H}{dt} &= \gamma p_H I - \delta_H H_H \\ \frac{dH_C}{dt} &= \gamma p_C I - \delta_C H_C; & \frac{dC_C}{dt} &= \delta_C H_C - \epsilon_C C_C \\ \frac{dR}{dt} &= \epsilon_C C_C + \delta_H H_H + \gamma p_M I\end{aligned}$$

#### B.2 Code availability

The code used to generate all figures and experiments in this paper can be found here:

[https://github.com/doubleBlindGit/COVID19\\_SocialDistance](https://github.com/doubleBlindGit/COVID19_SocialDistance)



## REFERENCES

- [AEO21] M. Agranov, M. Elliott, and P. Ortoleva. “The importance of social norms against strategic effects: the case of Covid-19 vaccine uptake.” *Econ Lett*, **206**, 2021.
- [AEP07] Andreas Argyriou, Theodoros Evgeniou, and Massimiliano Pontil. “Multi-task feature learning, in ‘Advances in Neural Information Processing Systems 19’.”, 2007.
- [An19] G. An. “Glaucoma diagnosis with machine learning based on optical coherence tomography and color fundus images.” *J. Healthc. Eng.*, **2019**, 2019.
- [Asa19] R. Asaoka. “Using deep learning and transfer learning to accurately diagnose early-onset glaucoma from macular optical coherence tomography images.” *Am. J. Ophthalmol.*, **198**, 2019.
- [ATW19] Susan Athey, Julie Tibshirani, and Stefan Wager. “Generalized random forests.”, 2019.
- [AW14] S. i. n. a. n. Aral and D. y. l. a. n. Walker. “Tie strength, embeddedness, and social influence: a large-scale networked experiment.” *Manag Sci*, **60**, 2014.
- [BBC15] C. Betsch, R. Böhm, and G. B. Chapman. “Using behavioral insights to increase vaccination policy effectiveness.” *Policy Insights Behav Brain Sci*, **2**, 2015.
- [Bon12] R. M. Bond. “A 61-million-person experiment in social influence and political mobilization.” *Nature*, **489**, 2012.
- [Bre17] N. T. Brewer. “Increasing vaccination: putting psychological science into action.” *Psychol Sci Public Interest*, **18**, 2017.
- [BRM12] Eytan Bakshy, Itamar Rosenn, Cameron Marlow, and Lada Adamic. “The role of social networks in information diffusion.” In *Proceedings of the 21st international conference on World Wide Web, WWW ’12*, pp. 519–528, New York, NY, USA, April 2012. Association for Computing Machinery.
- [Bru13] E. K. Brunson. “The impact of social networks on parents’ vaccination decisions.” *Pediatrics*, **131**, 2013.
- [Bru19] W. Bruine de Bruin. “Reports of social circles’ and own vaccination behavior: a national longitudinal survey.” *Health Psychol*, **38**, 2019.
- [BW15] Marco A Bonini Filho and Andre J Witkin. “Outer retinal layers as predictors of vision loss.” *Rev Ophthalmol*, **15**, 2015.

