

# Enhanced Critical Congenital Cardiac Disease Screening by Combining Interpretable Machine Learning Algorithms

Zhengfeng Lai<sup>1</sup>, Pranjali Vadlaputi<sup>2</sup>, Daniel J. Tancredi<sup>2</sup>, Meena Garg<sup>3</sup>, Robert I. Koppel<sup>4</sup>, Mera Goodman<sup>4</sup>, Whitnee Hogan<sup>5</sup>, Nicole Cresalia<sup>6</sup>, Stephan Juergensen<sup>6</sup>, Erlinda Manalo<sup>7</sup>, Satyan Lakshminrusimha<sup>2</sup>, Chen-Nee Chuah<sup>1</sup>, *Fellow, IEEE*, and Heather Siefkes<sup>2</sup>

**Abstract**—Critical Congenital Heart Disease (CCHD) screening that only uses oxygen saturation (SpO<sub>2</sub>), measured by pulse oximetry, fails to detect an estimated 900 US newborns annually. The addition of other pulse oximetry features such as perfusion index (PIx), heart rate, pulse delay and photoplethysmography characteristics may improve detection of CCHD, especially those with systemic blood flow obstruction such as Coarctation of the Aorta (CoA). To comprehensively study the most relevant features associated with CCHD, we investigated interpretable machine learning (ML) algorithms by using Recursive Feature Elimination (RFE) to identify an optimal subset of features. We then incorporated the trained ML models into the current SpO<sub>2</sub>-alone screening algorithm. Our proposed enhanced CCHD screening system, which adds the ML model, improved sensitivity by approximately 10 percentage points compared to the current standard SpO<sub>2</sub>-alone method with minimal to no impact on specificity.

**Clinical relevance**— This establishes proof of concept for a ML algorithm that combines pulse oximetry features to improve detection of CCHD with little impact on false positive rate.

## I. INTRODUCTION

Congenital heart disease (CHD) is the most common birth defect, affecting nearly 0.8% of all newborn infants [1]. Critical congenital heart disease (CCHD) is a subset of CHD accounting for almost 20% of these infants [1], representing the most severe forms of CHD. CCHD lesions require surgical or catheter-based intervention soon after birth, often including pre-procedural hospitalization and medical management. Late or missed detection of CCHD can lead to significant, preventable morbidity, as well as death [2]–[4]. Prior to mandatory oxygen-saturation (SpO<sub>2</sub>) based CCHD screening, 25% of newborns with CCHD were diagnosed

after going home from the hospital [4], [5]. Although, SpO<sub>2</sub> screening has helped with earlier diagnosis and reduced CCHD mortality, it still misses approximately 900 newborns with CCHD annually in the United States [6], [7]. The majority of the missed types of CCHD defects are those with obstructed systemic blood flow that do not commonly cause low SpO<sub>2</sub>, or hypoxemia [6].

To address this problem, we created an automated real-time data collection system [8] to collect additional pulse oximetry data in newborns, allowing us to analyze other pulse oximetry features that may augment the current screening process when added to the SpO<sub>2</sub> screening component.

As such, it is necessary to design an interpretable machine learning model that can be directly incorporated into our current SpO<sub>2</sub>-alone screening system [9] with automatic feature selection to further improve the sensitivity of CCHD detection with little impact on specificity (at least 99%). To our best knowledge, our work is the first to analyze the feature relevance of CCHD screening by using machine learning (ML) algorithms as well as the first to incorporate ML into current standard SpO<sub>2</sub> screening.

## II. METHODS

### A. Subjects

By using an automated collection system [8], we enrolled 335 newborns, including 236 newborns that have a final diagnosis (with or without CCHD) confirmed. Patients were excluded if they required vasoactive infusions other than Prostaglandin E<sub>1</sub>. The majority of patients were enrolled at University of California (UC), Davis. Four other hospitals, Sutter Medical Center in Sacramento, UC Los Angeles, UC San Francisco, and Cohen Children's Medical Center in New York, are also enrolling patients for this study. We recorded at least 5-minute dual limb (right hand and any foot) pulse oximetry measurements at three time periods: within 24 hours, 24-48 hours, and after 48 hours following the baby's birth. We analyzed healthy newborns (defined as those without any CHD) vs those with CCHD (newborns who require a surgical or catheter-based intervention within 30 days of age). We divided all measurements into two groups: (G1) 0-48 hours, which included 158 healthy and 27 CCHD newborns; and (G2) over 48 hours, which included 50 healthy and 36 CCHD newborns. The experimental procedures involving human subjects described in this paper were approved by the Institutional Review Board.

