

# Diagnosing and remediating harmful data shifts for the responsible deployment of clinical AI models

Vallijah Subasri<sup>1,2,3</sup>, Amrit Krishnan<sup>3</sup>, Azra Dhalla<sup>3</sup>, Deval Pandya<sup>3</sup>, David Malkin<sup>1,2,7,8</sup>, Fahad Razak<sup>4,5,6</sup>, Amol A. Verma<sup>4,5,6</sup>, Anna Goldenberg<sup>1,3,9</sup>, Elham Dolatabadi<sup>3,5</sup>

- <sup>1</sup>. Genetics and Genome Biology Program, The Hospital for Sick Children, Toronto, ON, Canada
- <sup>2</sup>. Department of Medical Biophysics, University of Toronto, Toronto, ON, Canada
- <sup>3</sup>. Vector Institute for Artificial Intelligence, Toronto, ON, Canada
- <sup>4</sup>. St Michael's Hospital, Unity Health Toronto, Toronto, ON, Canada
- <sup>5</sup>. Institute of Health Policy, Management and Evaluation (IHPME), University of Toronto, Toronto, ON, Canada
- <sup>6</sup>. Department of Medicine, University of Toronto, Toronto, ON, Canada
- <sup>7</sup>. Division of Hematology/Oncology, The Hospital for Sick Children, Toronto, ON, Canada
- <sup>8</sup>. Department of Pediatrics, University of Toronto, Toronto, ON, Canada
- <sup>9</sup>. Department of Computer Science, University of Toronto, Toronto, ON, Canada

## Abstract

Harmful data shifts occur when the distribution of data used to train a clinical AI system differs significantly from the distribution of data encountered during deployment, leading to erroneous predictions and potential harm to patients. We evaluated the impact of data shifts on an early warning system (EWS) for in-hospital mortality that uses electronic health record (EHR) data from patients admitted to a general internal medicine service. We found model performance to differ across subgroups of clinical diagnoses, sex and age. To explore the robustness of the model, we evaluated potentially harmful data shifts across demographics, hospital types, seasons, times of hospital admission, and whether the patient was admitted from an acute care institution or nursing home, without relying on model performance. Interestingly, we found that models trained on community hospitals experience harmful data shifts when evaluated on academic hospitals, whereas the models trained on academic hospitals transfer well to the community hospitals. To improve model performance across hospital sites we employed transfer learning, a strategy that stores knowledge gained from learning one domain and applies it to a different but related domain. We found hospital type-specific models that leverage transfer learning, perform better than models that use all available hospitals. Furthermore, we monitored data shifts over time and identified model deterioration during the COVID-19 pandemic. Typically machine learning models remain locked after deployment, however, this can lead to model deterioration due to data shifts that occur over time. We used continual learning, the process of learning from a continual stream of data in a sequential manner, to mitigate data shifts over time and improve model performance. Overall, our study is a crucial step towards the deployment of clinical AI models, by providing strategies and workflows to ensure the safety and efficacy of these models in real-world settings.

## Introduction

AI systems have leveraged clinical data to predict mortality<sup>1–5</sup>, length of stay (LOS)<sup>6</sup>, sepsis<sup>7–9</sup> and the occurrence of specific disease diagnoses<sup>10</sup>. As a growing number of AI systems are sought to be deployed in clinical settings, a defining challenge for AI in healthcare is how to responsibly deploy models that have been developed<sup>11,12</sup>. Building robust clinical machine learning (ML) models has proven to be difficult<sup>13</sup>, in part attributed to *data shifts* (or *data drift*)—changes in the data distribution over time and/or space that leads to spurious predictions<sup>14</sup>. This can occur due to changes in the *features* of the input data or due to changes in the *labels*, which represent the outcome the model is predicting. Data shifts are harmful when they result in *model drift*—a significant decrease in the model’s predictive power due to changes in the real world environment. A key barrier to the safe deployment of clinical AI systems is attributed to system malfunction due to harmful data shifts<sup>15</sup>. Data shifts occur when the underlying distribution of the data used to build a predictive model differs from the distribution of the data encountered during deployment. In healthcare, these shifts can exist along the axes of institutional differences (e.g., staffing, instruments and data-collection workflows), epidemiological changes (e.g. diseases, catastrophic events)<sup>16</sup>, temporal shifts (e.g. policy changes, changes in clinician or patient behaviours over time)<sup>17</sup> and differences in patient demographics (e.g. race, sex, age, socioeconomic background, and types of presenting illnesses and comorbidities)<sup>18–20</sup>. When the difference between the training and test data distribution is sufficient to deteriorate the model’s performance, clinical decision-making may be impaired. As a result, it is imperative to identify these potentially harmful shifts a priori, to inform clinical end-users and prevent harm to patients.

Rigorous evaluations across time, hospital sites, and patient characteristics are critical for identifying model degradation and ensuring equitable and quality patient care. The impact of distributional shifts on model performance<sup>21</sup> has been explored for the prediction of sepsis<sup>22</sup>, mortality<sup>19,23</sup>, ER admissions<sup>16</sup>, LOS<sup>19</sup> and *Clostridioides difficile* infections<sup>17</sup>. Model deterioration has previously been associated with transitions in EHRs systems over time<sup>13</sup> and across patient demographics in chest X-rays<sup>24</sup>, skin lesions<sup>25</sup> and sepsis prediction<sup>26</sup>. However, in many clinical prediction problems, the lead time to acquire labels is lengthy, and the process is resource-intensive. Labels like death or sepsis are rare; this causes a delay in the ability to detect a statistically significant change in model performance, at which point model deterioration may have already occurred, and it may be too late to take steps for remediation. This suggests retraining based on recognizing deterioration in model performance is impractical, and emphasizes the importance of detecting potentially harmful data shifts in a label-agnostic manner<sup>27–29</sup>. Furthermore, it is necessary to design effective strategies for model updating that proactively minimize model degradation in the presence of data shifts. Failure to correct for harmful data shifts can lead to the perpetuation of algorithmic biases, missing critical diagnoses and unnecessary clinical interventions that can be detrimental to patient outcomes and burden the healthcare system<sup>11,12</sup>.

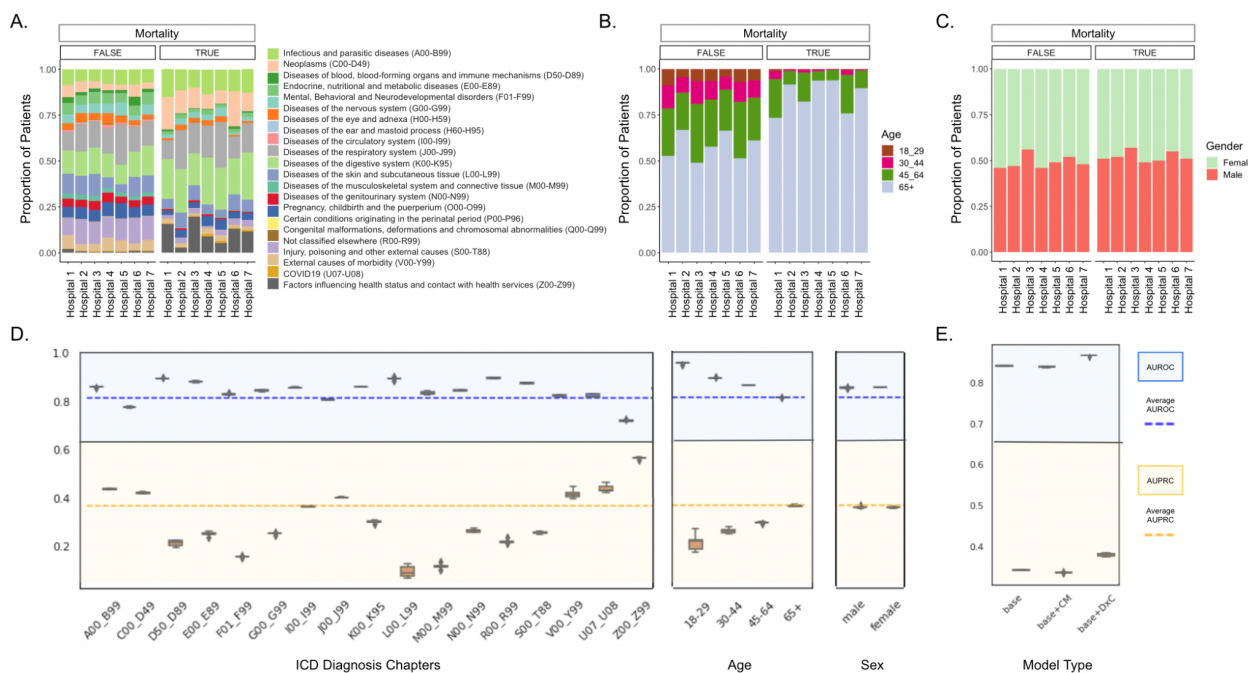
In this study, we developed an evaluation and monitoring pipeline to prepare clinical AI systems for deployment.<sup>30</sup> We used our pipeline to monitor for harmful data shifts in a label-agnostic manner using an early warning system (EWS) for all-cause in-hospital mortality. In doing so, we proactively identified harmful data shifts across various real-life scenarios, including institutional differences, time of hospital admission, whether a patient was admitted from an acute care institution or nursing home and the COVID-19 pandemic. In the presence of harmful data shifts across institutions, we leveraged transfer learning to identify strategies for improving model performance<sup>31–33</sup>. Lastly, we conducted a prospective evaluation, whereby we monitored for temporal data shifts and used continual learning to proactively update clinical AI models under harmful data shifts.

