

Interpreting a Recurrent Neural Network’s Predictions of ICU Mortality Risk

Long V. Ho^a, Melissa Aczon^a, David Ledbetter^a, Randall Wetzel^a

^aThe Laura P. and Leland K. Whittier Virtual Pediatric Intensive Care Unit
Children’s Hospital Los Angeles, 4650 Sunset Blvd, Los Angeles, CA 90027, United States

Abstract

Deep learning has demonstrated success in many applications; however, their use in healthcare has been limited due to the lack of transparency into how they generate predictions. Algorithms such as Recurrent Neural Networks (RNNs) when applied to Electronic Medical Records (EMR) introduce additional barriers to transparency because of the sequential processing of the RNN and the multi-modal nature of EMR data. This work seeks to improve transparency by: 1) introducing Learned Binary Masks (LBM) as a method for identifying which EMR variables contributed to an RNN model’s risk of mortality (ROM) predictions for critically ill children; and 2) applying KernelSHAP for the same purpose. Given an individual patient, LBM and KernelSHAP both generate an attribution matrix that shows the contribution of each input feature to the RNN’s sequence of predictions for that patient. Attribution matrices can be aggregated in many ways to facilitate different levels of analysis of the RNN model and its predictions. Presented are three methods of aggregations and analyses: 1) over volatile time periods within individual patient predictions, 2) over populations of ICU patients sharing specific diagnoses, and 3) across the general population of critically ill children.

Keywords: Model Interpretation, Recurrent Neural Networks, Feature Importance, Feature Attribution, Electronic Medical Records, Deep Learning

1. Introduction

Deep learning has demonstrated promising results in a wide variety of healthcare domains including radiology [6, 10, 21], oncology [27, 14, 33], and intensive care [7, 52, 32, 3, 58]. This is due to the increasing availability of large clinical datasets such as the Electronic Medical Records (EMR) [22] and advances in computing technology that enable practical training of deep learning models [15]. The promise of deep learning is its ability to learn complex interactions directly from high-volume, high-dimensional, and multi-modal data without the need for hand selecting and engineering features specific to a modeling technique or problem [29]. Unfortunately, this flexibility comes at a price: a model with millions of parameters and hundreds of operations, opaquely optimized from large complex datasets. As a result, how a particular input feature contributes to or affects a prediction is not immediately obvious.

This lack of transparency, especially in clinical settings where decisions may be lifesaving, has inspired research efforts to interpret these highly accurate and complex models [46]. Despite this growing interest, *interpreting* a model remains a nebulous concept and is usually defined specifically for the problem and application of the model [13]. Consequently, methods for interpreting deep learning models are very diverse in method and purpose; for example, aggregating and visualizing the model’s neuron activations to extract concepts learned [39] or using a simpler model

Email addresses: loho@chla.usc.edu (Long V. Ho), maczon@chla.usc.edu (Melissa Aczon), dledbetter@chla.usc.edu (David Ledbetter), rwetzel@chla.usc.edu (Randall Wetzel)

such as a decision tree to approximate the predictions of the original model and interpreting the simpler model as a proxy [5]. The general goal of methods for interpreting is to understand the model’s *decision making process*. In this work, we use a simplified definition for *interpreting* a model – that is understanding which inputs contributed to the model’s predictions.

In particular, we are interested in determining which input features contributed to the predictions of a previously well described recurrent neural network (RNN) model that uses electronic medical data to continuously assess the status of a critically ill child based on their risk of mortality (ROM) in a pediatric intensive care unit [3]. The ability to determine how input features impact these predictions is important for several reasons.

First, it may provide useful information for clinical intervention. ICU mortality predictions for an individual patient serve as a proxy for a child’s severity of illness (SOI) [41, 54, 40, 31]. They have intrinsic value, but understanding which of the child’s features underlie the acute changes of their state, as indicated by a change in the predicted SOI, would add further value. If the clinician already knows this information, then it provides corroboration and trust, and the clinician would know what to do with it (e.g. blood pressure is important, therefore administer a therapy to optimize blood pressure). If not, then it may propel further investigations that otherwise may have been overlooked without the model’s prediction and interpretation.

Second, it facilitates an environment in which users can interact with the model and learn its strengths and weaknesses. Users can then compare the extracted input contributions with their own clinical experience [37, 13].

Third, understanding which or how inputs contributed to predictions can be used to improve the model. Combining this understanding with clinical knowledge can identify when the model improperly leverages information. Determining and correcting such characteristics are especially important in deep learning models where parameters are optimized to large biased datasets commonly found in healthcare. The bias comes from the observational nature of healthcare data, where counterfactual events do not occur. For example, if a drug is given only to terminally ill patients during end-of-life withdrawal of support, a model may inadvertently leverage the use of the drug as an indicator of mortality, which would contradict the intended purpose of the model, e.g. to find meaningful features that can be changed to improve survival probability.

The above reasons motivated the following primary goals of this work, which also describe this paper’s contributions to the still limited but growing body of literature on the methods of interpreting deep learning models applied to EMR:

- Introduce the Learned Binary Mask (LBM), a new occlusion-based method, to identify which inputs contributed to the predictions of a many-to-many RNN model that continuously generates ROM scores for an individual child from multi-modal time series EMR data. The LBM is able to manage the mixed data modalities in Electronic Medical Records.
- Modify KernelSHAP to make it compatible with a many-to-many RNN model for risk of ICU mortality whose inputs are multi-modal EMR.
- Use both the LBM and KernelSHAP to compute attribution matrices across all individual patients. Aggregate the attribution matrices in various ways for different level of analysis of the RNN model and its predictions: 1) within volatile time periods in individual patient predictions, 2) across cohorts of children diagnosed with specific diseases (sepsis and brain neoplasm), and 3) across the general population of critically ill children. These use cases of interpreting the many-to-many RNN model’s predictions with the LBM and Kernel SHAP

emphasize the importance of understanding the mathematics and assumptions of each method when applied to real data to understand the analysis of the results.

- Introduce a novel matrix representation which reflects hour-to-hour *state changes* that a patient undergoes during their stay in the ICU. This matrix is the input to the RNN model that generates dynamically evolving predictions of the patient’s ICU mortality risk (severity of illness). The state change representation enables use of the LBM and KernelSHAP.

We emphasize that this work presents the LBM and KernelSHAP as complementary, not competing, methods. Their formulations and purposes are different; therefore, one method’s results should not be regarded as better or more accurate than the other. Importantly, evaluating methods for model interpretation is inherently qualitative because there are no *hard truths* against which to compare the outputs of such methods. The evaluations rely on clinical insights and experience, which are not necessarily quantifiable. Nevertheless, these qualitative analyses and evaluations are important for the reasons stated earlier. Finally, note that this study is not about feature selection for model development.

2. Related Works

Using both RNNs *and* multi-modal EMR data poses challenges to current interpretation methods: RNNs introduce complexities associated with the time dimension while EMR data complicate comparisons among features that have different distributions and clinical meaning. Many methods for interpreting deep learning models rely on sensitivity analysis which measures feature attributions by analyzing how the prediction changes when inputs are perturbed [20]. These methods are limited to single-modal inputs such as images or text in which changes in inputs and outputs can be compared readily among features. In contrast, data in EMRs contain an eclectic collection of data modalities which cannot be trivially compared, including continuous physiology (heart rate), categoricals (Glasgow Coma Score), binary (cultures), and unstructured texts (clinical notes) [19]. In addition, many methods rely on special visualizations for presenting the extracted information, e.g. heatmaps highlighting important pixels in an image. Such visualizations for EMR data and RNNs can be intractable because even a single patient’s data can contain hundreds of different measurements and thousands of time steps.

One approach to address these issues has been to adapt the problem to current methods of interpretability. For example, Rajkomar et al. [44] converted the problem into the familiar problem of text processing by tokenizing the EMR data as single-sensor text sequences and interpreting the RNN model with methods developed for natural language processing. Another way to leverage existing interpretability techniques is to use mimic learning wherein a simpler model approximates the complex model. This approach was taken by Che et al. [5] who approximated their RNN-EMR model with a gradient boosted trees model (GBM) trained to predict the RNN’s predictions; interpreting the GBM was a proxy for interpreting the original RNN-EMR model. Such methods often require non-trivial manipulation of the data or model, and this process introduces additional layers of complexity that further muddles interpretation.

Other algorithms that use RNNs and EMR data are interpreted by embedding certain components of the algorithm with interpretable constituents. Cerna et al. [4] aggregated the response from models trained specifically on individual modalities of the EMR and interpreted the weight of a final linear layer as contributions from each of the mixed modalities. Choi et al. [8] and Zhang et al. [59] used attention RNNs to interpret their models, modifying RNNs with an attention component that imposes an explicit attribution of the inputs to the outputs via weights. Such methods

only interpret *parts* of the deep network and require complex visualization techniques to distill the information. Furthermore, Poursabzi-Sangdeh et al. [43] found that more transparent models, when compared to the same models that were presented as black-boxes, had no benefits in application and actually had detrimental effects due to information overload.

Another approach has been to use *explanation models*, which are interpretable meta-models trained in addition to the original model. Compared to other interpretability methods, explanation models are constructed to investigate specific properties of the original model (e.g. rotational invariance of the model). Both KernelSHAP and LBM fall into this category. Explanation models are specific to a particular problem. To the authors' knowledge, the only publication that uses explanation models on RNNs and EMR data is in Suresh et. al [52], which examined the impact of features by comparing differences in predictions when individual features were included or excluded. Further, the authors are not aware of KernelSHAP applications to RNN models using EMR data. This is likely due to the limitations of the current KernelSHAP implementation [48]. To facilitate discussion, the formulation of KernelSHAP and some implementation choices made for this study are described in Section 3.4.

The LBM extends the methodology of Fong & Vedaldi [17] for interpreting CNNs and images to RNNs and EMRs. In Fong & Vedaldi, the fundamental concept is to find a mask that identifies which set of pixels of an image removes evidence for being in the class of interest. Similarly, the LBM method finds a *binary* mask, instead of a real-valued mask, that identifies which set of input features when set to zero removes evidence for mortality. The LBM's significant departure from established occlusion-based methods is its *binary* constraint on the mask, which is nontrivial to obtain in practice but essential for comparing the contributions of multi-modal features in the EMR data.

The aforementioned explanation methods generate what are known as local explanations: for each individual prediction, they compute the contribution of each input feature. These local explanations can be aggregated over different predictions to provide global insights. For example, Lundberg et al.[34] aggregated local explanations across entire datasets to compute traditional model feature importance, revealing the average contribution of features and avoiding problems associated with traditional global explanations. Similarly, we extended this process by aggregating and normalizing across sub-populations. This was used to identify which features had relatively high contributions to mortality risk predictions in different diagnosis groups and the general ICU population.

3. Methods

3.1. Data

This study used de-identified EMR data collected in the Pediatric Intensive Care Unit (PICU) from Children's Hospital Los Angeles (CHLA). The CHLA Institutional Review Board (IRB) reviewed the study protocol and waived the need for IRB approval. The data consisted of 9855 PICU encounters (7358 patients) from 2009 to 2017 (4% mortality rate), where an encounter is defined to be a patient's contiguous stay in the ICU. A patient can have multiple ICU encounters. Data for each patient encounter included irregularly sampled physiologic observations, laboratory results, drugs, and interventions (e.g. intubation parameters). Also collected were the patient's demographics, diagnoses, and outcomes (e.g. ICU mortality). The encounters were partitioned into training (60%), validation (20%), and test (20%) sets, where the partitioning was done such that all encounters from a single patient belonged to only one of the sets. Statistics of the datasets are included in Appendix A.

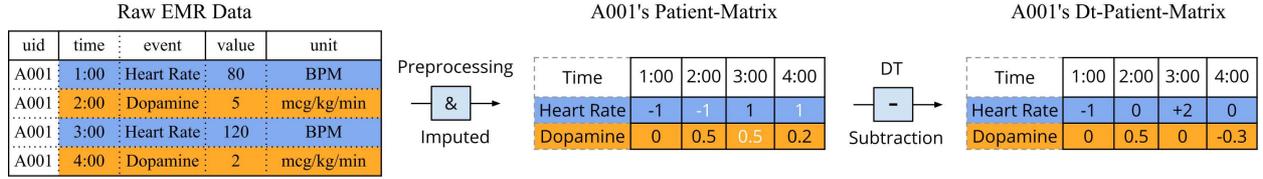


Figure 1: Pre-processing from raw EMR data to the *dt-patient-matrix*. Raw EMR Data is converted to a *patient-matrix* by pivoting the data to wide-format, imputing missing data with forward fill, converting physiologic variables to z-scores using means and standard deviations (Heart Rate with $(\mu, \sigma) = (100, 20)$), and normalizing exogenous variables such as drugs to $[0, 1]$ using minimum and maximum values (Dopamine with $(\min, \max) = (0, 10)$) computed from the training set. Note the normalization parameters used here were chosen for illustrative purposes. Finally, data from the *patient-matrix* is converted to the *dt-patient-matrix* by subtracting values from subsequent time-steps. Text colored in white (in A001's Patient-Matrix) indicates imputed data.

3.1.1. Patient Representation: The “*patient-matrix*”

Pre-processing steps converted the EMR data to matrices that facilitated machine learning while maintaining an interpretable patient representation. In collaboration with physicians, lower and upper limits for all variables were defined, and the entries were curated to remove observations not compatible with life. Values from different methods of measurements of the same variable were combined into a single variable when medically appropriate (e.g. invasive and non-invasive blood pressure). Physiologic variables such as labs and vitals were converted to z-scores using means and standard deviations computed from the *training set*. In the z-score space, the features associated with physiologic variables represent how *far* the patient’s physiologic measurements are from the ICU *averages*, and this distance is measured in terms of standard deviations. Exogenous variables such as interventions and drugs were normalized to values from 0 to 1 using the lower and upper limits of values/dosages computed from the *training set*. Consequently, features associated with exogenous variables in the transformed space represent a percentage of therapies administered in the ICU, with 0 indicating no therapy and 1 denoting maximum therapy possible within the dataset.