- [CDD15] Yashin Dicente Cid, Oscar Alfonso Jiménez Del Toro, Adrien Depeursinge, and Henning Müller. “Efficient and fully automatic segmentation of the lungs in CT volumes.” In *VISCERAL Challenge@ ISBI*, pp. 31–35, 2015.
- [CF13] N. A. Christakis and J. H. Fowler. “Social contagion theory: examining dynamic social networks and human behavior.” *Stat Med*, **32**, 2013.
- [CGC18] Noel C. F. Codella, David Gutman, M. Emre Celebi, Brian Helba, Michael A. Marchetti, Stephen W. Dusza, Aadi Kalloo, Konstantinos Liopyris, Nabin Mishra, Harald Kittler, and Allan Halpern. “Skin lesion analysis toward melanoma detection: A challenge at the 2017 International symposium on biomedical imaging (ISBI), hosted by the international skin imaging collaboration (ISIC).” In *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, pp. 168–172, 2018.
- [CLR] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. “Decision Transformer: Reinforcement Learning via Sequence Modeling.”
- [COP18] N. Cuenca, I. Ortuño-Lizarán, and I. Pinilla. “Cellular characterization of OCT and outer retinal bands using specific immunohistochemistry markers and clinical implications.” *Ophthalmology*, **125**, 2018.
- [CQY16] Hao Chen, Xiaojuan Qi, Lequan Yu, and Pheng-Ann Heng. “DCAN: deep contour-aware networks for accurate gland segmentation.” In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 2487–2496, 2016.
- [CV] Sébastien Cortaredona and Bruno Ventelou. “The Extra Cost of Comorbidity: Multiple Illnesses and the Economic Burden of Non-Communicable Diseases.” **15**:216.
- [DDA18] J. Dolz, C. Desrosiers, and I. B. Ayed. “3D fully convolutional networks for subcortical segmentation in MRI: a large-scale study.” *Neuroimage*, **170**, 2018.
- [DDG20] Ensheng Dong, Hongru Du, and Lauren Gardner. “An interactive web-based dashboard to track COVID-19 in real time.” *Lancet Infect. Dis.*, **20**(5):533–534, May 2020.
- [DDS09] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. “Imagenet: A large-scale hierarchical image database.” In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- [DWB11] Jan Dörr, Klaus D Wernecke, Markus Bock, Gunnar Gaede, Jens T Wuerfel, Caspar F Pfueller, Judith Bellmann-Strobl, Alina Freing, Alexander U Brandt, and Paul Friedemann. “Association of retinal and macular damage with brain atrophy in multiple sclerosis.” *PLoS One*, **6**(4):e18132, 2011.

- [EBC09] Dumitru Erhan, Yoshua Bengio, Aaron Courville, and Pascal Vincent. “Visualizing higher-layer features of a deep network.” *University of Montreal*, **1341**(3):1, 2009.
- [Fac] A I Facebook. “CS 4803 / 7643: Deep learning guest lecture: Embeddings and world2vec.” [https://www.cc.gatech.edu/classes/AY2020/cs7643\\_spring/slides/L13\\_Embedding\\_world2vec\\_final\\_version.pdf](https://www.cc.gatech.edu/classes/AY2020/cs7643_spring/slides/L13_Embedding_world2vec_final_version.pdf). Accessed: 2023-9-13.
- [Fau18] J. D. Fauw. “Clinically applicable deep learning for diagnosis and referral in retinal disease.” *Nat. Med.*, **24**, 2018.
- [FMF22] A. Fronczak, M. J. Mrowinski, and P. Fronczak. “Scientific success from the perspective of the strength of weak ties.” *Sci Rep*, **12**, 2022.
- [Fre86] Dieter Frey. “Recent Research on Selective Exposure to Information.” In Leonard Berkowitz, editor, *Advances in Experimental Social Psychology*, volume 19, pp. 41–80. Academic Press, January 1986.
- [GAL21] S. Graupensperger, D. A. Abdallah, and C. M. Lee. “Social norms and vaccine uptake: College students’ COVID vaccination intentions, attitudes, and estimated peer norms and comparisons with influenza vaccine.” *Vaccine*, **39**, 2021.
- [GBB20] Giulia Giordano, Franco Blanchini, Raffaele Bruno, Patrizio Colaneri, Alessandro Di Filippo, Angela Di Matteo, and Marta Colaneri. “Modelling the COVID-19 epidemic and implementation of population-wide interventions in Italy.” *Nature Medicine*, pp. 1–6, 2020.
- [GGM20] Matthew H Goldberg, Abel Gustafson, Edward Maibach, Sander van der Linden, Matthew T Ballew, Parrish Bergquist, John Kotcher, Jennifer R Marlon, Seth A Rosenthal, and Anthony Leiserowitz. “Social norms motivate COVID-19 preventive behaviors.” May 2020.
- [GGS16] H. Greenspan, B. Ginneken, and R. M. Summers. “Guest editorial deep learning in medical imaging: overview and future promise of an exciting new technique.” *IEEE Trans. Med. Imaging*, **35**, 2016.
- [GMW10] S. Goel, W. Mason, and D. J. Watts. “Real and perceived attitude agreement in social networks.” *J Personal Soc Psychol*, **99**, 2010.
- [Gra] Mark S. Granovetter. “The Strength of Weak Ties.” **78**(6):1360–1380.
- [GT13] D. Grewal and A. Tanna. “Diagnosis of glaucoma and detection of glaucoma progression using spectral domain optical coherence tomography.” *Curr. Opin. Ophthalmol.*, **24**, 2013.
- [Het00] Herbert W Hethcote. “The mathematics of infectious diseases.” *SIAM review*, **42**(4):599–653, 2000.