## Results

### All-cause in-hospital mortality early warning system (EWS)

We developed a dynamic EWS to predict the risk of in-hospital mortality within the next two weeks, every 24 hours, using EHR data consisting of lab results, transfusions, imaging reports and administrative

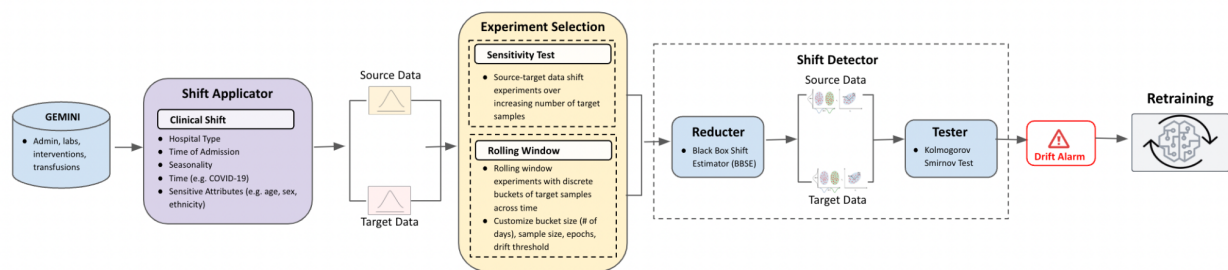
features (**Supplementary Table 1**) from 109,802 patient encounters admitted to general internal medicine (GIM) inpatient units at seven large hospitals in the greater Toronto area (GTA, in Canada). Given the varying distribution of diagnoses and demographics across hospitals (**Figure 1A-C**), we assessed the fairness of our model by evaluating the area under the receiver operating characteristic (AUROC) and area under the precision-recall curve (AUPRC) for subgroups of diagnoses, sex and age<sup>34</sup>. We defined diagnostic subgroups using the ICD-10 diagnosis chapters<sup>35</sup>—groupings of ICD-10 diagnosis codes assigned to patients during admission based on affected body systems and health conditions. We found that the model performed particularly well on certain diagnoses, including diseases of the circulatory system (I00-I99), respiratory system (J00-J99), COVID-19 (U07-U08) and certain other infectious and parasitic diseases (A00-B99). However, it had a much lower AUROC on individuals with benign or malignant neoplasms (C00-D49) and factors influencing health status and contact with health services (Z00\_Z99; **Figure 1D**). These primarily consisted of patients receiving palliative care ( $n_{Z515}=2042$ ; **Supplementary Figure 1**), including patients with cancer, heart failure, chronic obstructive pulmonary disease (COPD), dementia, and Parkinson disease. This is in accordance with what we know about palliative care as encompassing complex diseases with evolving needs, caused by a combination of genetic, environmental and lifestyle factors, which may make it more difficult to accurately predict<sup>36</sup>. We also found AUROC increased and AUPRC decreased across groups with decreasing age, this may be in part driven by the lower mortality rates in the younger age groups. Alternatively, performance was fairly consistent across sex (**Figure 1D**). Lastly, we compared the performance of our model, which included no prior information of patient history, to models that included comorbidities and ICD-10 diagnosis codes as features (**Supplementary Table 1**). In doing so, we found that including ICD-10 diagnosis codes as features in our model slightly improved overall performance (**Figure 1E**), but significantly increased the performance gap between many diagnostic subgroups (**Supplementary Figure 2**).



**Figure 1. Model fairness across subgroups. (A)** Distribution of ICD-10 diagnosis codes across hospitals by mortality status (true/false). **(B)** Distribution of age groups across hospitals by mortality status (true/false). **(C)** Distribution of sex across hospitals by mortality status (true/false). **(D)** AUROC and AUPRC of cross-site EWS across ICD-10 diagnosis codes, age and sex. **(E)** Overall AUROC and AUPRC of the model without prior information (base), with comorbidities (base+CM) and with diagnosis codes (base+DxC) as features.

Mortality	False							True						
	Hospital	1(A)	2(A)	3(A)	4(C)	5(C)	6(A)	7(A)	1(A)	2(A)	3(A)	4(C)	5(C)	6(A)
# of Encounters	16620	27394	16100	15096	22691	16238	14072	2405	2306	1589	1415	3524	1791	1808
LOS (days)	8.47	8.44	7.90	8.73	10.05	7.55	9.03	14.59	17.19	17.62	19.02	17.72	12.80	14.12
# of Prev Encounters from 2010-2020	0.72	0.60	1.19	0.68	0.60	0.82	0.90	1.28	1.12	1.53	1.36	1.16	1.31	1.63
From Acute Care (%)	2	0	1	1	0	1	1	7	1	1	3	0	1	0
From Nursing Home (%)	6	8	4	11	12	4	10	13	18	14	31	30	8	22
Palliative Care (n)	260	151	50	52	39	160	80	373	56	313	123	191	235	210

**Table 1. Patient characteristics.** Number of patient encounters, average length of stay (LOS), average number of previous encounters from 2010-2020, percentage of patient encounters from acute care institutions, percentage of patient encounters from nursing homes, and number of patients receiving palliative care, across hospitals and mortality status. A = academic hospital, C = community hospital.



**Figure 2. Monitoring and evaluation pipeline.** An end-to-end pipeline where EHR data is first sent to the *Shift Applicator*, which outputs a source and target data based on the clinical shift of choice. The source and target data can then be leveraged by the *Shift Detector* to conduct a drift sensitivity test or a rolling window analysis—if drift is detected, retraining is triggered.

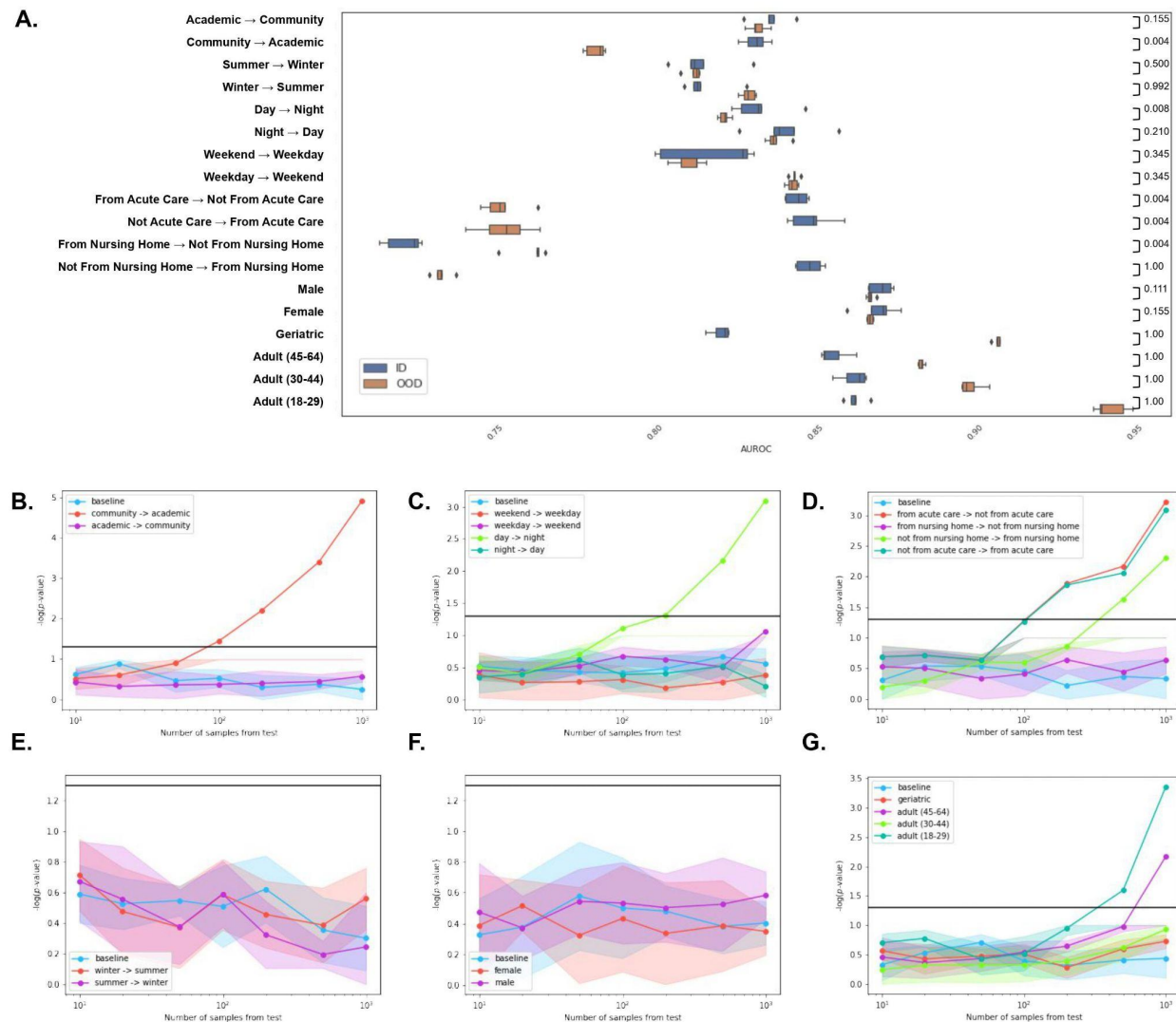
## Detection of harmful data shifts for evaluation and monitoring of clinical AI systems

In the clinical setting, there are a myriad of factors that can contribute to a model drifting and making erroneous predictions, such as changes in behaviour, technology, population or policy<sup>18</sup>. Using our monitoring and evaluation pipeline (**Figure 2**)<sup>37</sup>, we detected data shifts in a label-agnostic manner across increasing sample sizes for scenarios that we would expect to pose a threat to clinical AI systems during deployment, due to fundamental differences in patient populations. These scenarios consist of differences in demographics, hospital type, seasonality, time of day of hospital admission (i.e. day vs. night), time of week of hospital admission (i.e. weekday vs. weekend), and whether patients were admitted from an acute care institution or nursing home. Harmful data shifts were defined as those statistically significant between the source and target datasets ( $p\text{-value} < 0.05$ ). We detected harmful data shifts and associated performance degradation in five scenarios: when transferring models trained on i) community hospitals to academic hospitals (**Figure 3AB**), ii) patients admitted during the day to patients admitted at night (**Figure 3AC**), iii) patients not admitted from nursing homes to patients admitted from nursing homes (**Figure 3AD**), iv) patients admitted from acute care institutions to patients admitted from non-acute care institutions (**Figure 3AD**) and v) patients admitted from non-acute care institutions to patients admitted from acute care institutions (**Figure 3AD**). Interestingly, we found many of these harmful data shifts were unidirectional, suggesting that there exists patterns among patient encounters in academic hospitals, during night admissions and among patients admitted from nursing homes that are not captured at

community hospitals, during day admissions, and among patients admitted from outside of nursing homes, respectively. Harmful data shifts were not detected across seasons or sex (**Figure 3EF**). Although a harmful data shift was identified when evaluating on the 45-64 year-old age group, an associated decrease in AUROC did not occur (**Figure 3AG**).

These data shifts can arise for a variety of reasons, including differences in patient subpopulations, staffing, and/or resources that are not adequately represented in the training data<sup>38</sup>. Across all the scenarios where harmful data shifts were identified, we found that there was decreased performance in numerous diagnostic subgroups between the source and target data (**Supplementary Figure 3**). The largest performance differences between patients from acute care and non-acute care institutions was for diseases of the nervous system (G00-G99), and musculoskeletal system and connective tissue (M00-M99). Between patients admitted during the day and night, the largest decrease in AUROC was seen in patients with neoplasms (C00-D49) and diseases of the musculoskeletal system and connective tissue (M00-M99) and genitourinary system (N00-N99). When transferring from community hospitals to academic hospitals, the largest performance decrease across diagnostic subgroups was for patients with neoplasms (C00-D49), which is also found at a much higher prevalence in academic hospitals (**Supplementary Table 2**). Alternatively, the hospital type shift may be due to differences in the 45-64 year-old age group, which suffered a significant decrease in performance when models were transferred from community hospitals to academic hospitals ( $p=0.0079$ ; **Supplementary Figure 3**). This could in part be driven by the increased number of individuals admitted from nursing homes in community hospitals compared to academic hospitals (**Table 1**). This is also supported by our finding that models transferred from patients not admitted to nursing homes to patients admitted to nursing homes—which primarily consist of long-term care residents over the age of 85, result in harmful data shifts (**Figure 3AD**). It is also worth noting that the hospital type groupings coincide with differences in location which may also be a contributing factor of the data shift; more specifically, the academic hospitals are located in the central city while the community hospitals are located in residential suburbs. Interestingly, we found the inclusion of ICD-10 diagnosis codes as features decreased model deterioration due to data shifts (**Supplementary Table 3**).