<sup>1</sup>Z. Lai, C.-N. Chuah are with the Electrical and Computer Engineering Department, University of California Davis, Davis, CA, USA. {lzhengfeng, chuah}@ucdavis.edu

<sup>2</sup>P. Vadlaputi, D. Tancredi, S. Lakshminrusimha, H. Siefkes are with the Pediatrics Department, University of California Davis, Sacramento, CA, USA. {ppvadlaputi, djtancredi, slakshmi, hsiefkes}@ucdavis.edu

<sup>3</sup>M. Garg is with David Geffen school of Medicine UCLA, Los Angeles, CA, USA. mgarg@mednet.ucla.edu

<sup>4</sup>R. I. Koppel, M. Goodman are with Cohen Children's Medical Center, NY, USA. {rkoppel, Mgoodman3}@northwell.edu

<sup>5</sup>W. Hogan is with University of Utah Health, Salt Lake City, UT, USA. whitnee.hogan@gmail.com

<sup>6</sup>N. Cresalia, S. Juergensen are with University of California San Francisco, San Francisco, CA, USA. {Nicole.Cresalia, Stephan.Juergensen}@ucsf.edu

<sup>7</sup>E. Manalo is with Sutter Health, Sacramento, CA, USA. manaloe@sutterhealth.org

## B. Spot SpO<sub>2</sub>-alone Screening

Single pre and post-ductal SpO<sub>2</sub> values were recorded during the measurements as well. These spot values were used to assign pass or fail to all newborns per the current standard SpO<sub>2</sub>-alone screen [9]. If the last recorded spot SpO<sub>2</sub> measurements resulted in a "repeat" assignment from this algorithm, we assigned them a "fail" in a "Conservative spot SpO<sub>2</sub>-alone" algorithm to bias towards the null for CCHD detection.

## C. Features Extraction & Analysis

Pulse oximetry features evaluated for discrimination of healthy vs CCHD include: heart rate (HR), perfusion amplitude index (PAI), also known as perfusion index (PIx), and oxygen saturation (SpO<sub>2</sub>). We removed values likely associated with artifact: HR larger than 250 and SpO<sub>2</sub> larger than 100 (the pulse oximeter assigns a value of 127 for SpO<sub>2</sub> when the measurement quality is poor). We then extracted variance, min, max, median and mean for HR, PAI and SpO<sub>2</sub>. To study the differentiation of these features, we visualized each individually and then their correlation with each other. Fig. 1 illustrates the distribution of the mean SpO<sub>2</sub> and its correlation with min HR: the mean SpO<sub>2</sub> for the healthy newborns is typically higher than that for the CCHD newborns and the min HR for healthy newborns is typically less than that for CCHD newborns.

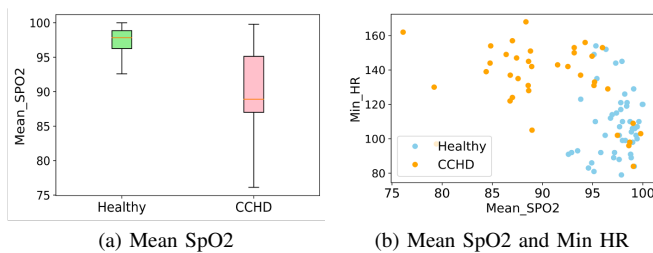


Fig. 1. Feature Analysis between Healthy vs CCHD Over 48 Hour of Age.

## D. Classification Algorithms

Several ML classifiers were tested for CCHD detection during this study, including Random Forest, Logistic Regression, and Multilayer Perceptron. Although a recent study [10] claimed random forest classifier has the best performance, we comprehensively investigated the above classification algorithms by Recursive Feature Elimination (RFE) [11] with 5-fold cross-validation on each algorithm separately. RFE can help determine the best performance of each model and the corresponding optimal feature set. RFE achieves this by searching for the most relevant subset of features to optimize our performance metric. We used cross-validation to optimize sensitivity by setting it as the score of RFE. To achieve the most optimal sensitivity, we started with all features from the training dataset as the input and fit the ML models, which ranked features by importance, discarded the least important features and refit the model. This process was repeated until the desired number of features resulted in the highest sensitivity (as shown in Fig. 2).