- [HG20] J. Howard and S. Gugger. “Fastai: A Layered API for Deep Learning.” *Information*, **11**, 2020.
- [HR] Hernán and Robins. “Selection bias.” *Causal Inference: What If*. Boca Raton: Chapman and.
- [HSV17] Xiaojie Huang, Junjie Shan, and Vivek Vaidya. “Lung nodule detection in CT using 3D convolutional neural networks.” In *2017 IEEE 14th International Symposium on Biomedical Imaging (ISBI 2017)*, pp. 379–383. IEEE, 2017.
- [HSZ] Shengchao Hu, Li Shen, Ya Zhang, Yixin Chen, and Dacheng Tao. “On Transforming Reinforcement Learning by Transformer: The Development Trajectory.”
- [Hum] Kyle Humphrey. “Using Reinforcement Learning to Personalize Dosing Strategies in a Simulated Cancer Trial with High Dimensional Data.”
- [HZR16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. “Deep residual learning for image recognition.” In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- [Imb] Guido W. Imbens. “The Role of the Propensity Score in Estimating Dose-Response Functions.” **87**(3):706–710.
- [IR] Guido W. Imbens and Donald B. Rubin. *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction*. Cambridge University Press.
- [Irv19] J. Irvin. “Chexpert: a large chest radiograph dataset with uncertainty labels and expert comparison.” *Proc. AAAI Conf. Artif. Intell.*, **33**, 2019.
- [Ive21] Mike Ives. “Celebrities Are Endorsing Covid Vaccines. Does It Help?” *The New York Times*, May 2021.
- [JLL] Michael Janner, Qiyang Li, and Sergey Levine. “Offline Reinforcement Learning as One Big Sequence Modeling Problem.”
- [JLT18] H. Jin, Z. Li, R. Tong, and L. Lin. “A deep 3D residual CNN for false-positive reduction in pulmonary nodule detection.” *Med. Phys.*, **45**, 2018.
- [JSB19] H. Jiang, T. Shi, Z. Bai, and L. Huang. “AHCNet: an application of attention mechanism and hybrid connection for liver tumor segmentation in CT volumes.” *IEEE Access*, **7**, 2019.
- [Kah18] T. A. Kah. “CuRRL syndrome: a case series.” *Acta Scientifica Ophthalmology*, **1**, 2018.
- [Kam17] K. Kamnitsas. “Efficient multi-scale 3D CNN with fully connected crf for accurate brain lesion segmentation.” *Med. Image Anal.*, **36**, 2017.