Each patient’s data were forward-filled and re-sampled to hourly observations, then pivoted to form a sparse $N \times T$ *patient-matrix*, where T is the number of distinct timestamps (varied across encounters), and N is the total number of input features (see Appendix Appendix A for a list of the $N = 398$ features). The training set mean was used to impute a physiologic or laboratory variable at all times prior to its first recorded value or at all times of the episode if it had no recorded value for the entire episode. Any missing measurement of a treatment indicates actual absence of that treatment; therefore it was set to zero. Note that patient diagnoses were not included as input features. A single column of this matrix contains measurements from all N features at a single time point, while a single row contains measurements of a single feature from all T time points. A similar matrix representation of patient encounter data, but without the hourly re-sampling, was used in previous work [3, 23, 28, 58].

3.1.2. Patient Representation: The “*dt-patient-matrix*”

To facilitate the application of both LBM and KernelSHAP, we further transformed the *patient-matrix* to represent the changes between each time step: the *dt-patient-matrix*. The *patient-matrix* was transformed to *dt-patient-matrix* with these steps: the first column of the *dt-patient-matrix* is exactly the same as the first column of the original patient matrix, which represents the patient’s state relative to the ICU population encapsulated in the training set. Values in subsequent columns of the *dt-patient-matrix*, each indexed by τ , were obtained by subtracting the patient state at time $\tau - 1$ from the state at time τ . Figure 1 illustrates the pre-processing steps from raw EMR data to the *dt-patient-matrix*. The *dt-patient-matrix* captures the hour-to-hour changes – physiologic and therapeutic – that a patient undergoes

during an ICU encounter.

Because sparsely measured features such as laboratory tests and rare treatments were included, pivoting the long-format EMR data to the original *patient-matrix* representation introduced a sparse matrix with over 94% of elements set to zero (measured across the 398 selected observations). In the *dt-patient-matrix* representation, the forward-filled values became zeros, indicating either no change from the population average (if the element is in the first column) or no change from the individual’s previous state (subsequent columns). This representation is consistent with the collection of data in clinical practice: new observations are typically recorded during state changes and are assumed to be the previous value until a new recording or entry is made [11, 12].

Importantly, this state change representation provides significant and distinctive advantages when using the LBM or KernelSHAP. The *dt-patient-matrix* removes ambiguities when occlusion-based methods of interpretation are used. To “delete” (i.e. set to zero) an element precisely means to have no change in that variable. This representation also facilitates practical use of KernelSHAP: the zero-valued elements of the *dt-patient-matrix* define a set of “missing” features, and having this set bypasses expensive computational steps otherwise required. These points are further discussed in Sections 3.3 and 3.4.

3.2. RNN Model for ICU Mortality

An RNN model using the *dt-patient-matrix* as input was trained to predict ICU mortality. Given a patient encounter, the model generates a risk of mortality (ROM) each time it receives new data (i.e. a single column from the *dt-patient-matrix*); see Figure 2A. The model is composed of three stacked Long Short-Term Memory (LSTM) [24] layers with hidden units 128, 256, and 128 respectively and a final logistic regression layer for classification. Each layer’s weights were initialized using Glorot uniform [18] and optimized using RMSprop [56]. Using an initial learning rate of $1e^{-4}$ and mini-batch size of 128, the model was trained to minimize the binary cross-entropy between the model’s predicted ICU-mortality risk and the patient’s true mortality response, repeated for every time-step. Performance on the validation set was evaluated after every epoch (a full cycle through the training set), and the best performing weights were saved. If the validation set’s binary cross-entropy did not decrease after 15 epochs, learning rate was decreased by a factor of 5 and terminated after 2 reductions. Model regularization included a 20% dropout of the output of each layer and an L2 penalty of $1e^{-5}$ against each layer’s weights. The Python package Keras [9] and the TensorFlow [1] backend were used to construct and train the model. Training the full model using a Titan V GPU took approximately 6 hours to complete.

The RNN model’s performance was evaluated by computing the area under the ROC curve (AUC) across the hold-out test set at various times: 1, 3, 6, 9, 12, -12, -9, -6, -3, -1, where a positive number indicates hours since ICU admission and negative indicates hours until ICU discharge (or death). In this test formulation, data up until $t - 1$ is given as input to generate a ROM prediction which is compared with the outcome at t . For example, data from $t = 0$ to $t = 11$ is used to evaluate the AUC at $t = 12$. This is presented in Section 4.1. Clinically used risk of mortality models, PIM2 [50] and PRISM3-12 [42], were also evaluated on the same test set to provide comparators. Note that PIM2 and PRISM-3 are static models (using logistic regressions) and use data within 12 hours of ICU admission to generate a single prediction per patient encounter. Consequently, their performances can only be fairly compared with the RNN’s 12th hour predictions. To ensure proper comparison of AUCs of the RNN predictions at different hours, patients with less than 24 hours of ICU data were excluded from the AUC computations. The AUC was chosen as the metric for performance evaluation because it is not sensitive to class imbalance. Preserving the class imbalance (96% vs 4%) during both model training and assessment is important because this imbalance would be present during actual deployment and would inform a full analysis of deployment benefits and costs.

3.3. Learned Binary Masks

The goal is to find a binary mask for each input element in the dt-patient matrix such that when applied, the RNN model predictions go to zero. See Figure 2B for a conceptual illustration and Appendix B for an algorithmic description. The reasoning for this formulation is analogous to the premise behind occlusion-based methods for interpreting convolutional neural networks that classify images: find the set of pixels such that when they are ‘deleted’ or masked, the class probability goes from non-zero to zero [17]. The binary requirement stems from the multi-modal nature of EMR data which include real-valued, integer-valued, categorical, and binary features. Multiplying a binary or integer-valued feature by a fractional mask value can lead to a clinically unrealistic transformation. For example, multiplying the variable indicating whether or not a chest X-Ray was taken by a 0.5 mask has no meaning. Similarly, a fractional multiplier for Glasgow Coma Score, which can only take on integer values from 3 to 15, is not realistic. For generalizability (ie. to accommodate models that use potentially different sets of EMR variables), the method that generates the weights needs to be independent of the exact structure of the inputs. This dictates the binary nature of the resulting weights.

We focus on the RNN ICU mortality model described in Section 3.2. Each patient encounter results in a matrix of binary masks indicating whether elements in the encounter’s dt-patient-matrix elements contributed to the trajectory of ROM predictions for that encounter. The formulation of the input dt-patient-matrix means that assigning a zero weight for a particular feature at a specific time step is equivalent to requiring no recorded change in that feature from the previous to the specified time step.

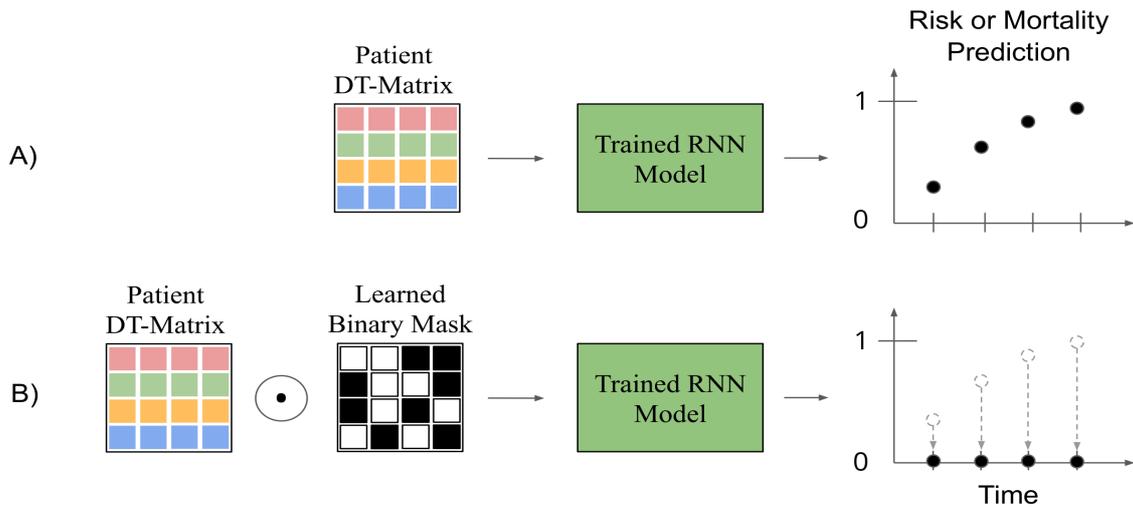


Figure 2: A) Every hour after ICU admission, a many-to-many RNN model acquires new patient data (single column of the patient dt-matrix) and generates a mortality prediction (Section 3.1.2). The Patient DT-matrix is color coded to denote the multi-modal nature of the input data: vitals in blue, labs in green, interventions in orange, and drugs in red. B) The Learned Binary Masks algorithm, detailed in Section 3.3 and Appendix B, computes a matrix of binary-valued weights such that when these weights multiply the input data, the RNN’s risk of mortality predictions go to zero at all time steps. The zeros in the binary matrix identify which of the input Patient DT-matrix elements made the ROM predictions go to zero. This is analogous to occlusion-based methods in images that find the set of pixels such that when they are ‘deleted’ or masked, the class probability goes from non-zero to zero [17].

3.3.1. Mathematical Formulation

The RNN ICU mortality model and derivation of LBM are mathematically formalized to facilitate understanding of how the LBM determines feature attributions. The mortality model is a many-to-many recurrent neural network f

with trained parameters Θ that maps $\mathbf{x}_{1:T} \in \mathbb{R}^{N \times T}$ (*dt-patient-matrix* with T time points and $N = 398$ features) to a T -length sequence of probabilities of mortality:

$$f(\Theta; \mathbf{x}_{1:T}) = [y_1, y_2, \dots, y_T] \equiv \mathbf{y}_{1:T}. \quad (1)$$

At any single time, τ , of an individual patient encounter, the notation $y_\tau = 0$ corresponds to zero risk of mortality, while $y_\tau = 1$ corresponds to 100% risk of mortality.

Let $\mathbf{0}_{1:T}$ denote a T -length sequence of zeros. If $\mathbf{y}_{1:T} \neq \mathbf{0}_{1:T}$, i.e. the prediction is not zero at every time step, then LBM will find a sparse, binary mask $\mathbf{M}_{1:T} \in \{0, 1\}^{N \times T}$ that satisfies

$$f(\Theta; \mathbf{M}_{1:T} \odot \mathbf{x}_{1:T}) = \mathbf{0}_{1:T}, \quad (2)$$

where \odot denotes the Hadamard product (element-wise matrix multiplication). This is illustrated in Figure 2B. If \mathbf{M} solves Equation 2, then the locations of \mathbf{M} 's zeros correspond to the elements of the input *dt-patient-matrix* that must be zeroed for the predictions to go to zero. If the inputs identified by \mathbf{M} are not zeroed (i.e. input x_τ^j retains its original non-zero value), then the ROM prediction at τ remains non-zero. Therefore, a solution with $M_\tau^j = 0$ indicates that the recorded change in feature j from time $\tau - 1$ to time τ contributed to the non-zero ROM prediction at time τ , i.e. the recorded change provided evidence for mortality.

The mask $\mathbf{M}_{1:T}^*$ that satisfies Equation 2 can be found by leveraging the same fundamental mechanics used to train the neural network – minimization of a regularized loss function through backpropagation. Mathematically, this is expressed as

$$\mathbf{M}_{1:T}^* = \underset{\mathbf{M}_{1:T}}{\operatorname{argmin}} f(\Theta; \mathbf{x}_{1:T} \odot \mathbf{M}_{1:T}) + R(\mathbf{M}_{1:T}), \quad (3)$$

where the first term, $f(\Theta; \mathbf{x}_{1:T} \odot \mathbf{M}_{1:T})$, ensures that the mask causes the predictions to go to 0, while the second term, $R(\mathbf{M}_{1:T})$, imposes sparsity on the mask via L_1 regularization: $R(\mathbf{M}_{1:T}) = \lambda_1 \|\mathbf{1} - \mathbf{M}_{1:T}\|_1$. This functional form for R pushes many of the mask entries to unity, which is equivalent to minimizing the number of zero entries in the mask. This means that only those features that provided evidence for non-zero ROM, i.e. mortality, are selected.

The non-differentiable nature of the binary constraint on $\mathbf{M}_{1:T}^*$ poses a challenge to optimization methods that leverage gradients such as backpropagation. A two-step process addresses this issue: (1) use backpropagation to solve Equation 3 for an intermediate non-binary mask, $\mathbf{m}_{1:T} \in [0, 1]^{N \times T}$; and (2) find another non-binary mask, $\boldsymbol{\eta}_{1:T} \in [0, 1]^{N \times T}$ to binarize $\mathbf{m}_{1:T}$. The final binary mask is given by $\mathbf{M}_{1:T} \equiv (\mathbf{m}_{1:T} > \boldsymbol{\eta}_{1:T})$.

The first mask $\mathbf{m}_{1:T}$ is found by applying backpropagation to a modified version of Equation 3:

$$\mathbf{m}_{1:T}^* = \underset{\mathbf{m}_{1:T} \in [0, 1]^{N \times T}}{\operatorname{argmin}} f(\Theta; \mathbf{x}_{1:T} \odot \mathbf{m}_{1:T}) + \lambda_1^1 \|\mathbf{1} - \mathbf{m}_{1:T}\|_1 + \lambda_2^1 H(\mathbf{m}_{1:T} > 0.5, \mathbf{m}_{1:T}), \quad (4)$$

where H is the binary cross-entropy function, and λ_1^1 and λ_2^1 are regularization constants. The first two terms parallel those in Equation 2, while the third term, $\lambda_2^1 H(\mathbf{m}_{1:T} > 0.5, \mathbf{m}_{1:T})$, encourages the values of $\mathbf{m}_{1:T}^*$ to be closer to 0 or 1, i.e. pushes $\mathbf{m}_{1:T}^*$ to be *near-binary*. Next, let $\mathbf{m}_{1:T} \equiv \sigma(A \times \mathbf{z}_{1:T})$, where A is a constant, σ is the sigmoid function $\frac{1}{1+e^{-x}}$, and $\mathbf{z}_{1:T} \in \mathbb{R}^{N \times T}$; and similarly define $\mathbf{m}_{1:T}^*$ from $\mathbf{z}_{1:T}^*$. Optimizing for $\mathbf{z}_{1:T}^*$ instead of $\mathbf{m}_{1:T}^*$ removes the $[0, 1]$ range constraint during the backpropagation process.