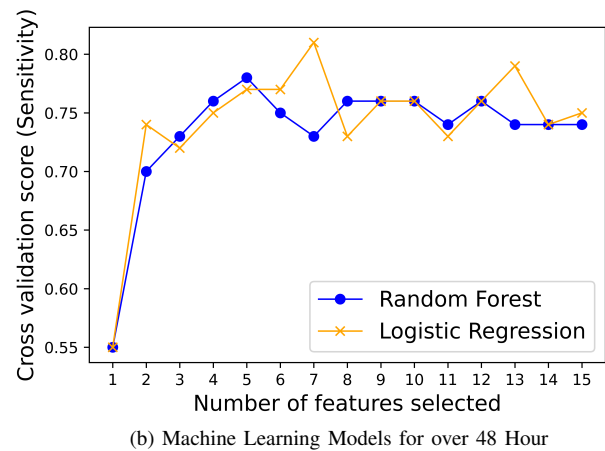
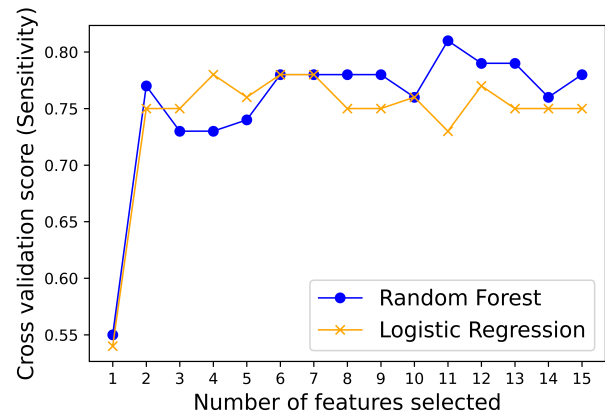


Fig. 2. Recursive Feature Elimination (RFE) by using Machine Learning Models for Healthy vs CCHD at Different Ages.

1) *0-48 Hours of age to evaluate for no-CCHD vs CCHD:* In this setting, we found that Random Forest Classifier resulted in the highest sensitivity in the 5-fold cross-validation for the 0-48 hour age. Fig. 2(a) shows the optimal sensitivity was achieved by using 11 features.

2) *Over 48 Hours of age to evaluate for no-CCHD vs CCHD:* We found that Logistic Regression performed best in this group, as Fig. 2(b) demonstrates, we found a subset of 7 optimal features using Logistic Regression.

The fewer number of CCHD cases compared to healthy cases resulted in a class imbalance problem, making it challenging to train a relatively unbiased ML model. To deal with this issue, we applied Random Oversampling (ROS), Random Undersampling (RUS), and Synthetic Minority Oversampling Technique (SMOTE) [12], separately. SMOTE is one type of oversampling method, where the minority class is oversampled by generating "composite" examples [12]. We tried different sampling ratios from 0.2 to 0.9, but the results had trivial improvements. Thus, we used the balanced loss function inside of the ML models, which led to the optimal results we could achieve in this study.

## E. Performance Evaluation

CCHD screening is a binary classification problem between healthy and CCHD, thus we used the following metrics

to comprehensively evaluate the performance of our model: Sensitivity (Sens) (1) and Specificity (Spec) (2).

$$\text{Sens} = \text{TPR} = \frac{TP}{TP + FN} \quad (1)$$

$$\text{Spec} = \frac{TN}{TN + FP} \quad (2)$$

$$\text{FPR} = \frac{FP}{FP + TN} \quad (3)$$

where:

- TP: the number of CCHD predicted as CCHD
- FP: the number of healthy predicted as CCHD
- TN: the number of healthy predicted as healthy
- FN: the number of CCHD predicted as healthy

We also calculated the Area Under the Receiver Operating Characteristics curve (AUROC) by plotting true positive rate (TPR) (1) against false positive rate (FPR) (3) with the discrimination threshold increasing from 0 to 1.