- [KB14] Diederik P Kingma and Jimmy Ba. “Adam: A method for stochastic optimization.” *arXiv preprint arXiv:1412.6980*, 2014.
- [KCB] Matthieu Komorowski, Leo A. Celi, Omar Badawi, Anthony C. Gordon, and A. Aldo Faisal. “The Artificial Intelligence Clinician Learns Optimal Treatment Strategies for Sepsis in Intensive Care.” **24**(11):1716–1720.
- [Kea12] P. A. Keane. “Evaluation of age-related macular degeneration with optical coherence tomography.” *Surv. Ophthalmol.*, **57**, 2012.
- [Ker18] D. S. Kermany. “Identifying medical diagnoses and treatable diseases by image-based deep learning.” *Cell*, **172**, 2018.
- [KGGK21] Pinelopi Konstantinou, Katerina Georgiou, Navin Kumar, Maria Kyprianidou, Christos Nicolaides, Maria Karekla, and Angelos P Kassianos. “Transmission of Vaccination Attitudes and Uptake Based on Social Contagion Theory: A Scoping Review.” *Vaccines (Basel)*, **9**(6), June 2021.
- [KM27] William Ogilvy Kermack and Anderson G McKendrick. “A contribution to the mathematical theory of epidemics.” *Proceedings of the royal society of london. Series A, Containing papers of a mathematical and physical character*, **115**(772):700–721, 1927.
- [KSB17] Sergey Korolev, Amir Safiullin, Mikhail Belyaev, and Yulia Dodonova. “Residual and plain convolutional neural networks for 3D brain MRI classification.” In *2017 IEEE 14th international symposium on biomedical imaging (ISBI 2017)*, pp. 835–838. IEEE, 2017.
- [KSE21] N. C. Krämer, V. Sauer, and N. Ellison. “The strength of weak ties revisited: further evidence of the role of strong ties in the provision of online social support.” *Soc Media+Soc*, **7**, 2021.
- [KSH12] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. “ImageNet Classification with Deep Convolutional Neural Networks.” In *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012.
- [KTG20] Stephen M Kissler, Christine Tedijanto, Edward Goldstein, Yonatan H Grad, and Marc Lipsitch. “Projecting the transmission dynamics of SARS-CoV-2 through the postpandemic period.” *Science*, 2020.
- [Kuw19] S. Kuwayama. “Automated detection of macular diseases by optical coherence tomography and artificial intelligence machine learning of optical coherence tomography images.” *J. Ophthalmol.*, **2019**, 2019.

- [KZG18] Daniel Kermany, Kang Zhang, Michael Goldbaum, et al. “Labeled optical coherence tomography (oct) and chest x-ray images for classification.” *Mendeley data*, **2**(2):651, 2018.
- [Lau22] B. H. P. Lau. “Understanding the societal factors of vaccine acceptance and hesitancy: evidence from Hong Kong.” *Public Health*, **207**, 2022.
- [Laz21] J. V. Lazarus. “A global survey of potential acceptance of a COVID-19 vaccine.” *Nat Med*, **27**, 2021.
- [LBA17] J. Lei, S. Balasubramanian, N. S. Abdelfattah, M. G. Nittala, and S. R. Sadda. “Proposal of a simple optical coherence tomography-based scoring system for progression of age-related macular degeneration.” *Graefe’s Arch. Clin. Exp. Ophthalmol.*, **255**, 2017.
- [LBH15] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. “Deep learning.” *Nature*, **521**(7553):436–444, May 2015. Publisher: Nature Publishing Group, a division of Macmillan Publishers Limited. All Rights Reserved.
- [LGW20] Ying Liu, Albert A Gayle, Annelies Wilder-Smith, and Joacim Rocklöv. “The reproductive number of COVID-19 is higher compared to SARS coronavirus.” *Journal of travel medicine*, 2020.
- [LL21] Y. Lu and J. K. Lee. “Determinants of cross-cutting discussion on Facebook: political interest, news consumption, and strong-tie heterogeneity.” *New Media Soc*, **23**, 2021.
- [LLG15] Mian Mian Lau, King Hann Lim, and Alpha Agape Gopalai. “Malaysia traffic sign recognition with convolutional neural network.” In *2015 IEEE international conference on digital signal processing (DSP)*, pp. 1006–1010. IEEE, 2015.
- [Loo21] S. Loomba. “Measuring the impact of COVID-19 vaccine misinformation on vaccination intent in the UK and USA.” *Nat Hum Behav*, **5**, 2021.
- [Lor21] Taylor Lorenz. “To Fight Vaccine Lies, Authorities Recruit an ‘Influencer Army’.” *The New York Times*, August 2021.
- [LPC20] Ruiyun Li, Sen Pei, Bin Chen, Yimeng Song, Tao Zhang, Wan Yang, and Jeffrey Shaman. “Substantial undocumented infection facilitates the rapid dissemination of novel coronavirus (SARS-CoV-2).” *Science*, **368**(6490):489–493, 2020.
- [LSN] Siqi Liu, Kay Choong See, Kee Yuan Ngiam, Leo Anthony Celi, Xingzhi Sun, and Mengling Feng. “Reinforcement Learning for Clinical Decision Support in Critical Care: Comprehensive Review.” **22**(7):e18477.