Next, the threshold mask $\boldsymbol{\eta}_{1:T}$ is found through a brute-force grid search of threshold values $\eta_t \in [0, 1]$ to minimize

$$\boldsymbol{\eta}_{1:T}^* = \underset{\boldsymbol{\eta}_{1:T}}{\operatorname{argmin}} f\left(\Theta; \mathbf{x}_{1:T} \odot \left(\mathbf{m}_{1:T}^* > \boldsymbol{\eta}_{1:T}\right)\right) + \lambda_1^2 \left\| 1 - \left(\mathbf{m}_{1:T}^* > \boldsymbol{\eta}_{1:T}\right) \right\|_1, \quad (5)$$

where λ_1^2 is a constant governing sparsity in the final mask $\mathbf{m}_{1:T}^* > \boldsymbol{\eta}_{1:T}$. Two simplifications make the brute-force optimization efficient while maintaining realistic representations. First, because the mask from the first optimization step is sparse and near-binary, the grid-search area can be limited to small regions around unique values found in $\mathbf{m}_{1:T}^*$. Second, the grid search can also be limited to optimizing *backwards* through time, applying the same threshold across all features at a given time, heavily reducing grid search from $\boldsymbol{\eta}_{1:T} \in [0, 1]^{N \times T} \rightarrow \eta_{1:T} \in [0, 1]^T$. This also urges the mask to maintain clinical validity by ensuring that binarization of features at times $t \leq \tau$ does not affect the optimization for $t > \tau$. Finally, the binary mask is obtained by defining

$$\mathbf{M}_{1:T}^* = \left(\mathbf{m}_{1:T} > \boldsymbol{\eta}_{1:T}\right). \quad (6)$$

3.3.2. Implementation Details

An algorithmic description summarizing LBM’s two-step process for generating a binary mask for an individual patient can be found in Appendix B. The LBM was implemented by adding to the trained RNN model an additional layer that multiplies the input data with a trainable mask (first term of the cost function in Equation 3). This enabled the leveraging of the same mechanics and infrastructure that were used to construct and train the RNN model. The trainable mask was initialized to all ones. Equation 4 was minimized using RMSProp [56] with an initial learning rate of 0.1. LBM hyperparameters were experimentally determined through trial and error and were chosen using qualitative examinations of the masks based on sparsity of the resulting mask and binarization of the first mask obtained from the first optimization step. The parameters were set as follows: $\lambda_1^1 = 0.005$, $\lambda_2^1 = 0.0005$, $A = -5$. The learning rate was reduced by a factor of 10 if the loss function did not improve after 5 iterations. Optimization was terminated when either 5000 iterations were reached or the training loss did not improve after reducing the learning rate 2 times. Equation 5 was optimized using brute-force grid search of threshold values backwards through time. Optimization was terminated when either each entry of the matrix $f(\Theta; \mathbf{x}_{1:T} \odot \mathbf{m}_{1:T})$ was less than 0.05 or three iterations of searching back through time were reached. Although there is no theoretical guarantee for the existence of a unique solution, our experiments indicate that the two-step approach was successful in finding a non-trivial binary mask that satisfies Equation 3 across all RNN patient predictions.

3.4. KernelSHAP

Introduced in the 1950s, Shapley values answer a question from cooperative game theory: if N players cooperated with each other for a collective reward, then how is the reward distributed fairly to each of the N players [49]? Equivalently, what is the marginal contribution of each player? Consider the payout when a subset of S players cooperate with each other without player j and the payout when these same S players cooperate with player j , then take the difference between these two payouts. This process is repeated over all possible 2^{N-1} subsets of players without player j , and the average of the resulting payout differences is the Shapley value for player j . Implicitly assumed in the formulation is the ability to observe the payout for each of the scenarios. The individual Shapley values sum to the collective reward.

This game theory principle was adapted to compute the “payout” of input features to a model’s prediction, i.e. contributions of input variables. The N players are the N specific feature values at a particular input data instance that

generated a specific prediction y^* , i.e. $f(x_1^*, \dots, x_N^*) = y^*$, where f denotes the model. Missing players correspond to input features whose values are unknown, and the payout in this scenario is the expectation of f over all the possible values that these features could take. The required computations have exponential time complexity and render the method impractical. Approximation methods, including sampling techniques described in [51], make the method tractable.

KernelSHAP incorporates the game theory principle of Shapley values with an existing model interpretation method, LIME (Local Interpretable Model-agnostic Explanation). The LIME framework finds a low-complexity function g that approximates f around a given point [45]. The inputs to g are $\{z_j\}_{j=1}^m$, where z_j is an interpretable combination of the input features $\{x_j^*\}_{j=1}^N$. For the purposes of this paper, $z_j = 1$ if the value of feature j is known, and $z_j = 0$ otherwise. In other words, z_j is a toggle for the presence or ‘missingness’ of feature j ; hence $m = N$. Let $h_x(z')$ denote the inverse mapping of $z'(x)$. If $z'_j = 1$, then $(h_x)_j = x_j^*$; if $z'_j = 0$, then $(h_x)_j$ is unknown or random. For g to locally approximate f , then $f(h_x(z')) = g(z')$ when $h_x(z')$ is near the data instance x . This means that if g is linear, then in the neighborhood of a specific data instance, x^* , f can be written as:

$$f(x^*) = \phi_0 + \sum_{j=1}^N \phi_j z'_j(x^*). \quad (7)$$

Note that if none of the feature values are known (i.e. $z'_j = 0$ for all j), then $f = \phi_0$. This means that ϕ_0 is the expected value of f over the entire space, i.e. the “background value.” If all of the feature values are known ($z'_j = 1$ for all j), then $f(x^*) - \phi_0 = \sum_{j=1}^N \phi_j$. This says that the summation on the right is how much the model prediction changes from the background value if all N input features take on the values in the data instance x^* . Lundberg and Lee showed that the coefficients ϕ_j ($j \geq 1$) are exactly the Shapley values if they solve the weighted least squares problem given by:

$$\underset{\phi \in \mathbb{R}^{N+1}}{\operatorname{argmin}} \sum_{z' \in \{0,1\}^N} \pi_x(z') \left[f(h_x(z')) - \left(\phi_0 + \sum_{j=1}^N \phi_j z'_j \right) \right]^2, \quad (8)$$

where $\pi_x(z')$ is the Shapley Kernel [36]. The term $f(h_x(z'))$ is the expected value of f conditioned on z' : if z' has ones in positions i, j, k and zeros everywhere else, then $f(h_x(z')) = E[f | x_i = x_i^*, x_j = x_j^*, x_k = x_k^*]$, where x_i^*, x_j^*, x_k^* are the actual values of features i, j, k in the specific input (data instance) that generated the prediction of interest. One can think of Equation 8 as an over-determined system of linear equations relating the ϕ_j ’s with the expected values of f when different combinations of features are unknown while the rest take their values from the data instance. The Shapley Kernel, $\pi_x(z')$, places higher weights to those equations corresponding to z vectors with either a very small or very large number of ones (i.e. a very small or very large number of features inherit their values from the given data instance x). Since each element of z' is either 0 or 1, then the outer summation has 2^N terms, which is the number of all possible z' vectors. Further, for all but one of these vectors, z' has at least one zero element, and computing $f(h_x(z'))$ involves taking the expectation of f over the corresponding subspace of feature values.

For practical use, implementations of Equation 8 must address two main problems: (1) efficiently estimate the expectation $f(h_x(z'))$ for a given z' ; and (2) reduce the number of terms in the outer summation. The first requires a background dataset from which to draw random values for x_i when $z'_i = 0$. For some classes of f , e.g. tree-based models, fast versions for estimating the expectation have been proposed, e.g. Tree SHAP [35]. KernelSHAP is the existing implementation of the general (model-agnostic) case, and in some of its applications, $f(h_x(z'))$ is approximated by evaluating f only once using a data vector where any ‘unknown’ feature value (ie. where $z_j = 0$) is

set to what is considered a ‘normal’ value for that feature, i.e. its median or mean [48]. Note that Equation 8 must be solved for each prediction that needs to be explained. In practice, if there are T predictions, then the T least squares problems are solved simultaneously. The total number of random z' samples for the collective outer summation is usually set to $2NT + 2048$ [48]. These samples are concentrated in regions where the the number of zeros in z' is very small or large since the kernel gives these regions higher weights.

Our implementation of KernelSHAP for the RNN mortality model computes the expectation $f(h_x(z'))$ by setting $x_\tau^j = 0$ when $z'_j = 0$ in the outer summation of Equation 8. This single-point evaluation simplifies the computation of $f(h_x(z'))$ and significantly reduces the run time. Recall that $x_\tau^j = 0$ means the measurement for feature j either did not deviate from the population mean (at the first timestep, $\tau = 0$) or did not change from the previous timestep ($\tau \geq 1$). The latter typically resulted from having no new measurements in feature j ; clinically, this meant the patient was considered stable [47]. Therefore, our expectation for f at a sample z' is equal to what the model would have predicted when the specified group of features (identified by zeros in z') had no new measurements. At the time of our implementation, the existing KernelSHAP libraries did not properly compute Shapley values for multi-modal time-series data such as EMR, so additional modifications were made as described in [16].

3.5. Interpreting RNN Predictions

This section provides three different use cases demonstrating how the LBM and KernelSHAP can be used to determine the contributions of input features to the RNN model’s ICU mortality predictions and identify which were most important at various scales or levels. These levels are 1) for an individual, 2) different subgroups defined by disease processes, and 3) the general ICU population. The common theme across the different levels is the repeated averaging of local information over specified time periods of a group of individuals. Section 4 will illustrate with specific examples.

3.5.1. Individual Attribution Matrices

Given the RNN model’s sequence of mortality risk predictions for an individual patient’s entire encounter, the LBM and KernelSHAP each generate an attribution matrix describing the contributions of input elements to those predictions. The attribution matrix, denoted by $\mathbf{a}_{1:T}^p \in \mathbb{R}^{N \times T}$ for patient encounter p , has the same dimensions as the encounter’s dt-patient-matrix, $\mathbf{x}_{1:T}^p$. For retrospective analysis presented in this study, T is the final hour before ICU discharge. Recall that N is the number of input features at each timestep.

The LBM attribution matrix is given by $\mathbf{a}_{1:T}^p = 1 - \mathbf{M}^*_{1:T}$, where $\mathbf{M}^*_{1:T}$ is the mask defined by Equation 6. Since $\mathbf{M}^*_{1:T}$ is binary, then $\mathbf{a}_{1:T}^p$ is also binary. For KernelSHAP, Equation 8 must be solved for each timestep $\tau \in [0, T]$, and the resulting coefficients ϕ_τ^j comprise the elements of the KernelSHAP attribution matrix. The different formulations of the LBM and KernelSHAP mean that they generate different attribution matrices that can provide complementary perspectives.

- The LBM answers the question, “Which of the non-zero elements of an individual’s input dt-patient matrix led to the non-zero mortality predictions for this individual in $[0, T]$?” The locations of ones in the LBM attribution matrix identify which of the non-zero entries in the dt-patient matrix were ‘zeroed’ by the LBM to drive all the ROM predictions for patient encounter p to zero. Equivalently, these non-zero changes in feature measurements from one timestep to the next provided evidence for the non-zero ROM predictions. For example, if j corresponds to heart rate, and $a_\tau^j = 1$, then the non-zero change in heart rate (e.g. decrease of 10 beats per minute) from time $\tau - 1$ to τ led to the non-zero mortality predictions. The LBM attribution matrix does not

describe how the prediction would change if the heart rate had increased by 30 bpm instead of decreased by 10 bpm. Neither does the LBM describe what would happen to the prediction if the heart rate had increased by 20 bpm instead of remained at the same value. It is important to note that the LBM will not highlight any input element x_j^τ that is already zero because setting the mask m_τ^j to zero for such elements does not change the value of f in Equation 4 but increases the regularization term $\|1 - \mathbf{m}_{1:T}\|$ in the loss function. Consequently, the LBM attribution matrix is sparse because the input dt-patient-matrix is sparse.

- KernelSHAP expresses the prediction at time τ as a sum of contributions from the inputs at that time: $y_\tau = \phi_0 + \sum_j \phi_\tau^j$. With ϕ_0 denoting what the prediction would have been if no new measurements were recorded at that time, then the attribution matrix element $a_\tau^j \equiv \phi_\tau^j$ (which can be positive or negative) reflects how much would be added to ϕ_0 given that the recorded change for feature j at that time was in fact x_τ^j . All the inputs at time τ cooperate with each other to generate y_τ , and a_τ^j is the marginal contribution from knowing that the measurement for feature j changed by exactly x_τ^j between $\tau - 1$ and τ . As a simple example, if the actual prediction at time τ was 0.45 and the background value was 0.15, then KernelSHAP expresses the difference, 0.30, as a sum of marginal contributions from the inputs: 0.25 was due to heart rate decreasing by exactly 10 bpm, -0.15 to systolic BP increasing by 30 mmHg, and 0.20 to the epinephrine dose remaining at 0.04 mcg/kg/minute from hour $\tau - 1$ to hour τ . Unlike the LBM, KernelSHAP’s a_τ^j can be non-zero even when x_τ^j is zero. KernelSHAP does not answer the question: “if x_τ^j had been different from its current value (regardless of what that current value is), then how would the prediction change?”

The attribution matrices generated by the LBM and KernelSHAP for a single patient encounter contain a lot of information. They can be aggregated in many different ways to facilitate different levels of analysis and enable display of more concise information. Below we describe three aggregation techniques that each answer a different question:

- Which features were important during an individual patient’s period of rapidly changing ROM predictions?
- Which features were important to the RNN’s ROM prediction in a specific diagnostic cohort of ICU patients *relative* to the rest of the ICU population without that diagnosis? Equivalently, which features contribute to risk of mortality in a specific diagnostic group relative to the other ICU population?
- Which features were important to the ROM predictions of the general population of critically ill children?

3.5.2. Identifying Important Features for Individual Patient Predictions

The goal is to identify which clinical features contributed to a critical period of illness of an individual patient, as indicated by rapidly changing ROM scores. Hence, we average the individual’s attribution matrix over the time window of increasing ROM scores. In a deployment setting, a clinician could highlight any time window of interest to understand which features were contributing to the RNN ROM predictions. The averaging process reduces the amount of information from an $N \times T$ matrix to an N -dimensional vector of real values reflecting which features contributed, on average, to the ROM predictions *in that time window*.