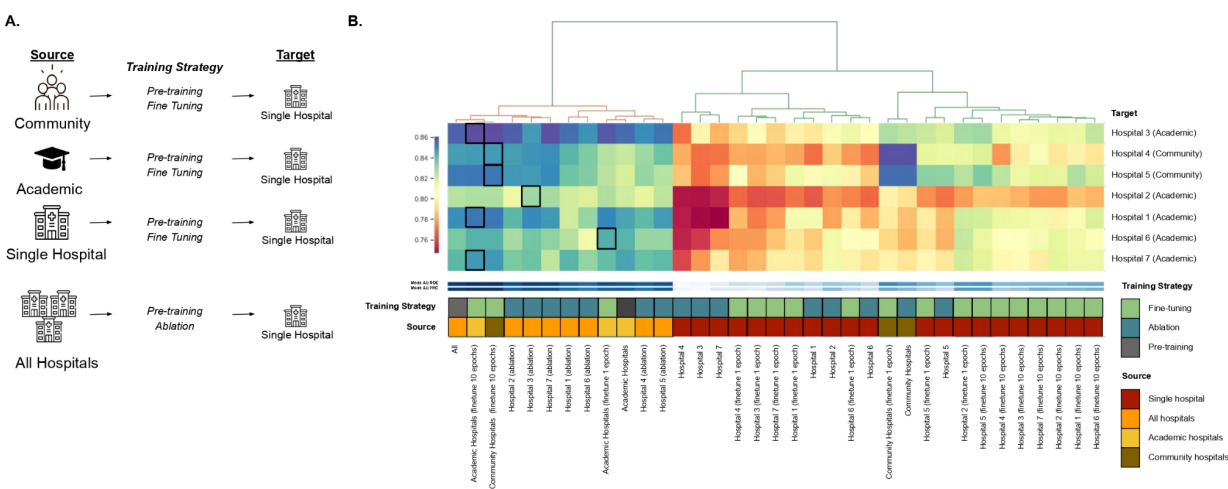


**Figure 3. Detection of harmful data shifts.** (A) Performance across in-distribution (ID) and out-of-distribution (OOD) data for demographics, hospital types, seasons, time of day of hospital admission, time of week of hospital admission, and whether the patient was admitted from an acute care institution or nursing home. P-values were calculated using a one-sided Mann-Whitney U test. Sensitivity of the data shift detection to increasing number of test samples was evaluated for (B) hospital types (C) time of day or week of hospital admission (D) seasons (E) admission from acute care institutions or nursing homes (F) sex and (G) age.

### Preventing harmful data shifts during cross-site deployment

It is common practice that an ML model is developed at one institution and transferred to other institutions for external validation. During cross-site evaluation, we found that differences in hospital type result in harmful data shifts that deteriorate model performance (Figure 3B). In order to address this, we developed EWSs for i) each individual hospital, ii) the combination of community hospitals, iii) the combination of academic hospitals and iv) the combination of all hospitals. We then compared strategies leveraging a) *pre-training* where we used a model pre-trained on source data and evaluated it on out-of-distribution data from the target hospital, b) *transfer learning* where we fine-tuned the performance on the target hospital prior to evaluating the target data and c) *ablation* where we excluded data from a single hospital prior to evaluating the target data. For each model, we evaluated the performance for each

individual hospital using a held out test set (**Figure 4A**). In general, cross-site training improved model performance; however, the use of all sites was never the optimal strategy suggesting that more data is not always helpful. We found training across all sites marginally improved model performance for academic hospitals but decreased performance for community hospitals (**Figure 4B**). Instead, using the model trained on both the community hospitals (Hospital 4 and 5) resulted in superior performance for community hospitals. Overall, fine-tuning the corresponding hospital type-specific model on the target hospital improved performance for all hospitals except Hospital 2. Interestingly, Hospital 2 is also the only hospital with a veteran's wing, where patients receiving palliative care were less likely to experience in-hospital mortality and where the number of previous hospital visits was negatively correlated with risk of in-hospital mortality (**Table 1**). In certain instances, the exclusion of a single hospital site improved model performance for another hospital. For Hospital 2, ablating Hospital 3 resulted in the best performing model (**Figure 4B**). It is worth noting, Hospital 2 and Hospital 3 also had the largest difference in the number of individuals with diseases due to factors influencing health status and contact with health services (17%), which is the diagnostic subgroup with the lowest performance (**Supplementary Table 3; Figure 1E**). These two hospitals also had the largest difference in patients receiving palliative care between mortality status; Hospital 2 had a 2.7-fold decrease and Hospital 3 had a 6.3-fold increase in palliative care among patients who died in the hospital. The population demographic and socioeconomic status (SES) between the two hospitals are also very different; Hospital 3 is an inner city urban and Hospital 2 is a suburban hospital. As a result, it is important that clinical AI systems be proactively evaluated for these differences so they are considered when transferring models across sites.



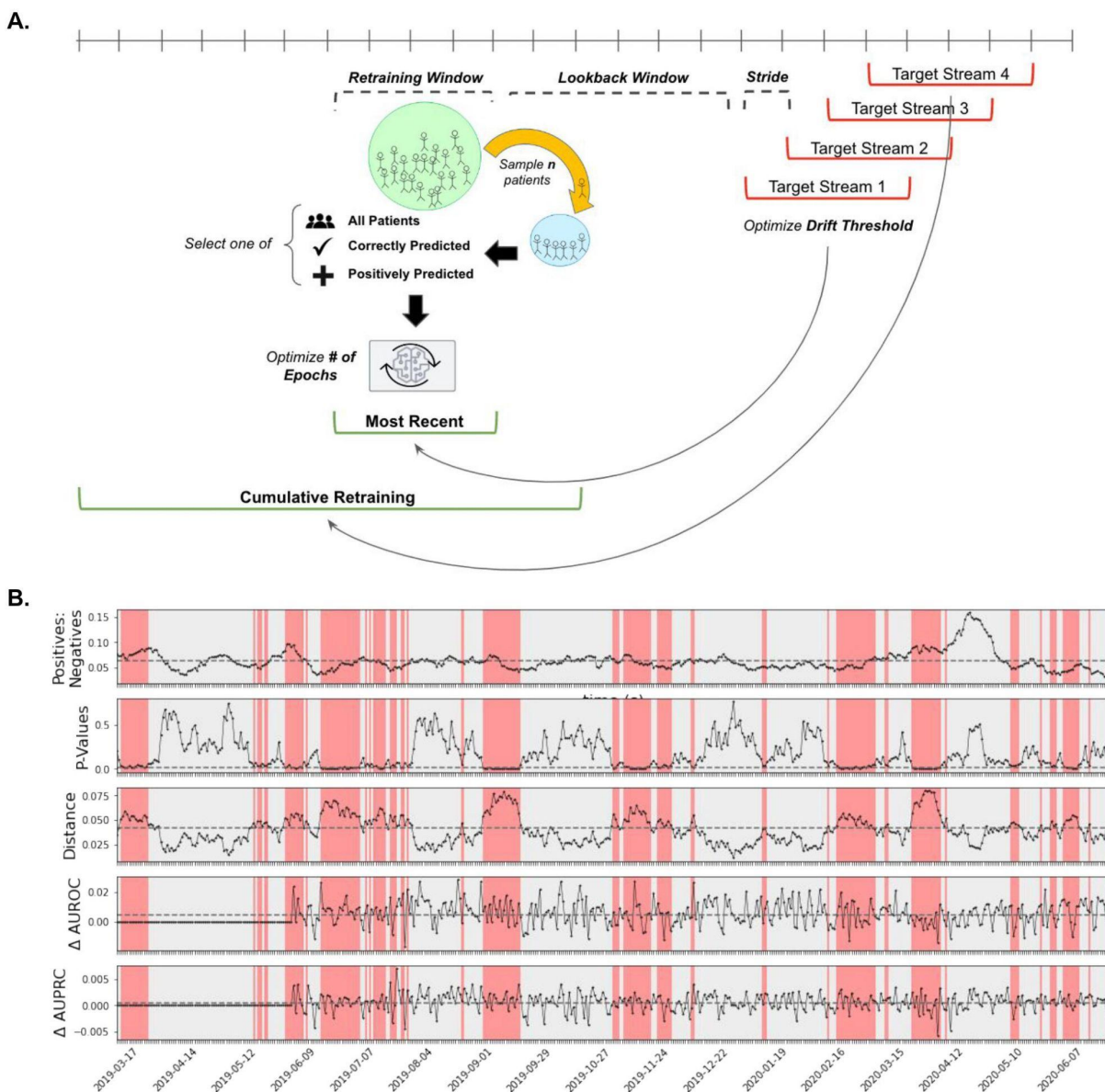
**Figure 4. Optimal training strategies for cross-site deployment. (A)** Pre-training, fine-tuning and ablation employed on single-site models, cross-site model and hospital-type models (community, academic). **(B)** Heatmap of AUROC of the training strategies across each test hospital site. Highlighted in black is the best performing Hospital model for each hospital.

### Detecting and mitigating model deterioration due to temporal data shifts

Lastly, we conducted a simulated prospective evaluation of an EWS for mortality prediction using GIM data from 2011-2018. In a real-time deployment scenario, labels are not always readily available at the time of prediction. Moreover, for outcomes like mortality, the problem with relying on model performance is that the event rate is relatively rare, so it can take many months to accrue a sufficient sample size for detecting model performance changes. As a result, label-agnostic drift detection is critical for identifying model degradation and triggering retraining procedures. We monitored our EWS for temporal shifts using a 14-day rolling window from March 2019 to August 2020. In the presence of drift, we used continual learning strategies to update our model and mitigate model deterioration (**Figure 5A**). First, we compared *periodic retraining*—whereby the model is updated at regular, pre-defined intervals and *drift-triggered retraining*—whereby the model is updated when there is significant data shift between the source data and

target data. We found drift-triggered retraining resulted in better overall performance (**Supplementary Figure 4**). To identify the optimal approach for drift-triggered retraining, we tuned various parameters including the retraining window size, lookahead window, sample size, drift threshold and number of epochs (**Supplementary Figures 5-9**). The retraining window represents how much previous data we want to use for updating the model. We found a larger retraining window improved AUROC and AUPRC, however, as the retraining window increased upwards of 180 days, the performance decreased, suggesting that greater amounts of past data are not always beneficial for model updating (**Supplementary Figure 5**). Due to the lead time for acquiring labels, it is possible that at the time model updating is triggered, labels for the most recent patient encounters are not available. As a result, we evaluated increasing lookback window sizes to determine how far back the data used to update the model can be, without sacrificing performance. We found lookback windows of up to 60 days were able to maintain similar model performance (**Supplementary Figure 6**). Although, the lookback window will differ depending on the frequency of the prediction outcome and the progression of the drift over time (i.e. gradual versus sudden). Given that the model updating is triggered by drift detection, the sensitivity of the drift test will influence the overall performance. We found that the optimal drift threshold was a p-value of 0.01 (**Supplementary Figure 7**) and the optimal number of encounters for the drift test was 1000 (**Supplementary Figure 8**). However, it is important to recognize each prediction task and domain is unique, and as a result the generalizability of the optimal threshold will need to be evaluated on a case-by-case basis. We also found that increasing the number of epochs during model updating resulted in catastrophic forgetting whereby the model overfit and model performance decreased over time (**Supplementary Figure 9**). We also compared updating whereby we only trained on encounters that were predicted correctly or positively; however, this was not as effective as using all the encounters (**Supplementary Figure 10**). Overall, the implementation of our drift-triggered continual updating strategy improved model performance over time and was more effective than maintaining a locked model during deployment (**Figure 5B**).