- [LSP] Mingyang Liu, Xiaotong Shen, and Wei Pan. “Deep Reinforcement Learning for Personalized Treatment Recommendation.” *41*(20):4034–4056.
- [Lut15] E. S. Lutkenhoff. “Thalamic and extrathalamic mechanisms of consciousness after severe brain injury.” *Ann. Neurol.*, **78**, 2015.
- [MCG21] Alex Moehring, Avinash Collis, Kiran Garimella, M Amin Rahimian, Sinan Aral, and Dean Eckles. “Surfacing norms to increase vaccine acceptance.” *SSRN Electron. J.*, 2021.
- [Met20] Meta. “Keeping people safe and informed about the Coronavirus.” <https://about.fb.com/news/2020/12/coronavirus/>, December 2020. Accessed: 2022-1-21.
- [Met21] Meta. “Encourage your friends to get a COVID-19 vaccine.” <https://about.fb.com/news/2021/04/encourage-your-friends-to-get-a-covid-19-vaccine/>, April 2021. Accessed: 2022-1-21.
- [MNA16] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. “V-net: Fully convolutional neural networks for volumetric medical image segmentation.” In *2016 fourth international conference on 3D vision (3DV)*, pp. 565–571. Ieee, 2016.
- [Nas18] M. Nassisi. “Quantity of intraretinal hyperreflective foci in patients with intermediate age-related macular degeneration correlates with 1-year progression.” *Invest. Ophthalmol. Vis. Sci.*, **59**, 2018.
- [Nas19] M. Nassisi. “OCT risk factors for development of late age-related macular degeneration in the fellow eyes of patients enrolled in the HARBOR study.” *Ophthalmology*, **126**, 2019.
- [Nit19] M. G. Nittala. “AMISH EYE STUDY: baseline spectral domain optical coherence tomography characteristics of age-related macular degeneration.” *Retina*, **39**, 2019.
- [NMS20] Hadis Najafimehr, Kosar Mohamed Ali, Saeed Safari, Mahmoud Yousefifard, and Mostafa Hosseini. “Estimation of basic reproduction number for COVID-19 and the reasons for its differences.” *International Journal of Clinical Practice*, 2020.
- [NZA16] D. Nie, H. Zhang, E. Adeli, L. Liu, and D. Shen. “3D deep learning for multi-modal imaging-guided survival time prediction of brain tumor patients.” *Med. Image Comput. Comput. Assist. Interv.*, **9901**, 2016.
- [OBL14] Maxime Oquab, Leon Bottou, Ivan Laptev, and Josef Sivic. “Learning and transferring mid-level image representations using convolutional neural networks.” In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1717–1724, 2014.