### III. RESULTS

#### A. Sensitivity and Specificity

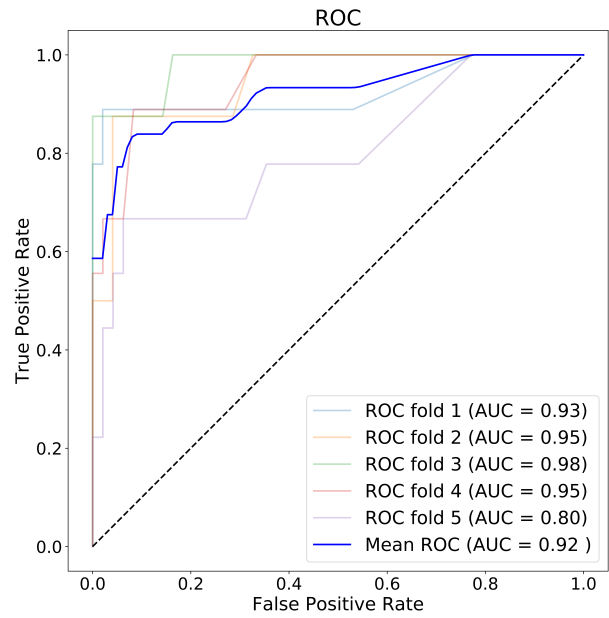
TABLE I  
COMPARED WITH CURRENT SPO2 SCREENING

Methods	Specificity	Sensitivity
True Spot SpO2-alone [9]	96.8	62.8
Conservative Spot SpO2-alone [9]	96.8	76.5
ML (0-48 hrs)	<b>97.5</b>	<b>81.5</b>
ML (>48 hrs)	<b>100</b>	<b>83.3</b>
Conservative Spot SpO2-alone + ML (0-48 hrs)	96.2	<b>85.2</b>
Conservative Spot SpO2-alone + ML (>48 hrs)	96	<b>88.9</b>
Conservative Spot SpO2-alone + Any ML	95.8	<b>86.4</b>

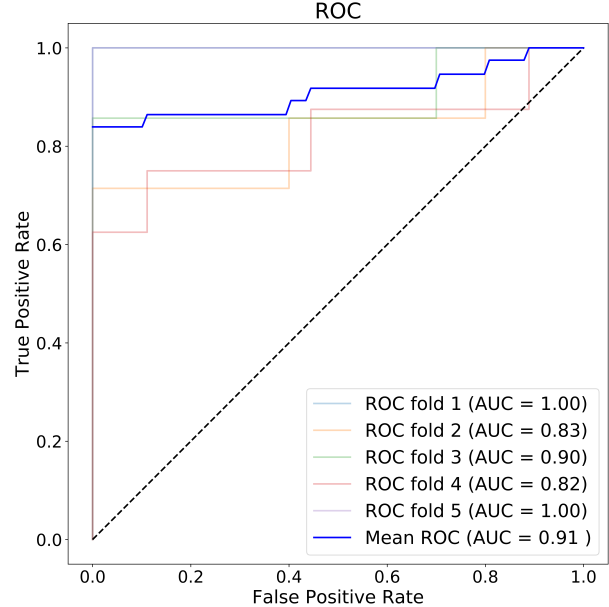
The sensitivity and specificity of our proposed method and the current SpO2 screening methods are summarized in Table. I. When comparing to the current standard CCHD screening, which includes spot right hand and any foot SpO2 measurements, the addition of the over 48 hour ML model did not lead to additional false positive results, but did detect 4 additional newborns with CCHD, hence increasing sensitivity from 76.5% to 88.9% (McNemar mid-p = 0.06) when tested on a sample of 50 healthy newborns and 36 newborns with CCHD. The addition of the 0-48 hour ML model, resulted in 3 additional false positive results and detected 3 additional newborns with CCHD, hence increasing sensitivity from 76.5% to 85.2% (McNemar mid-p = 0.13) when tested on a sample of 158 healthy newborns and 27 newborns with CCHD. Overall, the ML models improved the sensitivity from 76.5% to 86.4%, nearly 10 percentage point improvement.

#### B. ROC Curve

ROC curve results from 5-fold testing are shown in Fig. 3. For 0-48 hours, the average AUROC for no-CHD vs CCHD was 0.92 by using the Random Forest classifier; for over 48 hours, the average AUROC was 0.91 by using the Logistic Regression model. The estimated AUROC for our



(a) 0-48 Hours



(b) Over 48 Hours

Fig. 3. Area Under the Receiver Operating Curves (AUROC) for Models on No-CHD vs CCHD.

ML algorithms combining pulse oximetry features (PAI and HR) appears similar to or better than the current SpO2-alone screen [13], [14].

#### C. Interpretable Machine Learning

Compared to the current standard SpO2-alone CCHD screening, our ML models have the potential to achieve better sensitivity by incorporating features related to HR and PAI (or PIx). The optimal subset for the 0-48 hours Random Forest classifier includes: HR (median, mean, max, variance), SpO2 (min, max, median, mean), PAI or PIx (mean, median, max). The optimal subset for the over 48 hours Logistic Regression ML model includes: HR (min, max, variance),