Let $[t_i, t_f]$ denote a time interval of interest, e.g. where the ROM prediction for an individual patient increases significantly. The attribution matrix $\mathbf{a}_{1:T}^p$ can be regarded as a sequence of N -dimensional attribution vectors, \mathbf{a}_τ^p , with τ denoting time. The vectors corresponding to the time points in $[t_i, t_f]$ are averaged as follows:

$$\bar{\mathbf{a}}^p = \frac{1}{(t_f - t_i)} \sum_{\tau=t_i}^{t_f} |\mathbf{a}_\tau^p|, \quad (9)$$

where the absolute value is applied element-wise. For the LBM, the j^{th} element of the resulting N -dimensional vector, $\bar{\mathbf{a}}^p$, reflects how often the j^{th} input feature provided evidence for mortality within the specified time interval. For KernelSHAP, this element reflects how much feature j , on average during that time interval, changed the prediction (in magnitude) from the background value.

The vector $\bar{\mathbf{a}}^p$ can be normalized so that its largest element is unity, i.e. the feature that contributed the most is assigned a value of one:

$$\bar{\mathbf{a}}^p \leftarrow \frac{\bar{\mathbf{a}}^p}{\|\bar{\mathbf{a}}^p\|_\infty}. \quad (10)$$

Section 4.2 illustrates with examples the computation of $\bar{\mathbf{a}}^p$ for two patient encounters from the test set. The two encounters were chosen to have different diagnoses but similar predicted ROM trajectories with a single inflection point, going from a low predicted ROM at ICU admission to a high predicted ROM near the end of the ICU stay. Two additional encounters are similarly analyzed in Appendix C, which were selected because their ROM predictions similarly had only one time period of substantial change and because they survived their ICU stay.

3.5.3. Average Feature Attributions Within Cohorts: Relative Attribution Features

Instead of averaging only over specific time windows within individual attribution matrices, we can average over their entire ICU stay. This can further be done over a specific group of patients to understand the top contributing features for that group. As before, temporal averaging reduces each patient encounter’s attribution matrix to a vector (parallel to Equation 9). Next, these vectors are averaged over a group of encounters. When the group is defined by a disease process or diagnosis, this aggregation process is akin to traditional clinical research wherein logistic regression models are used to determine risk factors for that disease. For example, we can identify the features affecting mortality predictions in sepsis patients (here denoted by the set S) by computing the average absolute value of the attributions across patients whose primary diagnosis was sepsis:

$$\bar{\mathbf{a}}_S = \frac{1}{n(S)} \sum_{p \in S} \frac{1}{t_p - t_0} \sum_{\tau=t_0}^{t_p} |\mathbf{a}_\tau^p|, \quad (11)$$

where t_0 and t_p are the first and last time points for patient p , and $n(S)$ is the number of patients in S . A similar computation was performed over all patients without a sepsis diagnosis, denoted by S^c , to yield $\bar{\mathbf{a}}_{S^c}$. Both $\bar{\mathbf{a}}_S$ and $\bar{\mathbf{a}}_{S^c}$ are N -dimensional vectors, where $N = 398$ is the number of features used by the RNN model at each hour of prediction. While $\bar{\mathbf{a}}_S$ highlights which features were important in predictions for sepsis patients, $\bar{\mathbf{a}}_{S^c}$ identifies which features were pertinent in non-sepsis patient predictions. Normalizing them to unit-length vectors and then subtracting one from the other yields

$$\tilde{\mathbf{a}}_S = \frac{\bar{\mathbf{a}}_S}{\|\bar{\mathbf{a}}_S\|} - \frac{\bar{\mathbf{a}}_{S^c}}{\|\bar{\mathbf{a}}_{S^c}\|}. \quad (12)$$

We refer to the elements of the resulting N -dimensional vector, $\tilde{\mathbf{a}}_S$, as the *relative attribution features* (RAFs) because they describe which features affect mortality prediction among sepsis patients more than they do non-sepsis ones. If the j^{th} element of $\tilde{\mathbf{a}}_S$ is positive, then feature j was more important in predicting mortality for sepsis patients than in non-sepsis patients. RAFs can be computed between two cohorts to determine which features affect mortality more in one cohort relative to the other. Two examples were explored: between patients with and without sepsis, and between patients with and without brain neoplasm.

3.5.4. Feature Contributions in All Critically Ill Children

We are also interested in the top contributing features for the general population including all critically ill children irrespective of diagnosis. For this, the aggregation in Equation 11 is done over all patient encounters in the (test) set instead of only a subset of them, and the resulting vector normalized in a manner similar to Equation 10:

$$\bar{a}_S \leftarrow \frac{\bar{a}_S}{\|\bar{a}_S\|_\infty}. \quad (13)$$

For the LBM, the j^{th} element of \bar{a}_S corresponds to the relative frequency that non-zero changes in feature j (e.g. Heart Rate) led to non-zero ROM predictions in the general ICU population. It is “relative” in the sense that the most frequent feature is assigned a value of unity. For KernelSHAP, this value reflects the average contribution, in magnitude, of feature j to the mortality probabilities for the population. In either case, \bar{a}_S shows population-level feature contribution and can be considered as the “weights” of the RNN, similar to the weights of a logistic regression developed from the entire population. The vector \bar{a}_S was computed using both the LBM and KernelSHAP attribution matrices of *test* set encounters and is presented in Section 4.4. This computation of feature importance differs from standard *permutation* feature importance computations which rank features based on their effects on model performance. Feature importance computations using the LBM and KernelSHAP attributions are derived from examining specific properties for each individual prediction; they are not optimized to improve model performance.

3.6. Compute Time

Lastly, to test the real-time clinical deployment potential of each method, the time-to-compute per patient was compared and is presented in Section 4.5. Also included for comparison is the time-to-compute for generating RNN predictions. Because length of stay in the ICU varies greatly (and therefore the amount of time-steps in the patient data), it is expected that computing attribution matrices may also vary. Timings were done using a cuda-enabled NVIDIA Titan V GPU and python 2.7 with Keras and Tensorflow backend.

4. Results

4.1. RNN Model Performance

Figure 3 shows the performance of the model at critical times in the ICU: the first twelve hours after admission and the last twelve hours before discharge (or death). At the 12th hour after admission, the RNN (0.93) significantly outperformed PIM2 (0.86) and PRISM3-12 (0.88). From ICU admission to discharge, the RNN’s AUC improved over time, approaching 1 at the end of stay. Over time as the RNN accrued more data and learned more about the patient and as lead time decreased [30], the model’s predictions became more accurate.

4.2. Interpreting Individual Patient Predictions

The LBM and KernelSHAP were applied to the RNN model’s ROM predictions for all encounters in the test set. The attribution matrices are analyzed and presented here for two individual patients who subsequently died. These two patient encounters were selected for individual analysis because their ROM predictions contained only one time window of substantial change from low to high mortality risk. Patients in the ICU can undergo multiple periods of volatility, as reflected in multiple shifts in the trajectory of ROM predictions, and limiting the patient-level analysis to encounters with only one such period simplifies the interpretation of results. See Appendix C for analysis of two additional patients, which were selected similarly but with an additional criteria of surviving their ICU stay. The first

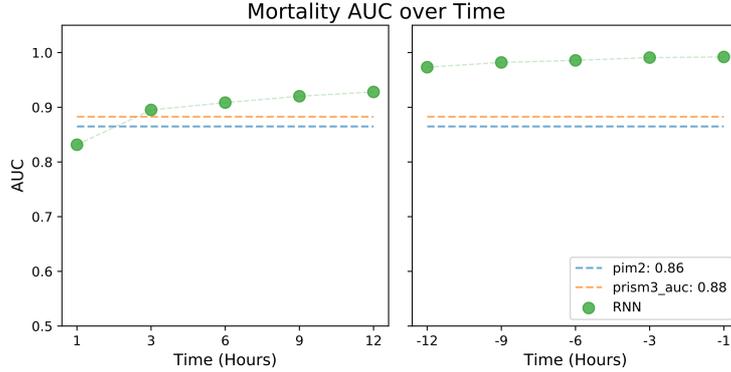


Figure 3: Performance of the RNN model’s predictions at significant times in the ICU

patient, p_1 , was an infant male weighing 6 kg and diagnosed with pneumonia and acute respiratory distress syndrome (ARDS) [57] caused by whooping cough (pertussis), severely affecting the respiratory system. The second patient, p_2 , was a 3 year old female weighing 27 kg with a diagnosis of Grand mal status epilepticus [2], a condition of continuous seizures, caused by intracranial hemorrhage (bleeding in the cranium with elevated intracranial pressure) and ARDS.

The RNN’s ROM predictions are shown in Figure 4A-a for p_1 and 4A-b for p_2 , while the attribution matrices are shown in Figure 4B (KernelSHAP) and 4C (LBM) as heatmaps, with time on the x-axis and features on the y-axis. In both patients, ROM started low but quickly increased over a 15 hour window (shaded region in Figures 4A-a and 4A-b) to a high value. Equation 9 was used to identify which features were responsible for this ROM increase. In other words, the attribution matrices were averaged across the time periods of interest: $\bar{\mathbf{a}}^{p_1} = \sum_{t=55}^{70} \mathbf{a}^{p_1}(t)/(70 - 55)$ and $\bar{\mathbf{a}}^{p_2} = \sum_{t=40}^{55} \mathbf{a}^{p_2}(t)/(55 - 40)$, then normalized according to Equation 10. The results are presented as bar plots in Figure 4D, where features are ranked in descending order by KernelSHAP attributions. To simplify the amount of information for analysis, the top 20 features from either method are presented in Figure 4E, which makes it easy to see which features were common to both KernelSHAP and LBM (blue and orange bars together), which were identified only by KernelSHAP (blue bars only), and which were identified only by LBM (orange bars only).

Despite the similar increase in ROM predictions for these two patients, the features that contributed to their ROM trajectories differed. In the 15-hour window of interest, the increasing ROM predictions for p_1 were attributed to 177 (45%) of the 398 input features, while the predictions for p_2 were attributed to 188 (47%). Despite the model not having any input of specific diagnoses, the identified features are clinically consistent with each patient’s diagnoses. The top features for the patient with pneumonia and ARDS (p_1) are either related to the respiratory system (EtCO₂, Pulse Oximetry, Respiratory Rate) [57] or associated with infections (Temperature, Heart Rate, Systolic/Diastolic/Mean Arterial Blood Pressure, and 5 blood gas measurements (ABG pH, ABG PO₂, ABG PCO₂, ABG HCO₃, ABG TCO₂) [55]. Two of the top ranked features for the patient with seizures (p_2) – Intracranial Pressure (rank 3) and Cerebral Perfusion Pressure (rank 5) – align with her diagnoses of intracranial hemorrhage. In addition, contributing features for this individual’s increase in predicted ROM included EtCO₂, Pulse Oximetry, and Respiratory Rate (rank 8, 11, and 10 respectively), which align with her diagnosis of ARDS [57]. Other top contributing features for this individual – Heart Rate, Systolic/Diastolic/Mean Arterial Blood Pressure – are general markers of critical illness [38, 42, 50].

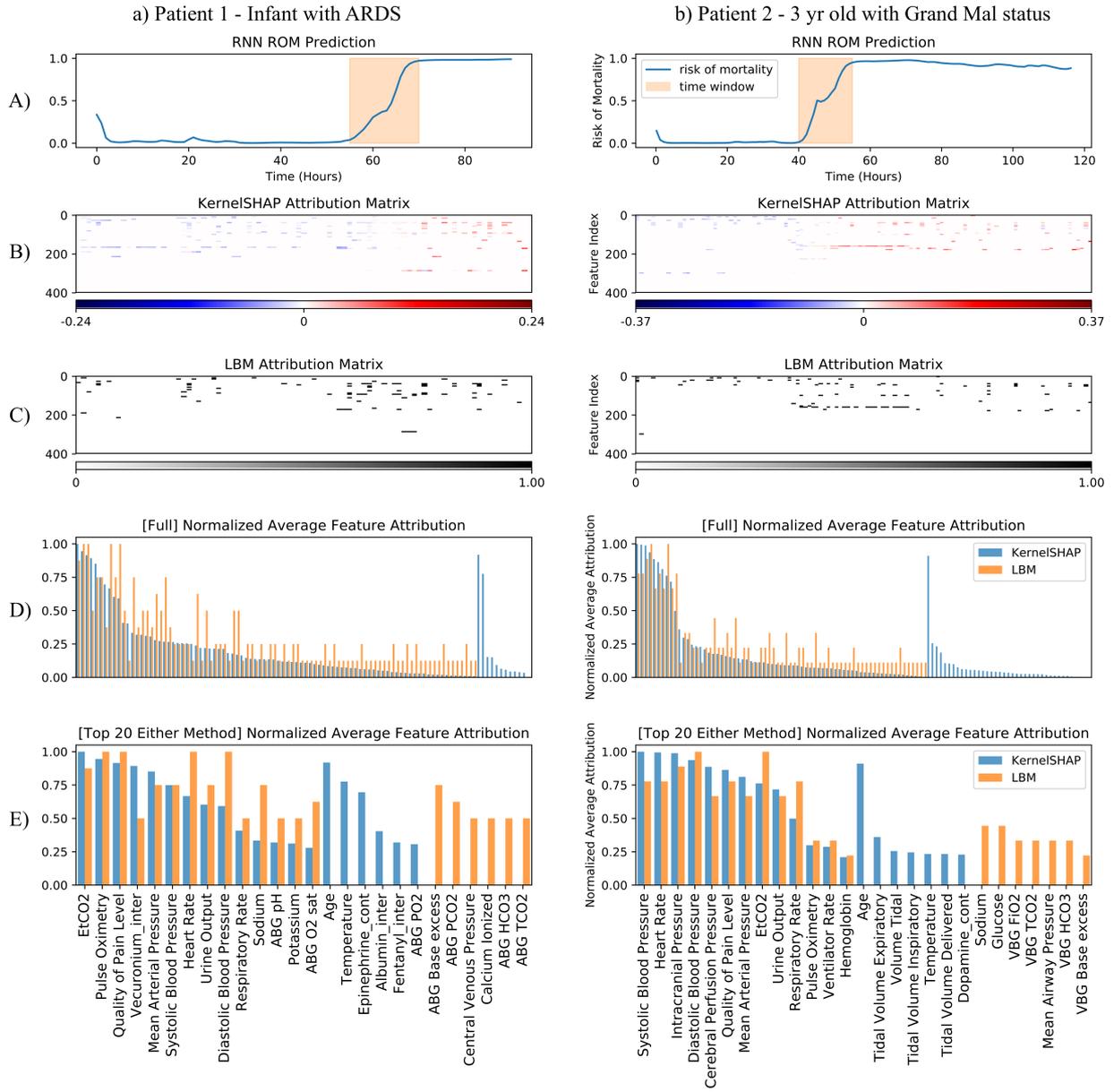


Figure 4: Predictions and explanations for two individual encounters: a) p_1 and b) p_2 . ROM predictions are visualized in A (top panel). KernelSHAP and LBM attribution matrices are shown as heatmaps in panels B and C, respectively. For KernelSHAP, the heatmap values range from negative to positive in probability units. For LBM, the heatmap is binary. Attribution matrices were averaged over time periods highlighted in panel A using Equations 9 & 10 to identify the features that contributed to the increasing ROM predictions in the highlighted time window. These features are visualized in panel D. Finally, panel E visualizes a subset of the features in panel D, presenting only the top 20 features from either method. Note that attribution matrices in B & C are plotted with time on the x-axis (corresponding to ROM plots in A) and features on the y-axis. Also note that there could be more than 20 variables in E as the selected top 20 features overlap between the methods.