**Figure 5. Prospective evaluation of EWS for mortality prediction. A.** Model updating strategies and parameters explored whereby the target stream is evaluated on i) the most recent  $n$  days or ii) all the encounters seen to date, with the option to select training on only the encounters predicted positively or correctly. Each strategy can be optimized for drift threshold, number of samples, number of epochs, stride length, lookback window, and retraining window. **B.** We monitored the proportion of positive:negative outcomes, drift p-values, and drift distance metric from 03/2019 to 08/2020 using a 14-day rolling window. In the event that drift is detected, model updating is triggered (red), for which we also monitored the change in AUROC and AUPRC between the retrained model and the baseline model, which implements no updating procedure.

## Discussion

Many widely implemented clinical AI systems<sup>26,39,40</sup>, have demonstrated poor generalizability upon external validation, as a result of harmful data shifts. However, these biases are rarely accounted for in a proactive manner, and are typically identified following deployment, while relying on ground-truth

labels<sup>18,41,42</sup>. In this study, we built a dynamic EWS that adapts to the ever-changing healthcare environment. We used our EWS to predict the risk of mortality to enable the effective triaging of patients admitted to GIM and performed robust evaluations for bias and data shifts across diagnostic subgroups, demographics, hospital sites, based on the when and where a patient was admitted, and over time. We accurately detected harmful data shifts in clinical data without relying on ground-truth labels by leveraging black box shift detection and two sample testing<sup>28</sup>; this permitted the proactive evaluation of ML models in clinical settings where labels can be costly, resource-intensive, and delayed. In doing so, we found models trained on patients admitted during the day do not generalize well to patients admitted at night, emphasizing the importance of careful cohort selection for model development. We also found harmful data shifts attributed to whether or not a patient was admitted from an acute care institution or nursing home, suggesting these settings have distinct patient populations. Institutional differences are among the most common causes of data shifts due to underlying differences in patient demographics, disease incidence and data-collection workflows<sup>2,11</sup>. We found models built on specific groups of hospitals such as community hospitals, undergo harmful data shifts when evaluated on academic hospitals and evaluated training strategies to mitigate model deterioration attributed to cross-site deployment. Lastly, we monitored data shifts over time and investigated key questions surrounding model updating like when to update a model, how much data to update on, and what data to use for the update. We found our drift-triggered continual updating strategy improved model performance and was more effective than maintaining a locked model during deployment.

However, it is unclear to what extent our findings will generalize, which is why it is critical to perform these experiments across several prediction tasks, patient populations and types of shifts. Likewise, many other sensitive attributes (e.g. socioeconomic status) and clinical scenarios (e.g. specialized hospitals) that merit evaluation remain. It is also imperative to characterize the extent to which other data modalities, like clinical notes, contribute to biases in clinical AI systems. There are a number of reasons these shifts could occur, including changes in the distribution of diagnoses, staffing, or resource allocation across patient populations. Identification of causal structures is a promising strategy to help explain the failures of fairness transfer across distribution shifts<sup>43</sup>. Given the sensitivity of clinical data, it is also important that future drift detection and retraining strategies consider privacy-preserving methods to ensure institutional boundaries are respected and autonomy is maintained over patient data<sup>44-46</sup>.

In this study, we developed a drift-triggered continual learning strategy to improve model performance over time. However, it is worth noting that continual learning is not without risks, including catastrophic forgetting and feedback loops<sup>47-49</sup>. Unfortunately, our dataset is unable to fully capture these long-term trends, but as more data is accumulated it will become possible to understand the impact of these model updating strategies over extended periods of time. Another caveat is that the current regulatory state of continual learning systems does not clearly define how and what aspects of a clinical AI system are permitted to change following authorisation<sup>41</sup>. There are also several other training and updating strategies we did not explore, which can be leveraged to improve model performance in the presence of data shifts, including domain generalization (DG)<sup>50,51</sup>, representation learning<sup>13,52</sup>, meta learning<sup>53,54</sup>, and multi-task learning<sup>55,56</sup>. For instance, consideration of other relevant prediction tasks (e.g. LOS, ICU transfer)<sup>55</sup> or patient populations<sup>56</sup> for pre-training could improve model generalization. Similarly, DG methods have been used as an alternative to baseline empirical risk minimization (ERM), to mitigate data shift<sup>57</sup>. However, many DG methods have repeatedly only been shown to improve performance in the context of extreme synthetic shifts and demonstrate poor performance on real world EHR data<sup>58,59</sup>. Instead, alternative ERM approaches (i.e. those that use stratified training, balanced subpopulation sampling, or worst-case model selection) outperform DG methods and show promise in mitigating model bias<sup>50,51,60</sup>. Unfortunately, many studies fail to consider strong and realistic ERM baselines.

Clinical AI systems are complex, and each will differ in its biases and optimal retraining and updating procedures. As such, we have developed a monitoring and evaluation pipeline as part of a broader ML operations (MLOps) framework for clinical AI systems<sup>37</sup> to facilitate robust evaluation and monitoring prior to deployment. Too often clinical ML models are reported with high performance metrics, while being developed in isolation. It is important to ensure that models are designed with deployment in mind, to ensure the responsible deployment of clinical AI systems. We hope our work permits the robust

evaluation and monitoring of clinical AI systems in an effort to bridge the gap between model development and deployment<sup>61–63</sup>.

## Methods

### Cohort Data

We conducted this study using de-identified Electronic Health Record (EHR) and hospital administrative data from 109,802 patients admitted to the general internal medicine (GIM) wards from 2015–2020 across 7 large hospitals in the Toronto, Canada-area. Of the 7 hospitals, 5 are academic hospitals (Hospital 1, Hospital 2, Hospital 3, Hospital 6, Hospital 7) and 2 are community hospitals (Hospital 4, Hospital 5).

### Ethics Approval

All patient data was collected and approved through GEMINI<sup>64,65</sup> under the oversight of the research ethics board (REB) at the Toronto Academic Health Science Network (REB reference number 15-087). The extension of the REB approval was issued by the Unity Health Toronto REB (reference number 15-087). A separate REB approval was obtained for Trillium Health Partners. All experiments were performed in accordance with institutional guidelines and regulations.

### Model Features

The *base* model consisted of 91 features comprising laboratory tests, blood transfusions, imaging reports and administrative features (**Supplementary Table 1**). The *base+CM* model consisted of the 91 features used in our base model, in addition to 18 comorbidities derived using ICD-10 codes (**Supplementary Table 2**). The *base+DxC* model consisted of the 91 features used in our base model, in addition to the 22 groupings of ICD-10 diagnosis codes (**Supplementary Table 3**). The input features used for time-series modelling were aggregated by taking the mean for 24-hour timesteps, over 144 hours.

### All-Cause In-Hospital Mortality Decompensation Prediction

Our goal was to predict whether the patient's health will rapidly deteriorate<sup>55</sup>. Each instance of this task is a binary classification instance and predictions are made every 24 hours for the risk of in-hospital mortality within the next two weeks starting 24 hours after admission using the target replication approach<sup>66</sup>. In addition to longitudinal clinical measures, demographics are included as static variables at every time step for the prediction task. Labels were encoded as 1 if a patient died within the next 2 weeks, 0 if they were alive within the next 2 weeks and -1 if they were discharged. Missing values were imputed using forward filling followed by backward filling. Unless a custom data split was applied (i.e. for the data shift experiments described below), a training/validation/test split of 8:1:1 was used. The training, validation and test data were normalized independently by subtracting the mean and scaling to unit variance. A long short-term memory (LSTM) recurrent neural network (RNN)<sup>66</sup> with 2 hidden layers, 64 hidden cells and a dropout rate of 0.2 was implemented using PyTorch<sup>67</sup>. The LSTM RNN was optimized for binary cross entropy with logits loss using Adagrad<sup>68</sup>, a step size of 128, gamma of 0.5, learning rate of  $3.0 \times 10^{-2}$ , weight decay of  $1.0 \times 10^{-6}$  and batch size of 64. To account for the class imbalance, we reweighted our loss function by the fraction of controls/cases in the training data. Each model was trained over 128 epochs with early stopping using a patience of 3 and delta of 0. We used a sigmoid activation function to obtain prediction probabilities. We generated standard errors by making a random choice of weight initializations and dataset splits for 10 repetitions. For consistency, model level parameters (e.g. number of cells, number of layers) were kept fixed across all experiments.

### Monitoring and Evaluation Pipeline

We detected distributional shifts between source and target data using our monitoring and evaluation pipeline (**Figure 2**) which consists of:

- 1) **Shift application:** EHR data is sent to the *Shift Applicator*, which outputs a source and target dataset based on the clinical data shift experiment of choice (e.g. hospital type, seasons, etc.).
- 2) **Dimensionality reduction:** Dimensionality reduction is performed using the *Shift Reductor* to obtain a latent representation of the source and target data. This was done using the softmax outputs of a LSTM neural network label classifier trained on source data (Black Box Shift Detector; BBSD)<sup>69</sup>. The architecture and training of the BBSD is described above as the base model.
- 3) **Statistical testing:** Univariate two-sample testing was performed with a Kolmogorov-Smirnov Test using the *Shift Tester*, in order to identify if a harmful data shift has occurred between the latent representation of the source and target data<sup>28</sup>.
- 4) **Sensitivity test:** A drift sensitivity test was conducted by performing step (2) and (3) to detect data shifts for  $n = \{10, 20, 50, 100, 250, 500, 1000\}$  patients from the target data.
- 5) **Rolling window analysis:** A 14-day rolling window was used to assess model stability over time by sampling 1000 patients and performing step (2) and (3) to test for drift every day. The drift detector was updated every day with the last 25000 patients.

## Clinical Data Shift Experiments

We used prior knowledge to devise data splits that reflect real-life scenarios that may result in harmful data shifts and model degradation of clinical AI systems. For all experiments we trained a model on the in-distribution (ID) data and evaluated on ID data as the *baseline* and out-of-distribution (OOD) data as the *shift experiment*. Sensitivity tests were performed for each scenario using the trained model as the BBSD. The scenarios are as follows:

**Winter - Baseline:** Patients admitted in the winter (Nov-Feb). **Shift Experiment:** Patients admitted in the winter (June-Aug).

