

# Change in machine learning model performance upon retraining after deployment into clinical practice: The real-world effect of model predictions on clinician actions, outcome labels, and the potential for contamination bias

## BACKGROUND

CHARTWatch is an AI model that was implemented at St. Michael's Hospital in October 2020 and predicts inpatient deterioration on medical wards.<sup>1</sup>

Al Models may degrade over time due to factors like data drift and may require retraining. However, retraining a deployed model using post-deployment data may worsen performance due to **contamination bias**—a phenomenon where the model changes outcomes it later uses to retrain, as shown in recent simulation studies.<sup>2-5</sup>

To inform whether to retrain CHARTwatch, we sought to quantify the degree of contamination bias and explore strategies to mitigate it.

### **METHODS**

#### Part 1: Prospective cohort study

- Evaluate how often CHARTwatch prevents patient deterioration (effective intervention rate) through realtime surveys and chart review for 100 consecutive alerts
- Quantify the overall outcome frequency and model recall



### **Part 2: Model retraining evaluations**

- Evaluate the change in model performance due to contamination bias when retraining with a range of effective intervention rates, recall frequencies, and ML model types, with values informed by the prospective cohort study
- Contamination bias mitigation: Evaluate the effect of removing potentially confounded outcomes prior to retraining on the magnitude of contamination bias and overall model performance

**Michael Colacci MD PhD<sup>1-3</sup>**, George-Alexandru Adam PhD<sup>4</sup>, Chloe Pou-Prom MSc<sup>1</sup>, Anna Goldenberg PhD<sup>4</sup>, Amol Verma MD, MPhil<sup>1-3</sup>, Muhammad Mamdani PharmD, MA, MPH<sup>1-3</sup>

[1] Li Ka Shing Knowledge Institute, St. Michael's Hospital, Unity Health Toronto, Toronto, Canada, [2] Department of Medicine, University of Toronto, Toronto, Canada [3] Institute of Health Policy, Management and Evaluation, Dalla Lana School of Public Health, University of Toronto, Toronto, Canada [4] Department of Computer Science, University of Toronto, Toronto, Canada







For effective intervention rate, the 3 vertical panels correspond to different probabilities of the patient surviving due to a CHARTwatch alert (10%, 30% and 50%) with a fixed recall of 70%. The top left plot (blue square, 0.1 EIR) is most representative of observed CHARTWatch performance. For model recall, the 3 vertical panels correspond to different model recall rates (25%, 50% and 75%), with a fixed 35% effective intervention rate.







# DISCUSSION

Clinician agreement with CHARTwatch predictions is associated with downstream clinical actions.

We present a framework for using clinician perception/actions and model parameters to estimate contamination bias.

For CHARTwatch, contamination bias at any retraining interval is limited ( $\triangle AUC < 2\%$ ) but the impact over time is summative.

Removing potentially confounded outcomes may help mitigate contamination bias during retraining.

#### **REFERENCES:**

Deployed

Not Deployed

1. Verma AA, Stukel TA, Colacci M, et al. Clinical evaluation of a machine learning-based early warning system for patient deterioration. CMAJ. 2024;196(30):E1027-E1037. doi:10.1503/cmaj.240132 2. Adam GA, Chang CHK, Haibe-Kains B, Goldenberg A. Hidden Risks of Machine Learning Applied to Healthcare: Unintended Feedback Loops Between Models and Future Data Causing Model Degradation. In: Proceedings of the 5th Machine Learning for Healthcare Conference; 2020. 3. Adam GA, Chang CHK, Haibe-Kains B, Goldenberg A. Error Amplification When Updating Deployed Machine Learning Models. In: Proceedings of Machine Learning Research. Vol 182.; 2022. 4. Vaid A, Sawant A, Suarez-Farinas M, et al. Implications of the Use of Artificial Intelligence Predictive Models in Health Care Settings : A Simulation Study. Ann Intern Med. 2023;176(10):1358-1369. doi:10.7326/M23-0949 5. Finlayson SG, Subbaswamy A, Singh K, et al. The Clinician and Dataset Shift in Artificial Intelligence. New England Journal of Medicine. 2021;385(3). doi:10.1056/nejmc2104626

6. Lenert MC, Matheny ME, Walsh CG. Prognostic models will be victims of their own success, unless... Journal of the American Medical Informatics Association 2019;26(12). doi:10.1093/jamia/ocz145