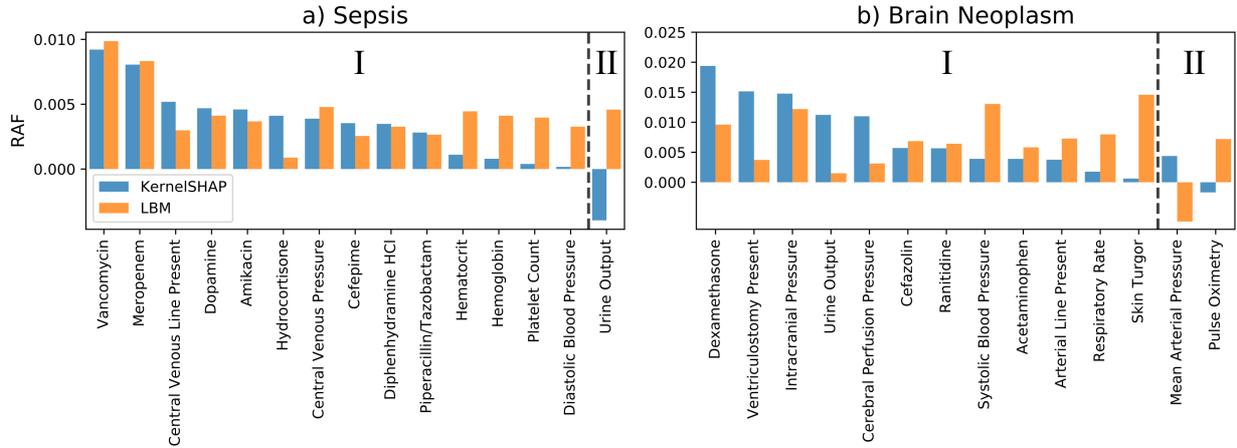


Figure 5: Top 10 Relative attribution features (RAFs) computed across two primary diagnoses: a) sepsis, b) brain neoplasm. Further, the plots are partitioned into two regions: I) when KernelSHAP and LBM align; and II) when KernelSHAP and LBM disagree. Note that there could be more than 10 variables as the top 10 features overlap between the methods.

4.3. Average Feature Attributions Within Cohorts (Relative Attribution Features)

The attribution matrices from KernelSHAP and LBM were aggregated using Equations 11 & 12 to compute relative attribution features (RAF) across two sub-populations of primary disease diagnoses: a) sepsis ($n(S) = 122$), b) brain neoplasm ($n(S) = 112$); recall $n(S)$ is the number of patients in sub-population S . The features with the top 10 RAF from either method are shown in Figure 5. Each bar plot is further partitioned into two regions: I) when both KernelSHAP and LBM align (i.e. both are positive); and II) when KernelSHAP and LBM disagree. Features in Region I are variables which both methods identified as affecting predicted ROM to a greater extent in the specific disease cohort relative to the general critically ill population. Note that features with negative RAF (denoting higher importance in the general critically ill population cohort, S^c , than in the specified disease cohort) from both methods were not included. Features in Region II are those where the methods disagreed. The following paragraphs provide a general description of each disease as well as relevant observations from Figure 5.

Sepsis. Sepsis is a condition in which the body responds to a severe infection by releasing chemicals which can also damage multiple organ systems, cause hemodynamic instability, and result in abnormal blood counts. Sepsis is treated by treating the underlying infection with antibiotics and maintaining blood flow, fluids, and organ function [55]. Not surprisingly, therefore, five of the features in Region I are antibiotics (Vancomycin, Meropenem, Amikacin, Cefepime, Piperacillin/Tazobactam). Another three are measures of blood counts (Hematocrit, Hemoglobin, and Platelet Count), three are related to blood pressure measurements (Central Venous Pressure/Line Present, Diastolic Blood Pressure), one is a drug to increase blood pressure (Dopamine), and another is a drug often used in sepsis to treat hemodynamic instability (Hydrocortisone).

Brain neoplasm. Brain neoplasms cause increased intracranial pressure and altered neurologic function [53]. Treatments include maintaining normal intracranial pressure using interventions such as drugs and shunts. Of the twelve variables in Region I, two are related to measurements of brain pressure and perfusion (Intracranial Pressure, Cerebral Perfusion Pressure), and three are interventions associated with elevated intracranial pressure (Ventriculostomy, Dexamethasone, Ranitidine).

4.4. Feature Contributions in All Critically Ill Children

Using Equations 11 & 13, KernelSHAP and LBM attribution matrices were averaged for *all* critically ill children at all times in the test set to compute a form of model feature importance. This population consisted of 2008 patient encounters. Because the importance from both methods decays rapidly, only the top 50 are investigated. The top 50 features from both methods are displayed in Figure 6. Features are colored by their variable type: vitals in blue, labs in green, interventions in orange, drugs in red, and statics (demographics data) in purple. Variables that were common to both KernelSHAP’s and LBM’s top 50 are indicated with darker text labels while those that are not are labeled in grey.

Because different mortality models often use different feature sets, it is difficult to compare the feature importances computed here for the RNN model with weights from other models such as PRISM3 or PIM2, which often use variables measured specifically for the algorithm and not recorded in the EMR. A notable observation from Figure 6 is that a majority of the top features extracted by both KernelSHAP and LBM were variables measured *internally* to the patient, i.e. vitals, labs, age, and gender. This is enumerated in Table 1 and is consistent with previously reported results [23], which showed little to no degradation in model AUC when external variables (interventions and drugs) were excluded from the input patient representation. In particular, such fundamental measures of the status of critically ill children such as heart rate, blood pressures, blood gases, pulse oximetry (which are routinely monitored in critically ill children), skin turgor, and observations reflecting the level of consciousness (pain perception, comfort, etc) [38, 42, 50] rank so highly. These observations give us confidence that the RNN model reasonably used the input features to generate risk of mortality predictions, and that the LBM and KernelSHAP extracted this information.

Variable Type	KernelSHAP	LBM
Vitals	30 (60%)	32 (64%)
Labs	5 (10%)	12 (23%)
Interventions	4 (8%)	2 (4%)
Drugs	10 (20%)	4 (8%)
Statics (Others)	1 (2%)	0 (0%)

Table 1: Count (and percentage) of each variable type in the top 50 feature importance from KernelSHAP and LBM in Figure 6.

4.5. Compute Time

Figure 7 shows the average computational time of each method across all patients in the validation and test set. The median (IQR) time to compute predictions for a patient is 0.002 (0.001) minutes. To interpret the same predictions, KernelSHAP takes a median (IQR) time of 0.95 (1.90) minutes and LBM takes 4.60 (7.33) minutes. It should be noted that KernelSHAP can be much more computationally expensive for large feature spaces (398 here), however, the *dt-patient-matrix* representation enables practical use of KernelSHAP to execute in real-time.

5. Discussion

Methods for interpreting deep learning models are very diverse in method and purpose. This study focused on determining the contribution of model inputs to model predictions. As a proof of concept, we have described in detail two methods that provide information about which clinical features made the most important contributions to risk of mortality predictions generated by a previously well described recurrent neural network (RNN) model using

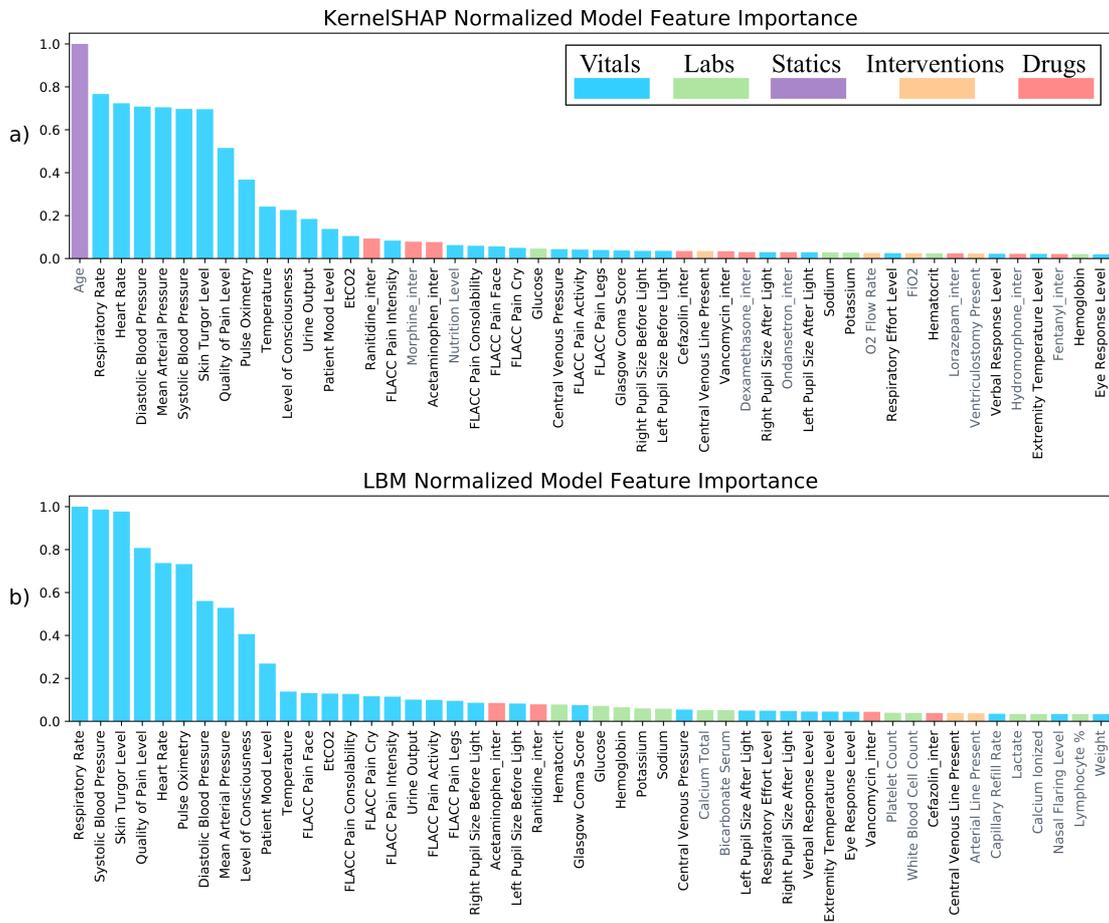


Figure 6: Top 50 features from model “feature importance” computed using a) KernelSHAP and b) LBM. Variables are colored by their variable types, denoted in the legend. Variables common to both KernelSHAP’s top 50 and LBM’s top 50 are indicated with darker text labels.

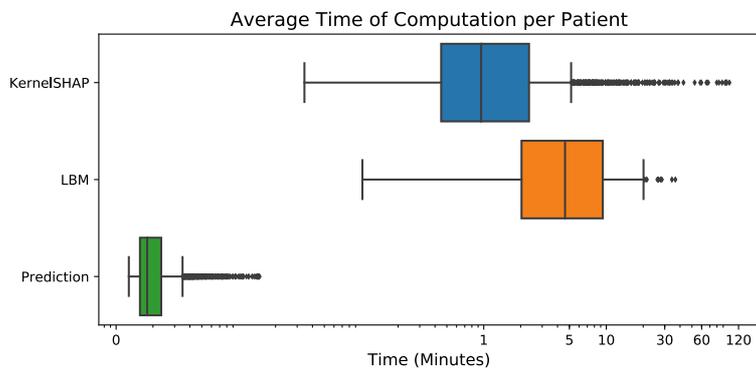


Figure 7: Box-and-whiskers plot showing distribution of time-to-compute per patient for KernelSHAP and LBM. For reference, prediction time-to-compute per patient is also included.

electronic health data of critically ill children. The first, Learned Binary Mask (LBM), is a newly described novel method based on occlusion principles to determine which inputs led to the mortality predictions of the many-to-many RNN ICU mortality model using multi-model EMR data. The second entailed modifying a previously described and familiar methodology, KernelSHAP, to perform on the same RNN model. While methods for interpreting deep learning models exist, most are not compatible with many-to-many RNN models whose inputs are multi-modal EMR data that include time series measurements. This motivated the approaches we took in this study with the LBM and KernelSHAP.

We focused this initial study of the LBM and KernelSHAP on the RNN mortality model. Given this model, its inputs, and the resulting trajectory of risk of mortality (ROM) predictions for an individual patient encounter, the LBM and KernelSHAP each estimate an attribution matrix that reflects each input’s contribution to those predictions. Because the LBM and KernelSHAP have different formulations, they generate different attribution matrices for the same patient encounter. Each matrix answers different questions. Having both perspectives affords a more comprehensive picture of the patient than having either one alone.

The analyses of individual patients in Section 4.2 and Appendix C show that the RNN model used patient-specific features, in addition to general markers for critical illness, for generating ROM predictions, highlighting the individual nature of the RNN predictions and what led to them. Aggregating individual attributions across disease cohorts and the entire population of critically ill children enabled demonstrations of the LBM and KernelSHAP, as well as the RNN model, on a larger scale than analyzing individual patients. The evaluations were all qualitative: features identified by the LBM and KernelSHAP as important contributors to the ROM predictions for specified cohorts (sepsis patients, brain neoplasm patients, general ICU) were consistent with generally described clinical characteristics of their respective disease processes [55, 53]. The results (Section 4.3, Section 4.4) support the notion that the LBM and KernelSHAP identify clinical features most relevant in an individual patient, a specific disease cohort and the entire critically ill population.