**Summer - Baseline:** Patients admitted in the summer (June-Aug) **Shift Experiment:** Patients admitted in the winter (Nov-Feb) .

**Community Hospitals - Baseline:** Academic hospitals (Hospital 1, Hospital 2, Hospital 3, Hospital 6, Hospital 7). **Shift Experiment:** Community hospitals (Hospital 4, Hospital 5).

**Academic Hospitals - Baseline:** Community hospitals (Hospital 4, Hospital 5). **Shift Experiment:** Academic hospitals (Hospital 1, Hospital 2, Hospital 3, Hospital 6, Hospital 7).

**Day Admission - Baseline:** Patients admitted during the day (7:30-19:30). **Shift Experiment:** Patients admitted during the night (0:00-7:30,19:30:23:59).

**Night Admission - Baseline:** Patients admitted during the night (0:00-7:30,19:30:23:59). **Shift Experiment:** Patients admitted during the day (7:30-19:30).

**Weekend Admission - Baseline:** Patients admitted on the weekend (i.e. Saturday and Sunday). **Shift Experiment:** Patients admitted on a weekday (i.e. Monday to Friday).

**Weekday Admission - Baseline:** Patients admitted on a weekday (i.e. Monday to Friday). **Shift Experiment:** Patients admitted on the weekend (i.e. Saturday and Sunday).

**Admitted from Nursing Home- Baseline:** Patients admitted from nursing homes. **Shift Experiment:** Patients not admitted from nursing homes.

**Not Admitted from Nursing Home - Baseline:** Patients not admitted from nursing homes. **Shift Experiment:** Patients admitted from nursing homes.

**Admitted from Acute Care Institution- Baseline:** Patients admitted from acute care institutions. **Shift Experiment:** Patients not admitted from acute care institutions.

**Not Admitted from Acute Care Institution - Baseline:** Patients not admitted from acute care institutions. **Shift Experiment:** Patients admitted from acute care institutions.

**Sex - Baseline:** Patients of all sexes. **Shift Experiments:** Patients that are i) males ii) females.

**Age - Baseline:** Patients of all ages. **Shift Experiments:** Patients that are i) 18-29 years ii) 30-44 years iii) 45-64 years iv) 65+ years.

## Transfer learning



To evaluate the optimal training strategy for each hospital, we compared models trained using i) a single hospital, ii) each hospital type (i.e. academic, community) and iii) all hospitals. We compared i) *pre-training* where we used a model pre-trained on source data and evaluated it on out-of-distribution data from the target hospital ii) *fine-tuning* where the single-site and hospital-type specific models were fine-tuned on the target hospital using 1 epoch or 10 epochs, and iii) *ablation* of a single hospital from the cross-site model, for each hospital. Each strategy was evaluated on a held out test set for each of the 7 hospital sites.

## Continual Learning

In order to mitigate model drift due to temporal data shifts, we compared the following continual learning strategies to a baseline where the model was kept locked and no changes or updates were made:

**Periodic Updating** - The model is updated at regular time intervals of  $n = \{7, 14, 30, 60\}$  days.

**Most Recent Updating**- When drift is detected, the model is updated using the most recent  $n$  number of days where  $n = \{7, 14, 30, 60, 120, 180, 270\}$  days.

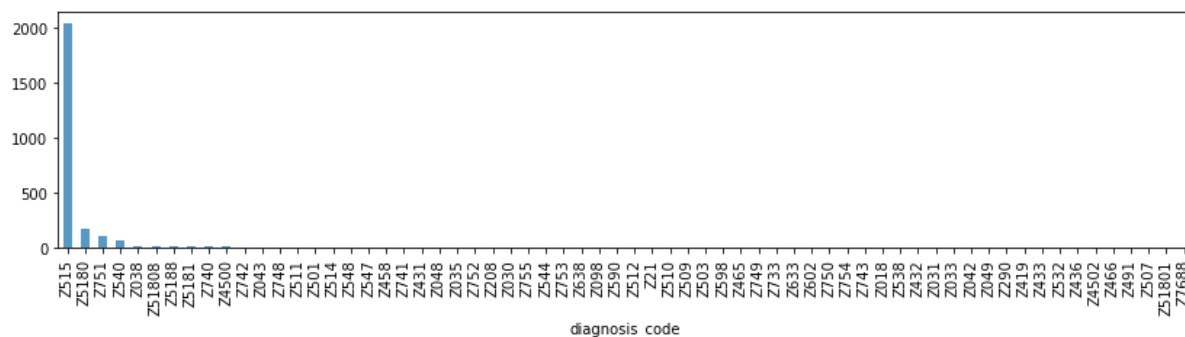
**Cumulative Updating** - When drift is detected, the model is updated using all the patient encounters seen to-date.

Model updating methods were optimized for the retraining window size, lookback window, sample size, drift threshold and number of epochs. We also compared sampling strategies where we used i) all the encounters in the retraining window ii) only the correctly predicted encounters in the retraining window and iii) only the positively predicted encounters in the retraining window.

## Acknowledgements

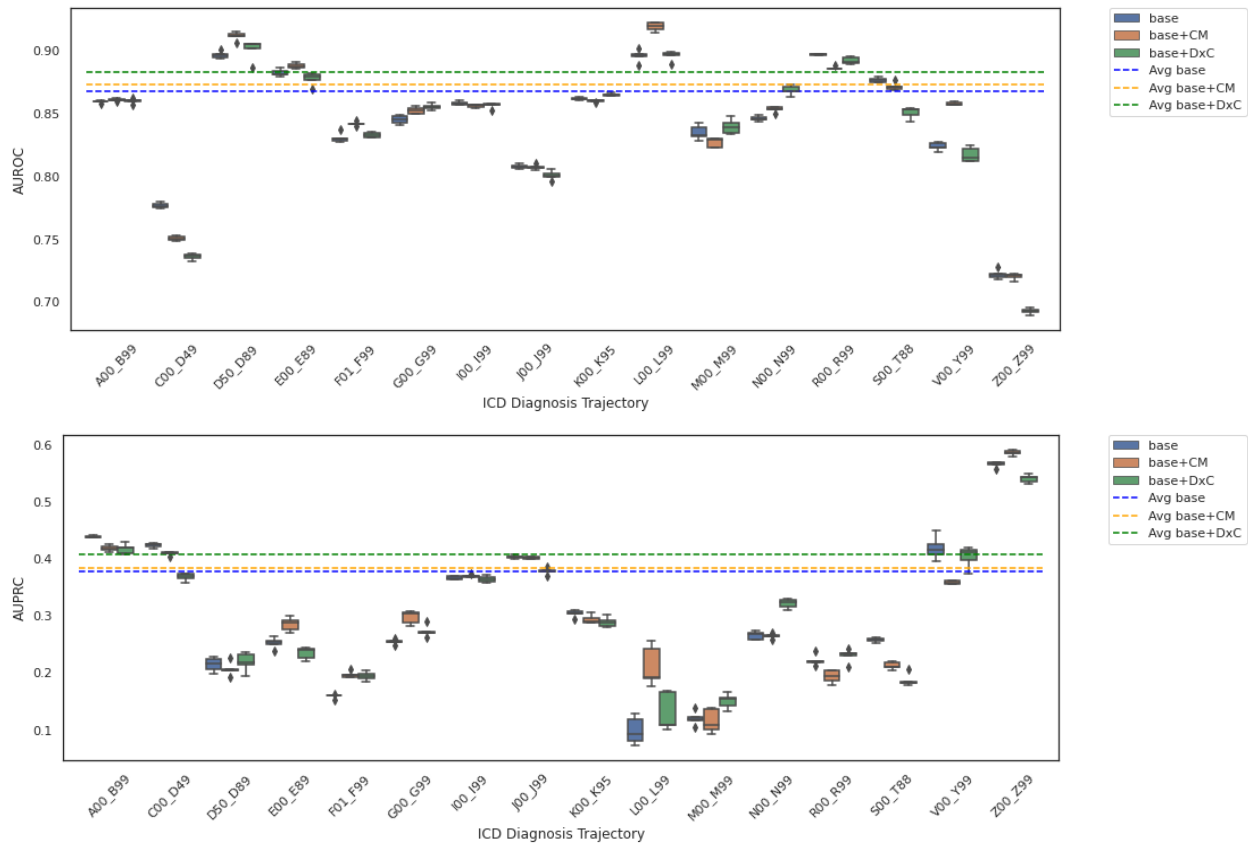
This work is made possible due to the data obtained from the General Medicine Inpatient Initiative (GEMINI), and we acknowledge the GEMINI team for their support. We also acknowledge HPC4Health, for enabling high performance computing environments which were involved in the development of this work. Finally, we acknowledge support from the Vector Institute and its vibrant community working at the intersection of health and machine learning. V.S. is supported by Ontario Graduate Scholarship and a Vector institute grant. A.V. is supported by the Temerty Professorship in Artificial Intelligence Research and Education in Medicine at the University of Toronto. D.M. is supported by the CIBC Children's Foundation Chair in Child Health Research. AG is supported by the Varma Family Chair and CIFAR AI Chair.

## Supplementary Materials

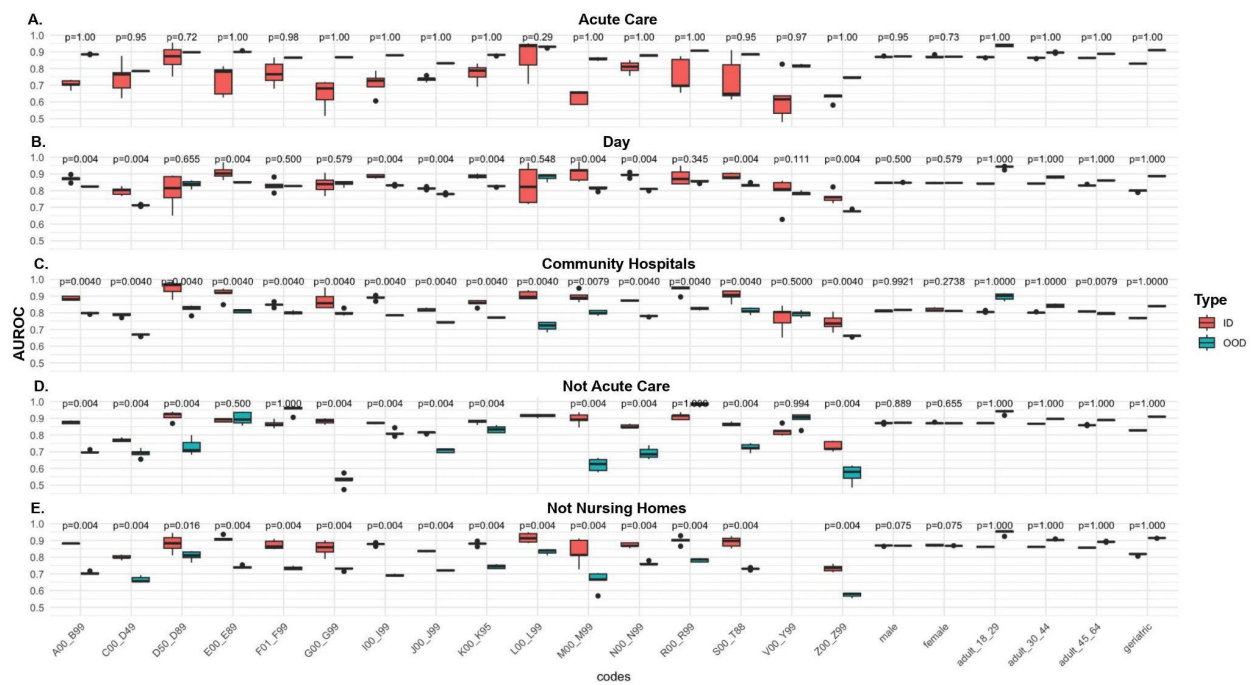


**Supplementary Figure 1.** Distribution of diagnosis codes across factors influencing health status and contact with health services (Z00-Z99).

