Exploring Edge Machine Learning based Stress Prediction using Wearable Devices

Sang-Hun Sim Arizona State University ssim6@asu.edu Tara Paranjpe
Arizona State University
tparanjp@asu.edu

Nicole Roberts Arizona State University nicole.a.roberts@asu.edu Ming Zhao Arizona State University mingzhao@asu.edu

Abstract—Stress is a central factor in our daily lives, impacting performance, decisions, well-being, and our interactions with others. With the development of IoT technology, smart wearable devices can handle diverse operations, including networking and recording biometric signals. It has also increased stress awareness among users with the wearables' enhanced data processing capability. Edge computing on such devices enables real-time feedback and in turn preemptive identification of reactions to stress. This can provide an opportunity to prevent severe consequences that might result if stress is left unaddressed. Edge computing can also strengthen privacy by implementing stress prediction on local devices without transferring personal information to the public cloud.

This paper presents a framework for real-time stress prediction, specifically for police training cadets, using wearable devices and machine learning with support from cloud computing. It includes an application for Fitbit and the user's accompanying smartphone to collect heart rate fluctuations and corresponding stress levels entered by users and a cloud backend for persisting data and training models. Real-world data for this study was collected from police cadets during a police academy training program. Machine learning classifiers for stress prediction were built using this data through classic machine learning models and neural networks. To analyze efficiency across different environments, the models were optimized using model compression and other relevant techniques and tested on cloud and edge environments. Evaluation using real data and real devices showed that the highest accuracy came from XGBoost and Tensorflow neural network models, and on-edge stress prediction models produced lower latency results.

I. INTRODUCTION

Stress often occurs when we become pressured by events that overwhelm our capacity to deal with them [1]. This can be particularly consequential in the workplace. Stress is inversely related to job performance and negatively affects decision-making processes [2]. Law enforcement work in particular comes with many unexpected threats and challenging demands. According to a study conducted by Griffin et al. [3], chronic stress increases the perception of threats and the aggression of responses to those threats. It becomes immediately clear that identifying and addressing symptoms of acute stress among officers is critical, as stress can otherwise become chronic, leading to loss in job performance, poor officer health, and potentially problematic behaviors [4]. Furthermore, despite wide variation from person to person [5], increased cardiovascular activation (e.g., heart rate) generally is associated with more intense physical and emotional states, including stress [6].

There is utility in understanding whether frequent or cumulative increases in heart rate can predict feelings of stress; which features or properties of the heart rate signal are predictive; and the time course for optimal predictions. This has real-world implications for police officers, who must be aware of and manage their stress and physiological arousal in real time. To address and ultimately mitigate the potentially detrimental effects of stress in police officers as they navigate their ever-changing field, we investigated stress prediction through the collection of biometric data via wearables, used that data to create machine learning models, and analyzed the success of those models on the edge.

As described in this paper, we leveraged machine learning and edge computing techniques for stress prediction. The machine learning approach allowed us to draw relevant predictions from a large volume of data collected over time. With the growth of machine learning techniques and models, this paper investigates the use of classic models as well as modern neural networks, giving us the ability to generalize results beyond our dataset. Edge computing, combined with a machine learning approach, creates a faster and more reliable stress prediction mechanism. With all the computing power moved to the local devices, users can achieve strong data privacy and fast prediction speed.

We developed a framework for real-time stress prediction using wearable devices and machine learning with support from cloud computing. We utilized Fitbit, a commercial wearable device, to collect biometrics including heart rate data and developed a custom application to collect stress level ratings that participants could manually input into the Fitbit. We used Amazon Web Services (AWS) to securely store data and train machine learning classifiers. Our custom Fitbit application generated prompts, triggered by heart rate fluctuation thresholds (35% increase above resting heart rate, based on previous research and pilot testing), with 5 stress level buttons to collect user input, as inspired by empirically-proven strategies. We created a secure web application with a dashboard to allow users and researchers to monitor each Fitbit device and view all the collected data.

We worked closely with a local metropolitan police training academy to collect biometric data and stress perception scores from police training cadets. Data was collected over the course of 4 months across 15 different cadets. Our team worked with each individual cadet to prepare and maintain their data

collection from Fitbit devices and to gather overall study feedback.

Using the data collected, we built classification models to categorize stress experiences from features of heart rate using machine learning. We employed binary classification methods (not stressed vs. stressed) using both classic models and a neural network. Then, we deployed the stress classification models on smartphones to predict stress ratings in proximity to the user without sending the data to the cloud. To optimize the neural network for edge use, we applied quantization to reduce the model size.

The most significant results of our study are as follows:

- The best accuracies for the stress prediction models were 96.98% from XGBoost and 95.98% from a Tensorflow neural network using an over-sampled data extraction approach.
- Quantization reduced the stress prediction model size by 2x and maintained the accuracy as before the quantization.
- Several trials proved that on-edge stress prediction with the quantized model saved nearly 200ms compared to incloud stress prediction.

The rest of the paper is organized as follows: Section 2 investigates the related works; Section 3 presents our data collection strategies; Section 4 details our approach to machine learning based stress prediction; Section 5 presents our edge computing based stress prediction solution; and Section 6 concludes the paper.

II. RELATED WORKS

A. Efficacy of Wearable Devices in Clinical Domains

Wearable devices as a health monitoring and management tool are increasing in popularity and functionality. Seamless integration with smartphone devices and real-time data collection have made it simple and convenient for users to learn more about their health, mental, and physiological states without any external involvement. As wearable devices continue to grow in popularity, studies have begun to leverage their health data collection capabilities to learn more about human behaviors and tendencies.