Some interesting observations arose from analyzing these feature attributions. Not surprisingly, vital signs dominated the top contributing features for the entire population (Figure 6) because they are general markers of critical illness. While a few drugs and interventions were also in the list of top 50 contributors, their contributions as measured by LBM or KernelSHAP were small compared to those of the vital signs. In contrast, several of the important contributors to individual patient ROM scores included drugs and interventions (Figure 4, Figure C.8). These contrasting observations likely result from the way the attribution matrices were aggregated and averaged. The population-level aggregation (Figure 6) included the entire ICU stay of each patient, while the patient-level aggregation (Figure 4) was confined to a limited time window of the patient’s ICU stay, selected for instability: where the ROM predictions were rapidly increasing. The window of increasing ROM may be regarded as a period of volatility for a patient, and during this period, drugs and interventions were important. The contributions of these drugs were diluted during the averaging process over the entire ICU stay of a patient who was stable for a majority of that stay. The difference in importance of drugs and interventions between Figure 4 and Figure 6, therefore, likely results from the fact that most patients survived their ICU encounter (96%) and spent the majority of their ICU stay in a “stable” (non-volatile) state.

The analyses of attributions in disease-specific cohorts (Figure 5) also had noteworthy findings. In the sepsis group, 13 out of the 15 features with the highest RAF were consistent with clinical expectations [55]; in the brain neoplasm group, it was 5 out of the top 11 [53]. This is significant because patients often have multiple diseases, yet the aggregated results identified features that were consistent with clinical knowledge on these cohorts’ diagnoses. We note, however, that this general alignment potentially results from the inherent bias in EMR data. For example,

antibiotics are more often given to sepsis patients, and the computed RAF may simply reflect that bias. One possible approach to distinguish between data biases and actual impact on mortality predictions would be to compare the RAFs with the frequency that the variables were measured or administered; this is left for future work.

The LBM and KernelSHAP identified features that are not specifically associated with the specified diagnoses of the individuals or the disease cohorts. It is important to note that these apparently non-specific features actually reflect critical illness and are in fact features common to many severity of illness scores [42, 50]. The LBM and KernelSHAP were formulated to determine features that contributed to patients’ risk of ICU mortality, which reflect their severity of critical illness; they were not formulated to identify risk factors for specific diagnoses. Therefore, it is not surprising – in fact, it is reassuring – that they identified general markers of critical illness. It is remarkable that they also identified features specific to a patient’s diagnoses despite having no information about those diagnoses.

Figure 6 shows that of the top 50 contributing features separately identified by the LBM and KernelSHAP, 39 (78%) were shared, illustrating remarkable consilience between the two methods. Consilience between the two methods was also seen in the analyses of individual patients (Section 4.2 and Appendix C) and disease-cohorts (Figure 5). Some concordance between the results from the two methods is not surprising and is, in fact, reassuring because they are conceptually similar in goal: to identify the set of features that contributed to the current prediction.

Nevertheless, the two methods formulate feature contribution differently. LBM identifies contributing features by finding the inputs that must be *zeroed* to drive the mortality predictions to zero (i.e. provide evidence for non-zero mortality predictions). KernelSHAP expresses a prediction as a sum of a background value and marginal contributions from the inputs. The different perspectives and objective functions of the LBM and KernelSHAP gave rise to the differences between their attribution matrices. In particular, they treat zero-valued or near zero-valued inputs differently. LBM does not mask inputs that are exactly zero-valued because masking such inputs adds no benefit but increases the objective function the LBM is minimizing. Masking inputs that are near zero also incurs a high cost in the LBM objective function; therefore, there is a very small chance of LBM selecting inputs that are near zero. In contrast, KernelSHAP has no such restrictions against zero-valued or near zero-valued inputs. The hour-to-hour changes of most drugs and interventions are zero for extended periods. Consequently, these features were not significant contributors in the LBM formulation, but were in the KernelSHAP formulation. Age is also treated differently by the two methods, as was seen in the individual case studies (Figure 4E, Figure C.8E) and in the general ICU population (Figure 6). All entries for age after the initial timestep, are exactly an hour, which is negligible in terms of years, the unit used to express age. In the initial timestep, the entry for age is the deviation of the patient’s age from the population’s mean age, and this varies across different patients. While the LBM can mask age at the initial time step, it very likely ignores the negligible 1-hour age change at all subsequent timesteps. Consequently, the aggregation over many timesteps dilutes any contribution that may come from the initial timestep. Mechanisms related to KernelSHAP alone, such as the ad hoc approximations to make the Shapley equation (Equation 8) computationally tractable, likely gave rise to some of the other differences observed between the two methods’ attribution matrices and the resulting aggregations.

While the results of applying the LBM and KernelSHAP to the RNN mortality model are promising, some questions remain. The LBM and KernelSHAP minimize objective functions that may have multiple local minima, and the actual numerical approximations of these minima depend on hyperparameters. The stability of these solutions, i.e. how the resulting feature contributions change with these solutions, has yet to be investigated. Understanding the sensitivity of the resulting clinical interpretations with respect to perturbations caused by artificial numerical issues is important. Closely related questions are the stability of the LBM or KernelSHAP outputs with respect to the RNN model itself, and with respect to the model inputs. For LBM, even a very slight change in an input’s value could

flip an attribution matrix element from 0 to 1 (or 1 to 0) because of the hard thresholding (Equation 6) to make the attributions binary. The aggregation of these elements over many time points can smooth this behavior.

Exploring generalizability of the LBM is another avenue for future work. This will mean applying the LBM to different clinical datasets, tasks, and models. Finally, the aggregation examples presented here – averaging over a narrow time period of an individual patient, averaging over the entire ICU stay of a patient then aggregating over many encounters – focused on retrospective use and represent only a small sampling of possibilities. Perhaps the most promising use for clinical deployment is demonstrated by time windows in individual patients. This would allow identifying periods of instability when the ROM is rapidly changing and determining what features make the major contribution to these changes. This approach could potentially direct attention to underlying contributions to this instability and direct therapeutic interventions. What and how to present information from the individual attribution matrices computed by the LBM or KernelSHAP will depend on the actual use case. Deploying the RNN model with the LBM and KernelSHAP for real-time clinical use has many considerations involving many different areas of understanding [26, 25], and these will require their own separate investigations.

We emphasize that this study is not about feature selection for model development. Feature selection methods typically are used to prune variables from a large subset of potential input features to improve model performance as measured by metrics such as area under the curve (for classification tasks) or mean absolute errors (for regression tasks). Although the LBM and KernelSHAP could potentially be used for this purpose, it was not the purpose here. Instead, they were used to determine and understand the contributions of predefined inputs to the predictions that have already been made by an existing (already trained) model.

6. Conclusion

This proof-of-concept study presented two methods for providing information about which features made the most important contributions to the risk of mortality predictions of a previously described RNN model using the EMR of critically ill children. The first, Learned Binary Mask (LBM), is a novel occlusion-based method developed here. The second is an existing explanation method, KernelSHAP, modified in this study to make it compatible with the same RNN model. A novel representation of patient data enabled practical application of both methods on the many-to-many RNN model. For any given patient, each method generated an attribution matrix that reflects how each input contributed to the RNN's predictions for that patient. The individual contributions were aggregated across different scales - a specified time window of an individual, or entire encounters of disease subgroups or the whole population of critically ill children in the ICU – to determine the clinical features which were most important to the mortality risk predictions for the specified cohort. Because the two methods have different formulations, their outputs provide complementary perspectives. While initial results show promise – some consistency with established clinical knowledge – several important questions about both the LBM and KernelSHAP remain open for investigation. These include, but are not limited to, more comprehensive clinical evaluations of their results; the stability of their solutions; and their generalization to other datasets, tasks, and models.

Funding

This work was supported by the L. K. Whittier Foundation.

References

- [1] Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., Zheng, X., 2015. TensorFlow: Large-scale machine learning on heterogeneous systems. URL: <https://www.tensorflow.org/>. software available from tensorflow.org.
- [2] Abend, N., Kessler, S., Licht, D., 2016. Evaluation of the comatose child, in: Shaffner, D., Nichols, D. (Eds.), *Roger's Textbook of Pediatric Intensive Care*. Wolters Kluwer Health, Philadelphia, PA. chapter 57, pp. 896–919.
- [3] Aczon, M., Ledbetter, D., Laskana, E., Ho, L., Wetzel, R., 2021. Continuous prediction of mortality in the picu: a recurrent neural network model in a single center dataset. *Pediatric Critical Care Medicine*, 2021 (in press) .
- [4] Cerna, A.E.U., Pattichis, M., VanMaanen, D.P., Jing, L., Patel, A.A., Stough, J.V., Haggerty, C.M., Fornwalt, B.K., 2019. Interpretable neural networks for predicting mortality risk using multi-modal electronic health records. *arXiv preprint arXiv:1901.08125* .
- [5] Che, Z., Purushotham, S., Khemani, R., Liu, Y., 2016. Interpretable deep models for icu outcome prediction, in: *AMIA Annual Symposium Proceedings*, American Medical Informatics Association. p. 371.
- [6] Cheng, J.Z., Ni, D., Chou, Y.H., Qin, J., Tiu, C.M., Chang, Y.C., Huang, C.S., Shen, D., Chen, C.M., 2016. Computer-aided diagnosis with deep learning architecture: applications to breast lesions in us images and pulmonary nodules in ct scans. *Scientific reports* 6, 24454.
- [7] Choi, E., Bahadori, M.T., Schuetz, A., Stewart, W.F., Sun, J., 2016a. Doctor ai: Predicting clinical events via recurrent neural networks, in: *Machine Learning for Healthcare Conference*, pp. 301–318.
- [8] Choi, E., Bahadori, M.T., Sun, J., Kulas, J., Schuetz, A., Stewart, W., 2016b. Retain: An interpretable predictive model for healthcare using reverse time attention mechanism, in: *Advances in Neural Information Processing Systems*, pp. 3504–3512.
- [9] Chollet, F., et al., 2015. Keras. <https://keras.io>.
- [10] Cicero, M., Bilbily, A., Colak, E., Dowdell, T., Gray, B., Perampaladas, K., Barfett, J., 2017. Training and validating a deep convolutional neural network for computer-aided detection and classification of abnormalities on frontal chest radiographs. *Investigative radiology* 52, 281–287.
- [11] Donabedian, A., 1966. Evaluating the quality of medical care. *The Milbank memorial fund quarterly* 44, 166–206.
- [12] Donabedian, A., 1988. The quality of care: how can it be assessed? *Jama* 260, 1743–1748.
- [13] Doshi-Velez, F., Kim, B., 2017. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608* .
- [14] Esteva, A., Kuprel, B., Novoa, R.A., Ko, J., Swetter, S.M., Blau, H.M., Thrun, S., 2017. Dermatologist-level classification of skin cancer with deep neural networks. *Nature* 542, 115.
- [15] Esteva, A., Robicquet, A., Ramsundar, B., Kuleshov, V., DePristo, M., Chou, K., Cui, C., Corrado, G., Thrun, S., Dean, J., 2019. A guide to deep learning in healthcare. *Nature medicine* 25, 24–29.
- [16] Ferreira, A., 2019. <https://andreconf.github.io>. URL: <https://andreconf.github.io/2019/07/31/InterpretingRecurrentNeuralNetworksOnMultivariateTimeSeriesData/>
- [17] Fong, R.C., Vedaldi, A., 2017. Interpretable explanations of black boxes by meaningful perturbation, in: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3429–3437.
- [18] Glorot, X., Bengio, Y., 2010. Understanding the difficulty of training deep feedforward neural networks, in: *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 249–256.
- [19] Goldstein, B.A., Navar, A.M., Pencina, M.J., Ioannidis, J., 2017. Opportunities and challenges in developing risk prediction models with electronic health records data: a systematic review. *Journal of the American Medical Informatics Association* 24, 198–208.
- [20] Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., Pedreschi, D., 2018. A survey of methods for explaining black box models. *ACM Comput. Surv.* 51, 93:1–93:42. URL: <http://doi.acm.org/10.1145/3236009>, doi:10.1145/3236009.
- [21] Gulshan, V., Peng, L., Coram, M., Stumpe, M.C., Wu, D., Narayanaswamy, A., Venugopalan, S., Widner, K., Madams, T., Cuadros, J., et al., 2016. Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *Jama* 316, 2402–2410.
- [22] Henry, J., Pylypchuk, Y., Searcy, T., Patel, V., 2016. Adoption of electronic health record systems among us non-federal acute care hospitals: 2008-2015. Retrieved from <http://dashboard.healthit.gov/evaluations/data-briefs/non-federal-acute-care-hospital-ehr-adoption-2008-2015.php> .
- [23] Ho, L.V., Ledbetter, D., Aczon, M., Wetzel, R., 2017. The dependence of machine learning on electronic medical record quality, in: *AMIA Annual Symposium Proceedings*, American Medical Informatics Association. p. 883.
- [24] Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural computation* 9, 1735–1780.
- [25] Keim-Malpass, J., Kitzmiller, R.R., Skeeles-Worley, A., Lindberg, C., Clark, M.T., Tai, R., Calland, J.F., Sullivan, K., Moorman, J.R., Anderson, R.A., 2018. Advancing continuous predictive analytics monitoring: Moving from implementation to clinical action in a learning health system. *Critical Care Nursing Clinics* 30, 273–287.