- [OPL] Sang Ho Oh, Jongyoul Park, Su Jin Lee, Seungyeon Kang, and Jeonghoon Mo. “Reinforcement Learning-Based Expanded Personalized Diabetes Treatment Recommendation Using South Korean Electronic Health Records.” **206**:117932.
- [Oz18] Mustafa Oz. *Discussing Controversial Issues on Social Media: Examining the Role of Affordances, Fear of Isolation and De-Individuation*. PhD thesis, The University of Texas at Austin, Ann Arbor, United States, 2018.
- [PCC] Niranjani Prasad, Li-Fang Cheng, Corey Chivers, Michael Draugelis, and Barbara E. Engelhardt. “A Reinforcement Learning Approach to Weaning of Mechanical Ventilation in Intensive Care Units.”.
- [PDW] Xuefeng Peng, Yi Ding, David Wihl, Omer Gottesman, Matthieu Komorowski, Li-wei H. Lehman, Andrew Ross, Aldo Faisal, and Finale Doshi-Velez. “Improving Sepsis Treatment Strategies by Combining Deep and Kernel-Based Reinforcement Learning.” **2018**:887–896.
- [Pfe18] A. Pfefferbaum. “Altered brain developmental trajectories in adolescents after initiating drinking.” *Am. J. Psychiatry*, **175**, 2018.
- [PGM19] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. “Pytorch: An imperative style, high-performance deep learning library.” *Advances in neural information processing systems*, **32**, 2019.
- [PLR20] Kiesha Prem, Yang Liu, Timothy W Russell, Adam J Kucharski, Rosalind M Eggo, Nicholas Davies, Stefan Flasche, Samuel Clifford, Carl AB Pearson, James D Munday, et al. “The effect of control strategies to reduce social mixing on outcomes of the COVID-19 epidemic in Wuhan, China: a modelling study.” *The Lancet Public Health*, 2020.
- [PY10] S. J. Pan and Q. Yang. “A survey on transfer learning.” *IEEE Trans. Knowl. Data Eng.*, **22**, 2010.
- [Qi 16] N. Qi Dou. “Automatic detection of cerebral microbleeds from MR images via 3D convolutional neural networks.” *IEEE Trans. Med. Imaging*, **35**, 2016.
- [Rab22] N. Rabb. “The influence of social norms varies with “others” groups: evidence from COVID-19 vaccination intentions.” *Proc Natl Acad Sci USA*, **119**, 2022.
- [Raj22] K. Rajkumar. “A causal test of the strength of weak ties.” *Science*, **377**, 2022.
- [Rav17] D. Ravi. “Deep learning for health informatics.” *IEEE J. Biomed. Health Inf.*, **21**, 2017.



- [RKA] Aniruddh Raghu, Matthieu Komorowski, Imran Ahmed, Leo Celi, Peter Szolovits, and Marzyeh Ghassemi. “Deep Reinforcement Learning for Sepsis Treatment.”
- [RLO19] D. B. Russakoff, A. Lamin, J. D. Oakley, A. M. Dubis, and S. Sivaprasad. “Deep learning for prediction of AMD progression: a pilot study.” *Investigative Ophthalmol. Vis. Sci.*, **60**, 2019.
- [Rob] Miguel A. Hernan Robins, James M. *Causal Inference: What If*. CRC Press.
- [Rot16] H. R. Roth. “Improving computer-aided detection using convolutional neural networks and random view aggregation.” *IEEE Trans. Med. Imaging*, **35**, 2016.
- [RR] PAUL R. ROSENBAUM and DONALD B. RUBIN. “The Central Role of the Propensity Score in Observational Studies for Causal Effects.” **70**(1):41–55.
- [SA15] Bogdan State and Lada Adamic. “The Diffusion of Support in an Online Social Movement: Evidence from the Adoption of Equal-Sign Profile Pictures.” In *Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing, CSCW ’15*, pp. 1741–1750, New York, NY, USA, February 2015. Association for Computing Machinery.
- [SA23] S. Sinclair and J. Agerström. “Do social norms influence young people’s willingness to take the COVID-19 vaccine?” *Health Commun.*, **38**, 2023.
- [Sch18] A. L. Schmidt. “Polarization of the vaccination debate on Facebook.” *Vaccine*, **36**, 2018.
- [She09] J. Sherman. “Photoreceptor integrity line joins the nerve fiber layer as key to clinical diagnosis.” *Optom. - J. Am. Optometric Assoc.*, **80**, 2009.
- [Shi16] H.-.. C. Shin. “Deep convolutional neural networks for computer-aided detection: CNN architectures, dataset characteristics and transfer learning.” *IEEE Trans. Med. Imaging*, **35**, 2016.
- [Sol21] J. S. Solís Arce. “COVID-19 vaccine acceptance and hesitancy in low- and middle-income countries.” *Nat Med*, **27**, 2021.
- [SP10] Skipper Seabold and Josef Perktold. “Statsmodels: Econometric and Statistical Modeling with Python.”, 2010.
- [ST20] T. Sun and S. J. Taylor. “Displaying things in common to encourage friendship formation: a large randomized field experiment.” *Quant Mark Econ*, **18**, 2020.
- [Taj20] N. Tajbakhsh. “Embracing imperfect datasets: a review of deep learning solutions for medical image segmentation.” *Med. Image Anal.*, **63**, 2020.