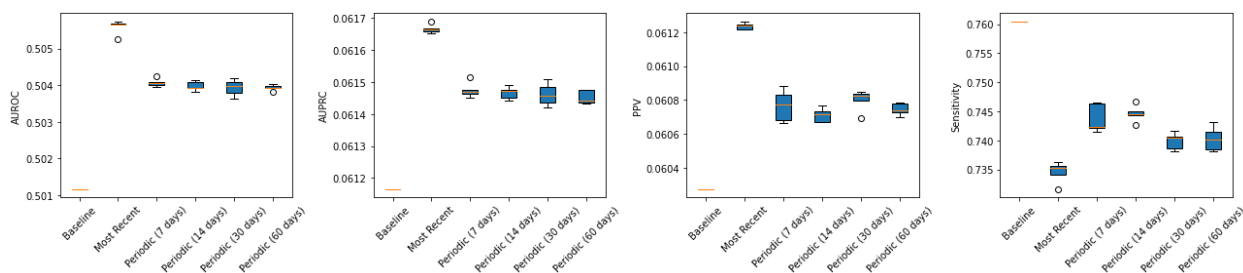




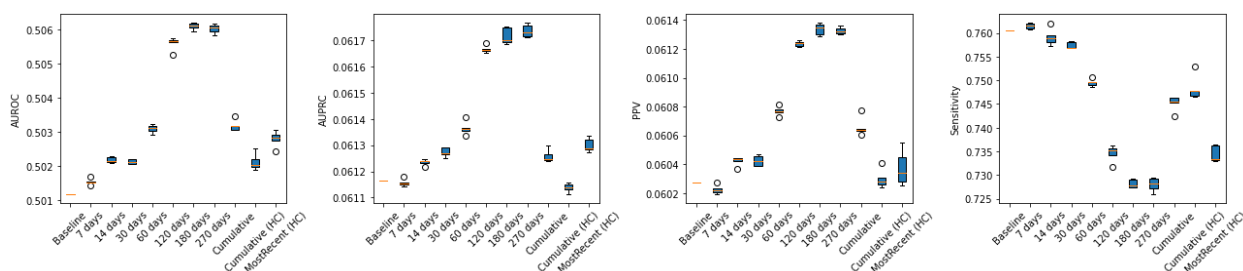
**Supplementary Figure 2.** Performance of model using no prior information (base), comorbidities (base+CM) and ICD-10 diagnosis codes (base+DxC) across diagnosis codes measured using (A) AUROC and (B) AUPRC.



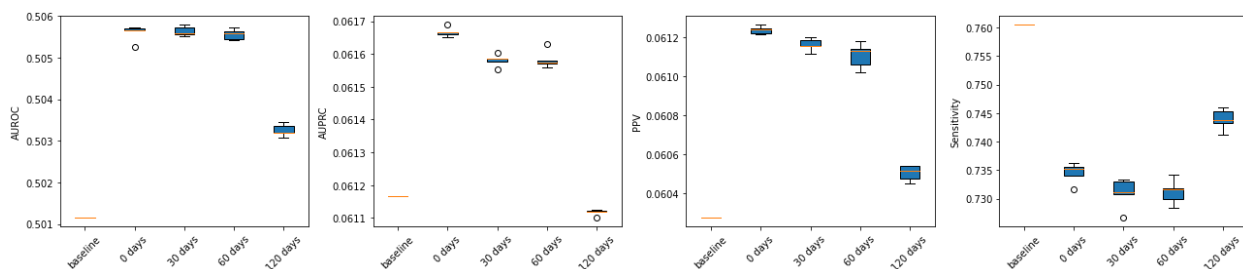
**Supplementary Figure 3.** AUROC for in-distribution (ID) and out-of-distribution (OOD) data, across ICD-10 diagnosis codes, age and sex, for for scenarios where harmful data shifts were detected: **(A)** model trained on patients admitted from acute care institutions **(B)** model trained on patients admitted during the day **(C)** model trained on patients admitted from community hospitals **(D)** model trained on patients admitted not from acute care institutions **(E)** model trained on patients not admitted from nursing homes. P-values were calculated using a one-sided Mann-Whitney U test.



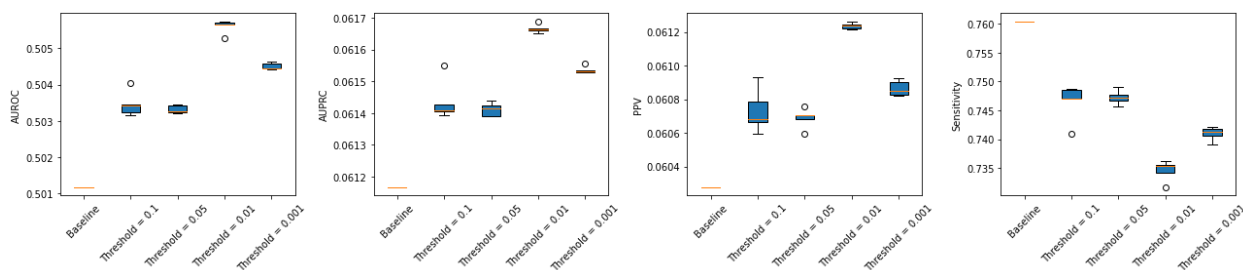
**Supplementary Figure 4.** Comparison of AUROC, AUPRC, PPV, and sensitivity when updating periodically every  $n = 7, 14, 30,$  and  $60$  days.



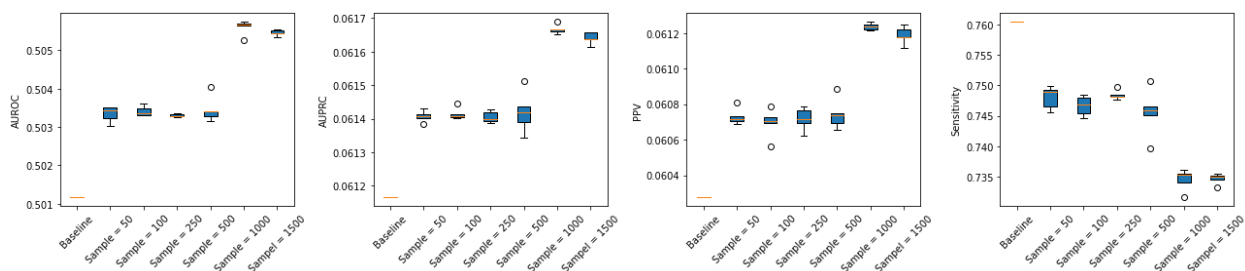
**Supplementary Figure 5.** Comparison of AUROC, AUPRC, PPV, and sensitivity across strategies retraining using a dynamic window of the most recent encounters ( $n = 7, 14, 30, 60, 120, 180$  days) and cumulatively.



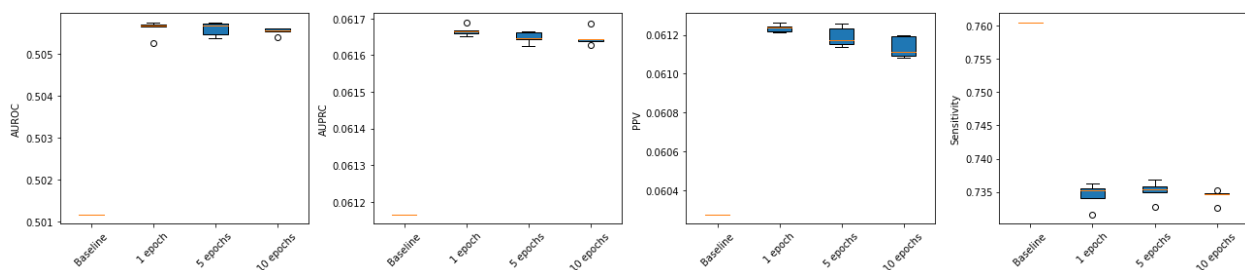
**Supplementary Figure 6.** Comparison of AUROC, AUPRC, PPV, and sensitivity across increasing lookback windows ( $n = 0, 30, 60, 120$  days).



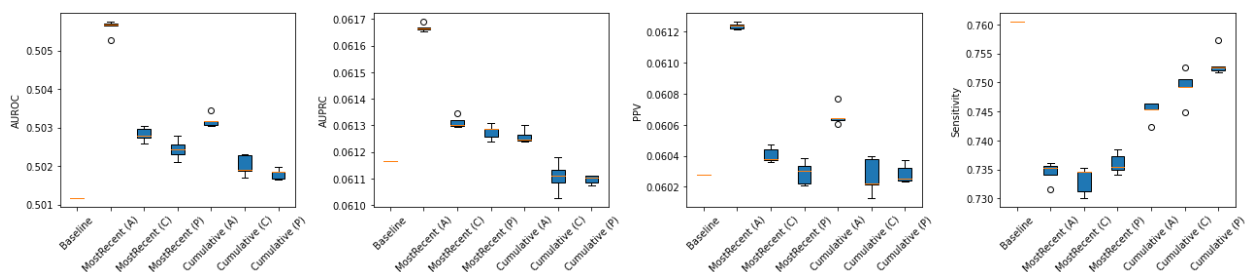
**Supplementary Figure 7.** Comparison of AUROC, AUPRC, PPV, and sensitivity across varying drift p-value thresholds for retraining ( $p = 0.1, 0.05, 0.01, 0.001$ ).



**Supplementary Figure 8.** Comparison of AUROC, AUPRC, PPV, and sensitivity across varying sample sizes for drift tests (n = 50, 100, 250, 500, 1000, 1500).



**Supplementary Figure 9.** Comparison of AUROC, AUPRC, PPV, and sensitivity across varying epochs for retraining (n = 1, 5, 10).



**Supplementary Figure 10.** Comparison of AUROC, AUPRC, PPV, and sensitivity when updating using the most recent and cumulative retraining strategy with all encounters (A), correctly predicted encounters (C) or positively predicted encounters (P).