- [26] Kitzmiller, R., Vaughan, A., Skeeles-Worley, A., Keim-Malpass, J., Yap, T., Lindberg, C., Kennerly, S., Mitchell, C., Tai, R., Sullivan, B., et al., 2019. Diffusing an innovation: Clinician perceptions of continuous predictive analytics monitoring in intensive care. *Applied clinical informatics* 10, 295.
- [27] Kooi, T., Litjens, G., Van Ginneken, B., Gubern-Mérida, A., Sánchez, C.I., Mann, R., den Heeten, A., Karssemeijer, N., 2017. Large scale deep learning for computer aided detection of mammographic lesions. *Medical image analysis* 35, 303–312.
- [28] Laksana, E., Aczon, M., Ho, L., Carlin, C., Ledbetter, D., Wetzel, R., 2020. The impact of extraneous features on the performance of recurrent neural network models in clinical tasks. *Journal of Biomedical Informatics* 102, 103351.
- [29] LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. *nature* 521, 436–444.
- [30] Leisman, D.E., Harhay, M.O., Lederer, D.J., Abramson, M., Adjei, A.A., Bakker, J., Ballas, Z.K., Barreiro, E., Bell, S.C., Bellomo, R., et al., 2020. Development and reporting of prediction models: guidance for authors from editors of respiratory, sleep, and critical care journals. *Critical care medicine* 48, 623.
- [31] Leteurtre, S., Duhamel, A., Deken, V., Lacroix, J., Leclerc, F., de Réanimation et Urgences Pédiatriques (GFRUP, G.F., et al., 2015. Daily estimation of the severity of organ dysfunctions in critically ill children by using the pelod-2 score. *Critical care* 19, 324.
- [32] Lipton, Z.C., Kale, D.C., Elkan, C., Wetzel, R., 2015. Learning to diagnose with lstm recurrent neural networks. *arXiv preprint arXiv:1511.03677* .
- [33] Liu, Y., Gadepalli, K., Norouzi, M., Dahl, G.E., Kohlberger, T., Boyko, A., Venugopalan, S., Timofeev, A., Nelson, P.Q., Corrado, G.S., et al., 2017. Detecting cancer metastases on gigapixel pathology images. *arXiv preprint arXiv:1703.02442* .
- [34] Lundberg, S.M., Erion, G., Chen, H., DeGrave, A., Prutkin, J.M., Nair, B., Katz, R., Himmelfarb, J., Bansal, N., Lee, S.I., 2019. Explainable ai for trees: From local explanations to global understanding. *arXiv preprint arXiv:1905.04610* .
- [35] Lundberg, S.M., Erion, G.G., Lee, S.I., 2018. Consistent individualized feature attribution for tree ensembles. *arXiv preprint arXiv:1802.03888* .
- [36] Lundberg, S.M., Lee, S.I., 2017. A unified approach to interpreting model predictions, in: *Advances in Neural Information Processing Systems*, pp. 4765–4774.
- [37] Molnar, C., 2019. *Interpretable Machine Learning*. <https://christophm.github.io/interpretable-ml-book/>.
- [38] Nichols, D.G., Shaffner, D.H., Argent, A.C., Arnold, J.H., Biagas, K.V., 2016. *Rogers' textbook of pediatric intensive care*. Wolters Kluwer Philadelphia.
- [39] Olah, C., Mordvintsev, A., Schubert, L., 2017. Feature visualization. *Distill* doi:10.23915/distill.00007. <https://distill.pub/2017/feature-visualization>.
- [40] Pollack, M.M., 2016. Severity of illness confusion. *Pediatric Critical Care Medicine* 17, 583.
- [41] Pollack, M.M., Holubkov, R., Funai, T., Dean, J.M., Berger, J.T., Wessel, D.L., Meert, K., Berg, R.A., Newth, C.J., Harrison, R.E., et al., 2016. The pediatric risk of mortality score: update 2015. *Pediatric Critical Care Medicine* 17, 2–9.
- [42] Pollack, M.M., Patel, K.M., Ruttimann, U.E., 1996. Prism iii: an updated pediatric risk of mortality score. *Critical care medicine* 24, 743–752.
- [43] Poursabzi-Sangdeh, F., Goldstein, D.G., Hofman, J.M., Vaughan, J.W., Wallach, H., 2018. Manipulating and measuring model interpretability. *arXiv preprint arXiv:1802.07810* .
- [44] Rajkomar, A., Oren, E., Chen, K., Dai, A.M., Hajaj, N., Hardt, M., Liu, P.J., Liu, X., Marcus, J., Sun, M., et al., 2018. Scalable and accurate deep learning with electronic health records. *NPJ Digital Medicine* 1, 18.
- [45] Ribeiro, M.T., Singh, S., Guestrin, C., 2016. " why should i trust you?" explaining the predictions of any classifier, in: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135–1144.
- [46] Samek, W., Wiegand, T., Müller, K.R., 2017. Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. *arXiv preprint arXiv:1708.08296* .
- [47] Schulman, C.S., Staul, L., 2010. Standards for frequency of measurement and documentation of vital signs and physical assessments. *Critical Care Nurse* 30, 74–76.
- [48] Scott M. Lundberg, S.I.L., 2019. *Shap*. <https://github.com/slundberg/shap>.
- [49] Shapley, L.S., 1953. A value for n-person games. *Contributions to the Theory of Games* 2, 307–317.
- [50] Slater, A., Shann, F., Pearson, G., Group, P.S., et al., 2003. Pim2: a revised version of the paediatric index of mortality. *Intensive care medicine* 29, 278–285.
- [51] Štrumbelj, E., Kononenko, I., 2014. Explaining prediction models and individual predictions with feature contributions. *Knowledge and information systems* 41, 647–665.
- [52] Suresh, H., Hunt, N., Johnson, A., Celi, L.A., Szolovits, P., Ghassemi, M., 2017. Clinical intervention prediction and understanding using deep networks. *arXiv preprint arXiv:1705.08498* .
- [53] Tasker, R., Aboy, M., Graham, A., Goldstein, B., 2016. Neurologic monitoring, in: Shaffner, D., Nichols, D. (Eds.), *Rogers' Textbook of*

- Pediatric Intensive Care. Wolters Kluwer Health, Philadelphia, PA. chapter 58, pp. 907–919.
- [54] Tasker, R.C., Randolph, A.G., 2016. Severity-of-illness scoring in pediatric critical care: Quo vadis? *Pediatric Critical Care Medicine* 17, 83–85.
- [55] Thomas, N., Tamburro, R., Rajasekaran, S., Fitzgerald, J., Weiss, S., Hall, M., 2016. Bacterial sepsis, in: Shaffner, D., Nichols, D. (Eds.), *Roger's Textbook of Pediatric Intensive Care*. Wolters Kluwer Health, Philadelphia, PA. chapter 87, pp. 1377–1396.
- [56] Tieleman, T., Hinton, G., 2012. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural networks for machine learning 4.
- [57] Ventre, K., Arnold, J., 2016. Acute lung injury and acute respiratory distress syndrome, in: Shaffner, D., Nichols, D. (Eds.), *Roger's Textbook of Pediatric Intensive Care*. Wolters Kluwer Health, Philadelphia, PA. chapter 49, pp. 766–793.
- [58] Winter, M.C., Day, T.E., Ledbetter, D.R., Aczon, M.D., Newth, C.J.L., Wetzell, R.C., Ross, P.A., 2020. Machine learning to predict cardiac death within 1 hour after terminal extubation. *Pediatric Critical Care Medicine Online First*. doi:10.1097/PCC.0000000000002612.
- [59] Zhang, J., Kowsari, K., Harrison, J.H., Lobo, J.M., Barnes, L.E., 2018. Patient2vec: A personalized interpretable deep representation of the longitudinal electronic health record. *IEEE Access* 6, 65333–65346.

Appendix A. Datasets Summary & Feature List

Table A.2: Summary statistics and demographics, partitioned into training (60%), validation (20%), and test (20%) sets. Presented are the total number of encounters per category with percentages rounded to the nearest integer in parenthesis.

		Train N (%)	Valid N (%)	Test N (%)
Mortality				
	Survived	5645 (96)	1885 (96)	1925 (96)
	Died	240 (4)	77 (4)	83 (4)
Age Group				
	[0, 1)	976 (17)	371 (19)	309 (15)
	[1, 5)	1449 (25)	481 (25)	496 (25)
	[5, 10)	1052 (18)	331 (17)	362 (18)
	[10, 18)	1941 (33)	644 (33)	683 (34)
	18+	467 (8)	135 (7)	158 (8)
Primary Diagnosis Category				
	Respiratory	1682 (29)	579 (30)	558 (28)
	Neurologic	884 (15)	232 (12)	289 (14)
	Oncologic	670 (11)	214 (11)	233 (12)
	Injury/Poisoning/Adverse Effects	583 (10)	228 (12)	185 (9)
	Orthopedic	448 (8)	158 (8)	153 (8)
	Infectious	417 (7)	155 (8)	153 (8)
	Other	339 (6)	125 (6)	118 (6)
	Gastrointestinal	288 (5)	91 (5)	115 (6)
	Genetic	243 (4)	80 (4)	73 (4)
	Cardiovascular	189 (3)	56 (3)	81 (4)
	Renal/Genitourinary	142 (2)	44 (2)	50 (2)

Table A.3: List of 398 features used as inputs to the RNN model

Pulse Oximetry	Heart Rate	Respiratory Rate	Weight	Systolic Blood Pressure
Diastolic Blood Pressure	Mean Arterial Pressure	Motor Response Level	Verbal Response Level	Eye Response Level
Glasgow Coma Score	Temperature	Right Pupillary Response Level	Level of Consciousness	Left Pupillary Response Level
Extremity Temperature Level	Patient Mood Level	Respiratory Effort Level	Capillary Refill Rate	Skin Turgor turgor
Right Pupil Size Before Light	Left Pupil Size Before Light	Nasal Flaring Level	Right Pupil Size After Light	Left Pupil Size After Light
Quality of Pain Level	Height	FLACC Pain Face	FLACC Pain Legs	FLACC Pain Activity
FLACC Pain Cry	FLACC Pain Consolability	FLACC Pain Intensity	Nutrition Level	Lip Moisture Level
Capillary Refill Delayed	Age	Sex F	Foley Catheter Volume	Sodium
Potassium	Glucose	Hematocrit	Hemoglobin	Creatinine
Bicarbonate Serum	Central Venous Pressure	Head Circumference	PaO2 to FiO2	Chloride
Calcium Total	BUN	Platelet Count	White Blood Cell Count	RBC Blood
MCH	MCV	MCHC	RDW	O2 Flow Rate
Lymphocyte Percent	Neutrophils Percent	Monocytes Percent	Basophils Percent	Oxygenation Index
Eosinophils Percent	Calcium Ionized	Lactate	Albumin Level	ALT
AST	Bilirubin Total	Alkaline phosphatase	Protein Total	PTT
INR	Abdominal Girth	PT	CBG PCO2	CBG pH
CBG PO2	CBG HCO3	CBG Base excess	CBG TCO2	CBG O2 sat
Phosphorus level	Magnesium Level	Bands Percent	ABG PO2	ABG pH
ABG PCO2	ABG O2 sat	ABG HCO3	ABG Base excess	ABG TCO2
VBG HCO3	VBG Base excess	VBG TCO2	VBG PO2	VBG pH
VBG PCO2	VBG O2 sat	Culture Blood	C-Reactive Protein	Fibrinogen
CBG FiO2	VBG FiO2	ABG FiO2	Culture Urine	Influenza Lab
Schistocytes	Metamyelocytes Percent	Culture Respiratory	Myelocytes Percent	Triglycerides
Lipase	Culture CSF	MVBG Base Excess	MVBG HCO3	MVBG PCO2
MVBG PO2	MVBG TCO2	MVBG pH	MVBG O2 Sat	MVBG FiO2
Oxygen Mode Level	FiO2	Central Venous Line Site	EiCO2	PEEP
Peak Inspiratory Pressure	Ventilator Rate	Inspiratory Time	Mean Airway Pressure	Chest X Ray
Arterial Line Site	Pressure Support	Tidal Volume Expiratory	Tidal Volume Inspiratory	Tidal Volume Delivered
Volume Tidal	CT Brain	EPAP	MRI Brain	NIV Set Rate
Ventriculostomy Site	Chest Tube Site	Abdominal X Ray	Hemofiltration Therapy Mode	ECMO Hours
Acetaminophen inter	Ranitidine inter	Gastrostomy Tube Volume	Morphine inter	Lorazepam inter
Ondansetron inter	Vancomycin inter	Fentanyl inter	Furosemide inter	Intracranial Pressure
Cefazolin inter	Diphenhydramine HCl inter	Pantoprazole inter	Fentanyl cont	Dexamethasone inter
Midazolam HCl inter	IPAP	Potassium Chloride inter	Ceftriaxone inter	Budesonide inter
Piperacillin/Tazobactam inter	Dopamine cont	Dexmedetomidine cont	Hydromorphone inter	Vecuronium inter
CSF RBC	CSF WBC	Methylprednisolone inter	Levetiracetam inter	Midazolam HCl cont
Bilirubin Conjugated	Bilirubin Unconjugated	D-dimer	Macrocystes	Ibuprofen inter
Diazepam inter	Alteplase inter	Ketorolac inter	Amylase	Rocuronium inter
Culture Fungus Blood	Ceftazidime inter	Spherocytes	Meropenem inter	Sodium Chloride inter
Famotidine inter	Cerebral Perfusion Pressure	Albuterol inter	Sodium Bicarbonate inter	Calcium Chloride inter
Albumin inter	Magnesium Sulfate inter	CSF Lymphs Percent	Oxycodone inter	Clindamycin inter
Metronidazole inter	Culture Wound	Phenobarbital inter	Fluconazole inter	Chlorothiazide inter
Lansoprazole inter	Glycopyrrolate inter	CSF Glucose	CSF Protein	Azithromycin inter
Potassium Phosphate inter	Atropine inter	Propofol inter	TSH	Reticulocyte Count
Trimethoprim/Sulfamethoxazole inter	CSF Segs Percent	Racemic Epi inter	Lactate Dehydrogenase Blood	Hydrocortisone inter
Acyclovir inter	Nifedipine inter	T4 Free	Baclofen inter	Acetaminophen/Hydrocodone inter
Cefotaxime inter	Methadone inter	Ampicillin/Sulbactam inter	Ferritin Level	Acetaminophen/Codeine inter
GGT	Tacrolimus inter	Ketamine inter	Nystatin inter	Gabapentin inter
Micafungin inter	ESR	Clonidine HCl inter	B-type Natriuretic Peptide	Ferrous Sulfate inter
Tobramycin inter	Prednisone inter	Enalapril inter	Amikacin inter	Osetamivir inter
Desmopressin inter	Insulin inter	Vitamin K inter	Naloxone HCL cont	Lactobacillus inter
Sodium Phosphate inter	Calcium Gluconate inter	Propofol cont	Ampicillin inter	Fluticasone inter
Olanzapine inter	Spironolactone inter	Aspirin inter	Isradipine inter	Acetazolamide inter
Metoclopramide inter	Amlodipine inter	Montelukast Sodium inter	Amphotericin B Lipid Complex inter	Immune Globulin inter
Heparin inter	Cefepime inter	Levocarnitine inter	Gentamicin inter	Lidocaine inter
Topiramate inter	Filgrastim inter	Labetalol inter	Ursodiol inter	Fosphenytoin inter
Voriconazole inter	CT Chest	Mycophenolate Mofetil inter	Valproic Acid inter	Anti-Xa Heparin
Clonazepam inter	Epinephrine cont	Levalbuterol inter	Sucralfate inter	Aminophylline cont
Oxcarbazepine inter	Ciprofloxacin HCL inter	Ipratropium Bromide inter	Furosemide cont	Levothyroxine Sodium inter
Hydromorphone cont	Insulin cont	Vasopressin cont	Heparin cont	Epoetin inter
Nitric Oxide	Blasts Percent	Amoxicillin inter	Epinephrine inter	Cephalexin inter
Lactic Acid Blood	Aminophylline inter	Sildenafil inter	Enoxaparin inter	Chloral Hydrate inter
Risperidone inter	Haptoglobin	Cefoxitin inter	Valganciclovir inter	Ganciclovir Sodium inter
Basiliximab inter	Amoxicillin/clavulanic acid inter	HFOV Amplitude	HFOV Frequency	CT Abdomen Pelvis
Oxacillin inter	Prednisolone inter	Lisinopril inter	Complement C3 Serum	Complement C4 Serum
Carbamazepine inter	Linezolid inter	Rifampin inter	Pentobarbital inter	Propranolol HCl inter
Milrinone cont	Azathioprine inter	Octreotide Acetate cont	Nitroprusside cont	Cisatracurium cont
CSF Bands Percent	Dornase Alfa inter	Bumetanide inter	Terbutaline cont	Allopurinol inter
Phenytoin inter	Digoxin inter	Cyclophosphamide inter	Calcium Chloride cont	Naloxone HCL inter
Ketamine cont	Levofloxacin inter	Isoniazid inter	Cisatracurium inter	Norepinephrine cont
Penicillin G Sodium inter	Factor VII inter	Erythromycin inter	Dobutamine cont	Varidol inter
Labetalol cont	Sodium Bicarbonate cont	Nitrofurantoin inter	Phenylephrine HCl cont	Vitamin E inter
Haloperidol inter	Esmolol Hydrochloride cont	Aminocaproic Acid inter	Calcium Gluconate cont	Calcium Glubionate inter
Warfarin Sodium inter	Amphotericin B inter	Metolazone inter	Pentobarbital cont	Doxycycline Hyclate inter
Atenolol inter	Cyclosporine inter	Doxorubicin cont	Morphine inter	Lidocaine cont
Aminocaproic Acid cont	Isoproterenol cont	Amiodarone inter	Naproxen inter	Nitroglycerine cont
Vecuronium cont	Etomidate inter	Captopril inter	Nitroglycerine cont	Alprostadil cont
Bumetanide cont	Nesiritide cont	Protamine inter	Sildenafil cont	Procainamide cont
Flecainide Acetate inter	Acetaminophen/Oxycodone inter	Itraconazole inter	Tacrolimus cont	Infliximab inter
Cefuroxime inter	Alteplase cont	Cromolyn Sodium inter		