- [The20] The New York Times. “Coronavirus in the U.S.: Latest Map and Case Count.” *The New York Times*, March 2020.
- [TLC] Huan-Hsin Tseng, Yi Luo, Sunan Cui, Jen-Tzung Chien, Randall K. Ten Haken, and Issam El Naqa. “Deep Reinforcement Learning for Automated Radiation Adaptation in Lung Cancer.” **44**(12):6690–6705.
- [TMS] Shengpu Tang, Maggie Makar, Michael Sjoding, Finale Doshi-Velez, and Jenna Wiens. “Leveraging Factored Action Spaces for Efficient Offline Reinforcement Learning in Healthcare.” **35**:34272–34286.
- [TPF18] Anna CS Tan, Matthew G Pilgrim, Sarah Fearn, Sergio Bertazzo, Elena Tsolaki, Alexander P Morrell, Miaoling Li, Jeffrey D Messinger, Rosa Dolz-Marco, Jianqin Lei, et al. “Calcified nodules in retinal drusen are associated with disease progression in age-related macular degeneration.” *Science translational medicine*, **10**(466):eaat4544, 2018.
- [TRK18] P. Tschandl, C. Rosendahl, and H. Kittler. “The HAM10000 dataset, a large collection of multi-source dermatoscopic images of common pigmented skin lesions.” *Sci. Data*, **5**, 2018.
- [Val17] S. Valverde. “Improving Automated Multiple Sclerosis Lesion Segmentation with a Cascaded 3D Convolutional Neural Network Approach.” *NeuroImage*, **155**, 2017.
- [VSP] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. “Attention Is All You Need.” In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, pp. 6000–6010. Curran Associates Inc.
- [Wen] John E Wennberg. “Unwarranted Variations in Healthcare Delivery: Implications for Academic Medical Centres.” **325**(7370):961–964.
- [Wil45] F. Wilcoxon. “Individual comparisons by ranking methods.” *Biometrics Bull.*, **1**, 1945.
- [WKV] Jeremy Watts, Anahita Khojandi, Rama Vasudevan, and Ritesh Ramdhani. “Optimizing Individualized Treatment Planning for Parkinson’s Disease Using Deep Reinforcement Learning.” In *2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pp. 5406–5409.
- [WKW16] K. Weiss, T. M. Khoshgoftaar, and D. Wang. “A survey of transfer learning.” *J. Big Data*, **3**, 2016.

- [WSD] Paul J. Wallace, Nilay D. Shah, Taylor Dennen, Paul A. Bleicher, and William H. Crown. “Optum Labs: Building A Novel Node In The Learning Health Care System.” **33**(7):1187–1194.
- [YCN15] Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson. “Understanding neural networks through deep visualization.” *arXiv preprint arXiv:1506.06579*, 2015.
- [YF14] Elad Yom-Tov and Luis Fernandez-Luque. “Information is in the eye of the beholder: Seeking information on the MMR vaccine through an Internet search engine.” *AMIA Annu. Symp. Proc.*, **2014**:1238–1247, November 2014.
- [YS] Gregory Yauney and Pratik Shah. “Reinforcement Learning with Action-Derived Rewards for Chemotherapy and Clinical Trial Dosing Regimen Selection.” In Finale Doshi-Velez, Jim Fackler, Ken Jung, David Kale, Rajesh Ranganath, Byron Wallace, and Jenna Wiens, editors, *Proceedings of the 3rd Machine Learning for Healthcare Conference*, volume 85 of *Proceedings of Machine Learning Research*, pp. 161–226. PMLR.
- [ZG] Siyan Zhao and Aditya Grover. “Decision Stacks: Flexible Reinforcement Learning via Modular Generative Models.”
- [ZZS] Yufan Zhao, Donglin Zeng, Mark A. Socinski, and Michael R. Kosorok. “Reinforcement Learning Strategies for Clinical Trials in Non-small Cell Lung Cancer.” **67**(4):1422–1433.