Feature Type	# of Features	Features
Administrative	16	sex, age, prev_encounter_count, triage_level_emergent, triage_level_no_info, triage_level_non-urgent, triage_level_resuscitation, triage_level_semi-urgent, triage_level_urgent, readmission_new_to_acute, readmission_nota, readmission_planned_from_acute, readmission_unplanned_7_day_acute, readmission_unplanned_7_day_day_surg, readmission_unplanned_8_to_28_day_acute, from_nursing_home_mapped, from_acute_care_institution_mapped
Interventions	6	unmapped_intervention, inv_mech_vent_mapped, endoscopy_mapped, dialysis_mapped, surgery_mapped, interventional
Labs	55	albumin, alp, alt, aptt, arterial_paco2, arterial_pao2, arterial_ph, ast, bicarbonate, bilirubin, blood_urea_nitrogen, calcium, calcium_ionized, creatinine, crp, d-dimer, esr, ferritin, fibrinogen, glucose_fasting, glucose_point_of_care, glucose_random, hba1c, hematocrit, hemoglobin, high_sensitivity_troponin, influenza, inr, ketone, lactate_arterial, lactate_venous, ldh, lipase, lymphocyte, mean_cell_volume, neutrophils, other, platelet_count, potassium, pt, serum_alcohol, serum_osmolality, sodium, troponin, tsh, urinalysis, urine_osmolality, urine_sodium, urine

		specific gravity, venous pco2, venous ph, vitamin b12, vitamin d, white blood cell count
Imaging Reports	5	ct, mri, x-ray, echo, ultrasound
Blood Transfusions	2	rbc, non-rbc
Comorbidities <i>(Only used in Base+CM)</i>	18	Kidney disease: N18, N19 Ischemic heart disease: I20-I52 Cerebrovascular disease: I60-69 Hypertension: I10-I15 Diabetes: E10-E13 Hyperlipidemia: E78 Hypertension: I10 Congestive heart failure: I50 Cancer: C00-D49 Dyspnea: R06 COPD: J44 Asthma: J45 Pulmonary embolism: I26 Connective tissue disease: I30-I36 Inflammatory bowel disease: K50, K51, Osteoarthritis: M15-M19 Rheumatoid arthritis: M05-M14 HIV: B20-B24
ICD-10 Diagnosis Codes <i>(Only used in Base+DxC)</i>	22	Certain infectious and parasitic diseases: A00-B99 Neoplasms: C00-D49 Diseases of the blood and blood-forming organs and certain disorders involving the immune mechanism: D50-D89 Endocrine, nutritional and metabolic diseases: E00-E89 Mental, Behavioral and Neurodevelopmental disorders: F01-F99 Diseases of the nervous system: G00-G99 Diseases of the eye and adnexa: H00-H59 Diseases of the ear and mastoid process: H60-H95 Diseases of the circulatory system: I00-I99 Diseases of the respiratory system: J00-J99 Diseases of the digestive system: K00-K95 Diseases of the skin and subcutaneous tissue: L00-L99 Diseases of the musculoskeletal system and connective tissue: M00-M99 Diseases of the genitourinary system: N00-N99 Pregnancy, childbirth and the puerperium: O00-O99 Certain conditions originating in the perinatal period: P00-P96 Congenital malformations, deformations and chromosomal abnormalities: Q00-Q99 Symptoms, signs and abnormal clinical and laboratory findings, not elsewhere classified: R00-R99 Injury, poisoning and certain other consequences of external causes: S00-T88 External causes of morbidity: V00-Y99 COVID19: U07-U08 Factors influencing health status and contact with health services: Z00-Z99

**Supplementary Table 1.** EHR features used for mortality risk prediction.

Mortality	False							True						
	1	2	3	4	5	6	7	1	2	3	4	5	6	7
<b>A00_B99</b>	0.086	0.065	0.065	0.078	0.082	0.083	0.072	0.152	0.115	0.099	0.135	0.112	0.119	0.129
<b>C00_D49</b>	0.065	0.059	0.039	0.036	0.029	0.066	0.035	0.177	0.125	0.111	0.080	0.073	0.190	0.102
<b>D50_D89</b>	0.032	0.018	0.019	0.015	0.016	0.048	0.018	0.011	0.004	0.003	0.004	0.003	0.015	0.007
<b>E00_E89</b>	0.063	0.050	0.063	0.058	0.051	0.055	0.066	0.017	0.025	0.019	0.028	0.018	0.020	0.023
<b>F01_F99</b>	0.047	0.047	0.053	0.052	0.072	0.026	0.057	0.020	0.032	0.031	0.030	0.048	0.013	0.023
<b>G00_G99</b>	0.038	0.051	0.038	0.065	0.039	0.026	0.028	0.011	0.031	0.028	0.028	0.031	0.008	0.012

H00_H59	0.004	0.006	0.003	0.005	0.001	0.001	0.002	NA	NA	NA	NA	NA	NA	NA
H60_H95	0.006	0.006	0.004	0.012	0.009	0.006	0.006	NA	NA	NA	NA	NA	NA	NA
I00_I99	0.101	0.147	0.143	0.138	0.219	0.138	0.131	0.103	0.211	0.165	0.175	0.249	0.119	0.158
J00_J99	0.127	0.132	0.144	0.130	0.106	0.137	0.163	0.217	0.238	0.175	0.256	0.229	0.197	0.258
K00_K95	0.108	0.103	0.118	0.059	0.046	0.106	0.094	0.051	0.081	0.079	0.048	0.039	0.085	0.064
L00_L99	0.028	0.026	0.029	0.024	0.021	0.022	0.022	0.003	0.005	0.003	0.004	0.006	0.004	0.001
M00_M99	0.044	0.041	0.049	0.053	0.036	0.036	0.043	0.009	0.008	0.011	0.011	0.007	0.006	0.006
N00_N99	0.059	0.066	0.062	0.076	0.085	0.057	0.063	0.025	0.042	0.030	0.054	0.057	0.035	0.039
O00_O99	0.002	0.001	0.001	0.003	0.001	0.000	0.001	NA	NA	NA	NA	NA	NA	NA
Q00_Q99	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
R00_R99	0.094	0.134	0.120	0.113	0.116	0.142	0.129	0.019	0.030	0.028	0.011	0.030	0.036	0.035
S00_T88	0.074	0.038	0.041	0.072	0.061	0.036	0.059	0.027	0.023	0.016	0.033	0.031	0.016	0.018
V00_Y99	0.002	0.002	0.004	0.005	0.003	0.003	0.004	0.004	0.003	0.004	0.013	0.011	0.005	0.009
Z00_Z99	0.020	0.008	0.004	0.006	0.005	0.011	0.008	0.156	0.027	0.197	0.089	0.054	0.131	0.117

**Supplementary Table 2.** Proportion of patient encounters across ICD-10 diagnosis codes, hospital and mortality status. Values for groupings of diagnosis codes with < 5 patient encounters have been omitted due to privacy preserving practices required by GEMINI.

			Base	Base + Comorbidities	Base + ICD-10 Diagnosis Codes
Source	Target	Metric	Target – Source (%)	Target – Source (%)	Target – Source (%)
Community Hospitals	Academic Hospitals	AUROC	-4.0	-2.6	-0.9
		AUPRC	-6.8	-5.0	-3.2
Academic Hospitals	Community Hospitals	AUROC	-1.1	-1.8	0.0
		AUPRC	-3.2	-4.8	0.0
Excl. Winter	Seasonal Winter	AUROC	+0.7	-0.3	-1.2
		AUPRC	-0.6	+0.9	-0.1
Excl. Summer	Seasonal Summer	AUROC	+0.8	+1.2	+0.3
		AUPRC	+1.0	+0.9	-0.3
Night Admission	Day Admission	AUROC	+1.6	-0.2	+0.2
		AUPRC	+3.0	+0.9	+2.1
Day Admission	Night Admission	AUROC	-1.9	-0.7	-0.8
		AUPRC	-5.6	-3.2	-3.7

**Supplementary Table 3.** Performance of models on in-distribution (source) and out-of-distribution (target) data without prior information (base), with comorbidities (base+CM) and with diagnosis codes (base+DxC) as features.

## References

1. Yan, L. *et al.* An interpretable mortality prediction model for COVID-19 patients. *Nature Machine Intelligence* **2**, 283–288 (2020).



2. Barish, M., Bolourani, S., Lau, L. F., Shah, S. & Zanos, T. P. External validation demonstrates limited clinical utility of the interpretable mortality prediction model for patients with COVID-19. *Nature Machine Intelligence* **3**, 25–27 (2020).
3. Choi, M. H. *et al.* Mortality prediction of patients in intensive care units using machine learning algorithms based on electronic health records. *Sci. Rep.* **12**, 7180 (2022).
4. Verma, A. A. *et al.* Characteristics and outcomes of hospital admissions for COVID-19 and influenza in the Toronto area. *CMAJ* **193**, E410–E418 (2021).
5. Rasmy, L. *et al.* Recurrent neural network models (CovRNN) for predicting outcomes of patients with COVID-19 on admission to hospital: model development and validation using electronic health record data. *Lancet Digit Health* **4**, e415–e425 (2022).
6. Iwase, S. *et al.* Prediction algorithm for ICU mortality and length of stay using machine learning. *Sci. Rep.* **12**, 12912 (2022).
7. Gao, Y. *et al.* Machine learning based early warning system enables accurate mortality risk prediction for COVID-19. *Nat. Commun.* **11**, 5033 (2020).
8. Adams, R. *et al.* Prospective, multi-site study of patient outcomes after implementation of the TREWS machine learning-based early warning system for sepsis. *Nat. Med.* **28**, 1455–1460 (2022).
9. Henry, K. E. *et al.* Factors driving provider adoption of the TREWS machine learning-based early warning system and its effects on sepsis treatment timing. *Nat. Med.* **28**, 1447–1454 (2022).
10. Rajkomar, A. *et al.* Scalable and accurate deep learning with electronic health records. *NPJ Digit Med* **1**, 18 (2018).
11. Zhang, A., Xing, L., Zou, J. & Wu, J. C. Shifting machine learning for healthcare from development to deployment and from models to data. *Nat Biomed Eng* (2022) doi:10.1038/s41551-022-00898-y.
12. Subbaswamy, A., Adams, R. & Saria, S. Evaluating Model Robustness and Stability to Dataset Shift. in *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics* (eds. Banerjee, A. & Fukumizu, K.) vol. 130 2611–2619 (PMLR, 13–15 Apr 2021).
13. Nestor, B. *et al.* Feature Robustness in Non-stationary Health Records: Caveats to Deployable Model Performance in Common Clinical Machine Learning Tasks. in *Proceedings of the 4th Machine Learning for Healthcare Conference* (eds. Doshi-Velez, F. *et al.*) vol. 106 381–405 (PMLR, 09--10

Aug 2019).