Appendix B. Learned Binary Mask Algorithmic Description

Algorithm 1: Generating a Learned Binary Mask for an Individual Patient

Step 1: Optimize for intermediate mask $\mathbf{m}_{1:T}$ leveraging back-propagation

Inputs: $\mathbf{x}_{1:T} \in \mathbb{R}^{N \times T}$: dt-patient-matrix defined in Section 3.1.2

$f(\Theta)$: RNN with trained weights Θ described in Section 3.2

Output: $\mathbf{m}_{1:T} \in \mathbb{R}^{N \times T}$: Intermediate mask, solution to Equation 4

1. Generate patient's ROM prediction: $\mathbf{y}_{1:T} = f(\Theta, \mathbf{x}_{1:T})$
2. Modify the RNN by inserting a custom layer, MaskedPerturbLayer, immediately after the input layer
 - with weight $\mathbf{m}_{1:T} = \sigma(A \times \mathbf{z}_{1:T}) \in \mathbb{R}^{N \times T}$, with $\mathbf{z}_{1:T}$ initialized to all ones, $\sigma(x) = \frac{1}{1+e^{-x}}$, and $A = 5$
 - input is $\mathbf{x}_{1:T}$
 - output is $\mathbf{m}_{1:T} \odot \mathbf{x}_{1:T}$
3. Define loss function for optimization of $\mathbf{m}_{1:T}$ as defined in Equation 4
 - with $\lambda_1 = 0.005$ and $\lambda_2 = 0.5$
 - "Freeze" training of weights of all other layers, Θ , such that only weight updated is $\mathbf{z}_{1:T}$
4. Minimize loss function using optimizer RMSProp with initial learning rate of 0.1
 - reduce current learning rate by a factor of 10 if loss has not been improved over the last 5 iterations
 - terminate training if max iteration = 5000, or if learning rate has been reduced 3 times
5. Get intermediate mask from weight of MaskedPerturbLayer: $\mathbf{m}_{1:T} = \sigma(A \times \mathbf{z}_{1:T})$

Step 2 Brute-force grid search to find threshold mask $\boldsymbol{\eta}_{1:T}$

Inputs: $\mathbf{x}_{1:T} \in \mathbb{R}^{N \times T}$: dt-patient-matrix defined in Section 3.1.2

$f(\Theta)$: RNN with trained weights Θ described in Section 3.2

$\mathbf{m}_{1:T} \in \mathbb{R}^{N \times T}$: Intermediate mask, solution to **Step 1**

Outputs: $\boldsymbol{\eta}_{1:T} \in \mathbb{R}^{N \times T}$: Threshold mask

$\mathbf{M}_{1:T} \in \mathbb{R}^{N \times T}$: Binary mask, solution to Equation 6

Initialize: $\boldsymbol{\eta}_{1:T} = \{0.5\}^{N \times T}$; $L_{\min} = \infty$; $\lambda_1^2 = 0.0001$; $s_{\min} = 0.05$; max_iter = 2

for $i \leftarrow 0$ **to** max_iter **do**

for $t = T$ **to** 0 *of each time-step* **do**

 Let u be an array containing all unique values in \mathbf{m}_t ; // intermediate mask at t^{th} timestep

for $u \in U$ **do**

$\mathbf{q}_{1:T} = \boldsymbol{\eta}_{1:T}$; // q is used as a temporary threshold

$\mathbf{q}_t = u$; // t^{th} timestep to a single threshold value

$l = f(\Theta; \mathbf{x}_{1:T} \odot (\mathbf{m}_{1:T} > \mathbf{q}_{1:T})) + \lambda_1^2 \|1 - (\mathbf{m}_{1:T} > \mathbf{q}_{1:T})\|$

if $l < L$ **then**

$L = l$

$\boldsymbol{\eta}_{1:T} = \mathbf{q}_{1:T}$

end

end

end

$\mathbf{s}_{1:T} = f(\Theta; \mathbf{x}_{1:T} \odot (\mathbf{m}_{1:T} > \boldsymbol{\eta}_{1:T}))$

if $\forall t \in [0, T], s_t < s_{\min}$ **then**

 | terminate

end

end

Appendix C. Interpreting Individual Predictions: Additional Cases

KernelSHAP and LBM were used to interpret the RNN ROM predictions for two additional patients, and results are presented in Figure C.8. In contrast to the two non-surviving patients analyzed in Section 4.2, the two patients here survived their ICU stay and do not have any illnesses related to respiratory or cardiovascular systems. The first patient, p_3 , was a 16 year old male with diagnoses of spine curvature disorder, cerebral palsy, and mental retardation. The second patient, p_4 , was a 6 year old female with diagnoses of acute hepatic failure and chronic pancreatitis. Similar to results presented in Section 4.2, the two patients' RNN ROM predictions were analyzed by averaging their LBM and KernelSHAP attribution matrices corresponding to periods of increasing ROM. These periods are highlighted by the beige window shown in Figure C.8A. As in Section 4.2, Equation 9 was applied to obtain $\bar{\mathbf{a}}^{p_3} = \sum_{t=150}^{165} \mathbf{a}^{p_3}(t)/(165 - 150)$ and $\bar{\mathbf{a}}^{p_4} = \sum_{t=90}^{125} \mathbf{a}^{p_4}(t)/(125 - 90)$ for the two patients, and these were subsequently normalized using Equation 10.

The top 20 features from both methods during the highlighted time windows of interest are presented in Figure C.8E. For p_3 , the top features from both methods were vitals and lab values. It is worth noting that no importance was placed on respiratory-related measurements (e.g. EtCO₂) or infection-related measurements (e.g. Temperature) which constitute a majority of the illnesses and treatments in the pediatric ICU. Interestingly, KernelSHAP and LBM highlighted that Lactate contributed to the increase in ROM predictions of patient p_3 . Lactate measurements before and during the time window of interest increased (average Lactate before and during: 9.2 to 24.5 mg/dL). What is interesting about these two patients is that these two patients survive and that there is less specificity to their diagnoses than the two patients presented in Section 4.2 who did not survive their ICU stay. While the features highlighted here are not specific to their diagnoses, many of the important features that are highlighted are those of general critical illness which are common to many severity of illness scores [42, 50]. Although abnormal levels of Lactate do not directly relate to p_3 's diagnoses, it is reassuring that Lactate was highlighted as being important in both LBM and KernelSHAP, as such levels clearly indicate patient decomposition.

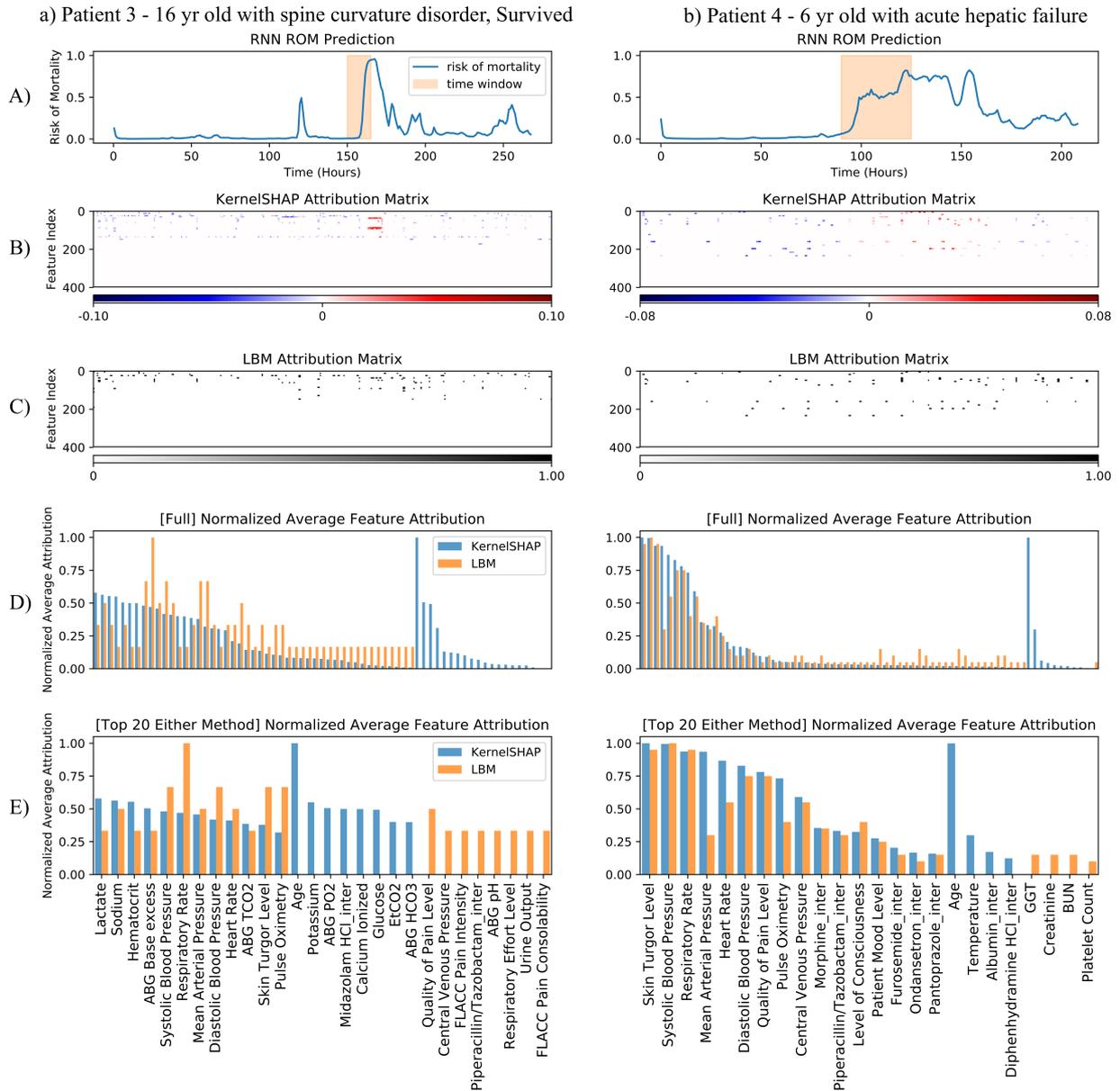


Figure C.8: Predictions and explanations for two additional individual encounters: a) p_3 and b) p_4 . ROM predictions are visualized in A (top panel). KernelSHAP and LBM attribution matrices are shown as heatmaps in panels B and C, respectively. For KernelSHAP, the heatmap values range from negative to positive in probability units. For LBM, the heatmap is binary. Attribution matrices were averaged over time periods highlighted in panel A using Equations 9 & 10 to identify the features that contributed to the increasing ROM predictions in the highlighted time window. These features are visualized in panel D. Finally, panel E visualizes a subset of the features in panel D, presenting only the top 20 features from either method. Note that attribution matrices in B & C are plotted with time on the x-axis (corresponding to ROM plots in A) and features on the y-axis. Also note that there could be more than 20 variables in E as the selected top 20 features overlap between the methods.