14. Futoma, J., Simons, M., Panch, T., Doshi-Velez, F. & Celi, L. A. The myth of generalisability in clinical research and machine learning in health care. *Lancet Digit Health* **2**, e489–e492 (2020).
15. Cohen, J. P. *et al.* Problems in the deployment of machine-learned models in health care. *CMAJ: Canadian Medical Association journal = journal de l'Association medicale canadienne* vol. 193 E1391–E1394 (2021).
16. Duckworth, C. *et al.* Using explainable machine learning to characterise data drift and detect emergent health risks for emergency department admissions during COVID-19. *Sci. Rep.* **11**, 23017 (2021).
17. Otlés, E. *et al.* Mind the Performance Gap: Examining Dataset Shift During Prospective Validation. in *Proceedings of the 6th Machine Learning for Healthcare Conference* (eds. Jung, K., Yeung, S., Sendak, M., Sjoding, M. & Ranganath, R.) vol. 149 506–534 (PMLR, 06--07 Aug 2021).
18. Finlayson, S. G. *et al.* The Clinician and Dataset Shift in Artificial Intelligence. *N. Engl. J. Med.* **385**, 283–286 (2021).
19. Avati, A. *et al.* BEDS-Bench: Behavior of EHR-models under Distributional Shift--A Benchmark. *arXiv [cs.LG]* (2021).
20. Chen, I. Y. *et al.* Ethical Machine Learning in Healthcare. *Annu Rev Biomed Data Sci* **4**, 123–144 (2021).
21. Koh, P. W. *et al.* WILDS: A Benchmark of in-the-Wild Distribution Shifts. in *Proceedings of the 38th International Conference on Machine Learning* (eds. Meila, M. & Zhang, T.) vol. 139 5637–5664 (PMLR, 18--24 Jul 2021).
22. Rahmani, K. *et al.* Assessing the effects of data drift on the performance of machine learning models used in clinical sepsis prediction. *medRxiv* (2022) doi:10.1101/2022.06.06.22276062.
23. Singh, H., Mhasawade, V. & Chunara, R. Generalizability Challenges of Mortality Risk Prediction Models: A Retrospective Analysis on a Multi-center Database. Preprint at <https://doi.org/10.1101/2021.07.14.21260493>.
24. Seyyed-Kalantari, L., Zhang, H., McDermott, M. B. A., Chen, I. Y. & Ghassemi, M. Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient

- populations. *Nat. Med.* **27**, 2176–2182 (2021).
25. Adamson, A. S. & Smith, A. Machine Learning and Health Care Disparities in Dermatology. *JAMA Dermatol.* **154**, 1247–1248 (2018).
  26. Wong, A. *et al.* External Validation of a Widely Implemented Proprietary Sepsis Prediction Model in Hospitalized Patients. *JAMA Intern. Med.* **181**, 1065–1070 (2021).
  27. Cao, T., Huang, C.-W., Hui, D. Y.-T. & Cohen, J. P. A Benchmark of Medical Out of Distribution Detection. *arXiv [cs.LG]* (2020).
  28. Rabanser, S., Günnemann, S. & Lipton, Z. Failing loudly: An empirical study of methods for detecting dataset shift. *Adv. Neural Inf. Process. Syst.* **32**, (2019).
  29. Maggio, S. & Dreyfus-Schmidt, L. Ensembling Shift Detectors: An Extensive Empirical Evaluation. in *Machine Learning and Knowledge Discovery in Databases. Research Track* 362–377 (Springer International Publishing, 2021).
  30. Sauer, C. M. *et al.* Leveraging electronic health records for data science: common pitfalls and how to avoid them. *The Lancet Digital Health* vol. 4 e893–e898 Preprint at [https://doi.org/10.1016/s2589-7500\(22\)00154-6](https://doi.org/10.1016/s2589-7500(22)00154-6) (2022).
  31. Feng, J. *et al.* Clinical artificial intelligence quality improvement: towards continual monitoring and updating of AI algorithms in healthcare. *NPJ Digit Med* **5**, 66 (2022).
  32. Lee, C. S. & Lee, A. Y. Clinical applications of continual learning machine learning. *Lancet Digit Health* **2**, e279–e281 (2020).
  33. Amrollahi, F., Shashikumar, S. P., Holder, A. L. & Nemati, S. Leveraging clinical data across healthcare institutions for continual learning of predictive risk models. *Sci. Rep.* **12**, 8380 (2022).
  34. Bozkurt, S. *et al.* Reporting of demographic data and representativeness in machine learning models using electronic health records. *J. Am. Med. Inform. Assoc.* **27**, 1878–1884 (2020).
  35. for Health Information, C. I. Canadian coding standards for version 2018 ICD-10-CA and CCI. Preprint at (2018).
  36. Health Organization, W. Palliative care. [https://apps.who.int/iris/bitstream/handle/10665/44024/9241547345\\_eng.pdf](https://apps.who.int/iris/bitstream/handle/10665/44024/9241547345_eng.pdf).
  37. Krishnan, A. *et al.* CyclOps: Cyclical development towards operationalizing ML models for health.

*bioRxiv* (2022) doi:10.1101/2022.12.02.22283021.

38. Han, S. S. *et al.* Classification of the Clinical Images for Benign and Malignant Cutaneous Tumors Using a Deep Learning Algorithm. *J. Invest. Dermatol.* **138**, 1529–1538 (2018).
39. Davis, S. E., Lasko, T. A., Chen, G., Siew, E. D. & Matheny, M. E. Calibration drift in regression and machine learning models for acute kidney injury. *J. Am. Med. Inform. Assoc.* **24**, 1052–1061 (2017).
40. Zech, J. R. *et al.* Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: A cross-sectional study. *PLoS Med.* **15**, e1002683 (2018).
41. Vokinger, K. N., Feuerriegel, S. & Kesselheim, A. S. Mitigating bias in machine learning for medicine. *Commun. Med.* **1**, 25 (2021).
42. Gianfrancesco, M. A., Tamang, S., Yazdany, J. & Schmajuk, G. Potential Biases in Machine Learning Algorithms Using Electronic Health Record Data. *JAMA Intern. Med.* **178**, 1544–1547 (2018).
43. Schrouff, J. *et al.* Maintaining fairness across distribution shift: do we have viable solutions for real-world applications? *arXiv [cs.LG]* (2022).
44. Kaissis, G. *et al.* End-to-end privacy preserving deep learning on multi-institutional medical imaging. *Nature Machine Intelligence* **3**, 473–484 (2021).
45. Warnat-Herresthal, S. *et al.* Swarm Learning for decentralized and confidential clinical machine learning. *Nature* **594**, 265–270 (2021).
46. Price, W. N., 2nd & Cohen, I. G. Privacy in the age of medical big data. *Nat. Med.* **25**, 37–43 (2019).
47. Adam, G. A., Chang, C.-H. K., 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* (eds. Doshi-Velez, F. *et al.*) vol. 126 710–731 (PMLR, 07–08 Aug 2020).
48. Armstrong, J. & Clifton, D. Continual learning of longitudinal health records. *arXiv [cs.LG]* (2021).
49. Perkonigg, M. *et al.* Dynamic memory to alleviate catastrophic forgetting in continual learning with medical imaging. *Nat. Commun.* **12**, 5678 (2021).
50. Guo, L. L. *et al.* Evaluation of domain generalization and adaptation on improving model robustness to temporal dataset shift in clinical medicine. *Sci. Rep.* **12**, 2726 (2022).
51. Zhang, H. *et al.* Improving the Fairness of Chest X-ray Classifiers. in *Proceedings of the Conference*

- on Health, Inference, and Learning* (eds. Flores, G., Chen, G. H., Pollard, T., Ho, J. C. & Naumann, T.) vol. 174 204–233 (PMLR, 07–08 Apr 2022).
52. Izmailov, P., Kirichenko, P., Gruver, N. & Wilson, A. G. On Feature Learning in the Presence of Spurious Correlations. *arXiv [cs.LG]* (2022).
  53. Zhang, M. *et al.* Adaptive risk minimization: Learning to adapt to domain shift. *Adv. Neural Inf. Process. Syst.* **34**, 23664–23678 (2021).
  54. Ajay, A., Gupta, A., Ghosh, D., Levine, S. & Agrawal, P. Distributionally Adaptive Meta Reinforcement Learning. *arXiv [cs.LG]* (2022).
  55. Harutyunyan, H., Khachatrian, H., Kale, D. C., Ver Steeg, G. & Galstyan, A. Multitask learning and benchmarking with clinical time series data. *Sci Data* **6**, 96 (2019).
  56. Suresh, H., Gong, J. J. & Guttag, J. V. Learning Tasks for Multitask Learning: Heterogenous Patient Populations in the ICU. in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* 802–810 (Association for Computing Machinery, 2018).
  57. Bellot, A. & van der Schaar, M. Accounting for Unobserved Confounding in Domain Generalization. *arXiv [stat.ML]* (2020).
  58. Zhang, H. *et al.* An empirical framework for domain generalization in clinical settings. in *Proceedings of the Conference on Health, Inference, and Learning* 279–290 (Association for Computing Machinery, 2021).
  59. Pfohl, S. R. *et al.* A comparison of approaches to improve worst-case predictive model performance over patient subpopulations. *Scientific Reports* vol. 12 Preprint at <https://doi.org/10.1038/s41598-022-07167-7> (2022).
  60. Gulrajani, I. & Lopez-Paz, D. In Search of Lost Domain Generalization. *arXiv [cs.LG]* (2020).
  61. Liu, X. *et al.* Reporting guidelines for clinical trial reports for interventions involving artificial intelligence: the CONSORT-AI Extension. *BMJ* **370**, m3164 (2020).
  62. DECIDE-AI: new reporting guidelines to bridge the development-to-implementation gap in clinical artificial intelligence. *Nat. Med.* **27**, 186–187 (2021).
  63. Rivera, S. C. *et al.* Guidelines for clinical trial protocols for interventions involving artificial intelligence: the SPIRIT-AI Extension. *BMJ* **370**, m3210 (2020).



64. Verma, A. A. *et al.* Patient characteristics, resource use and outcomes associated with general internal medicine hospital care: the General Medicine Inpatient Initiative (GEMINI) retrospective cohort study. *CMAJ Open* vol. 5 E842–E849 Preprint at <https://doi.org/10.9778/cmajo.20170097> (2017).
65. Verma, A. A. *et al.* Assessing the quality of clinical and administrative data extracted from hospitals: the General Medicine Inpatient Initiative (GEMINI) experience. *Journal of the American Medical Informatics Association* vol. 28 578–587 Preprint at <https://doi.org/10.1093/jamia/ocaa225> (2021).
66. Lipton, Z. C., Kale, D. C., Elkan, C. & Wetzell, R. Learning to Diagnose with LSTM Recurrent Neural Networks. *arXiv [cs.LG]* (2015).
67. Paszke, A. *et al.* PyTorch: An imperative style, high-performance deep learning library. *arXiv [cs.LG]* (2019).
68. Stochastic Optimization. Adaptive Subgradient Methods for. <https://www.jmlr.org/papers/volume12/duchi11a/duchi11a.pdf> (2011).
69. Lipton, Z., Wang, Y.-X. & Smola, A. Detecting and Correcting for Label Shift with Black Box Predictors. in *Proceedings of the 35th International Conference on Machine Learning* (eds. Dy, J. & Krause, A.) vol. 80 3122–3130 (PMLR, 10–15 Jul 2018).