A study conducted by Beniczky et al. [7] investigated the use of wearable devices to detect and predict seizures for patients with epilepsy. Through the use of the sensors on the device, the researchers were able to collect electrocardiogram (ECG), heart rate variability (HRV), and accelerometer data. These data were then fed into machine learning models to create predictive models for generalized tonic-clonic seizures. Results showed that the noninvasive wearable devices were able to detect such seizures with 90%-96% accuracy. Rykov et al. [8] were able to leverage wearable devices in the same capacity to screen for depression-related biomarkers, collecting sleep patterns, physical activity, and psychological measurements to evaluate mental states. Models created in this depression detection domain were able to reach 82% accuracy.

The ability to make diagnoses and predictions like these in real-world contexts rather than typical evaluations in controlled environments demonstrate the feasibility of using wearable devices to extend clinical research. It becomes clear that wearable devices have capabilities that far surpass traditional data collection techniques, allowing for data collection in broader and more extensive contexts.

B. Stress Prediction using Machine Learning and Wearable Devices

With the promising results from wearable devices in clinical domains, several studies have leveraged data collected from wearable devices in machine learning based stress prediction. Lawanont et al. [9] and Dai et al. [10] both investigated the use of wearable devices to detect stress in controlled stress scenarios. Using the Fitbit and the Fossil Gen4 Explorist, commercially available wearable devices, both studies collected several different data points during controlled experiments. Some inputs included heart data—specifically root mean square of successive differences (RMSSD) and inter-beat interval metrics, sleep data-including rapid eye movement (REM), deep and light sleep time, and other collected metrics like calories, steps, and intensity of movement. Can et al. [11] similarly leveraged wearable devices in the stress prediction domain, collecting heart rate variability data, accelerometer data, and temperature data, but took the research a step further: attempting to detect stress in situ. Research in real-life contexts is rather rare; stress can be all-consuming. Simply put, as researchers, we do not have much context when we gather data in daily-life contexts, and rely on self-reported information as ground truth. Our work is among the first of its kind to investigate stress prediction in real-world settings.

These studies employed classification-based models to detect stress. Among these, the Lawanot et al. study [9] achieved an 84.10% accuracy with the Decision Tree model (no F1-score reported), the Dai et al. study [10] achieved an 82.3% accuracy and 62.3% F1-score with the Support Vector Machine (SVM) model, and the Can et al. [11] study achieved 92.19% accuracy and 90.30% F1-score among other significant results. Through such key findings, we can learn that there is great promise in leveraging empirically-chosen classification models to predict stress-controlled situations.

Furthermore, with increased demand for stress management in law enforcement domains, a handful of studies involving wearable devices, machine learning, and behavior prediction have been conducted. Tiwari et al. [12] proposes stress prediction through the use of heart rate variability, breathing analyses from three varying "waves" of data collection. Each wave consisted of a different assessment—beginning with daily wear to gauge baseline levels, and moving on to shooting range exercises and intervention simulation exercises. Similarly, in Erickson et al. [13], heart rate and sleep data were collected through in-class, day, and night field training. Both studies used the collected data as inputs into classification models, including SVM (most accurate for [12]), Logistic Regression, Random Forest (most accurate for [13]), and Adaboost. Though there are stress prediction and detection studies in

the law enforcement field, most data is collected in controlled environments with an understanding of the adjoining context.

Built upon the findings of these stress prediction studies in both general and law enforcement domains, our work makes several new contributions: 1) we developed an end-to-end system with a custom wearable app to prompt users for real-time inputs based on physiological signals; 2) we used this system to collect multiple months of data from real users in real-world scenarios; and 3) we proposed the use of edge computing for real-time stress prediction.

C. Machine Learning on Edge Devices

Due to the immediate and dynamic nature of stress and how it manifests, it is important to look into studies that implement real-time, on-edge machine learning to guide our stress prediction initiative.

Chen et al. [14] and Mauldin et al. [15] provided in-depth research about the capability of mobile devices for training and inferencing with deep learning models. Chen et al. leveraged the CIFAR-10 dataset, containing millions of images used to train computer-vision models. These data were fed into various deep learning models, significantly, the Convolutional Neural Network (CNN). Mauldin et al., proposing a real-time fall detection system, leveraged external wearable datasets in several models, including SVM, Naive Bayes Classifier, and a Deep Neural Network (DNN) which was the most accurate. The former study found that training operations contribute most significantly to the latency, especially for the gradient calculation of the backward path. As a result, the study concluded that it is possible to run both training and inference on mobile devices provided that models' complexity is reduced. The latter study similarly concluded that it is possible to run machine learning inference on mobile devices, provided a lightweight model. It becomes clear from these works that mobile-inference is possible in real-time machine learning deployments.

Ogden et al. [16] and Guo et al. [17] delved deeper into onedge vs. in-cloud latency with their respective deployments. Ogden et al. mainly focused on which model compression techniques to use, which model to choose for mobile inference, and when to depend on servers. They proved that the quantized models have significantly smaller model size; moreover, loading the 8-bit quantized model did not contribute significantly to inference time, compared with other models. Additionally, through an analysis of various devices and internet connection strength, they determined that mobile-based inference only proved effective on high-performance smartphones, while onedge inference was more reliable, with lower latency observed even in poor network conditions.

Guo et al. [17] tested on-edge vs. on-cloud latency, but observed varying results. The evaluation metrics used in this study were latency, power consumption, and resource usage. They found that the on-cloud approach outperforms on-edge. They share, however, that this conclusion is dependent on the performance capabilities of the smartphone and the model size. With manipulation of these two parameters, there is a

high possibility for low latency through on-edge deployment. Though both conclusions differ, this lays the groundwork for additional research between on-edge and in-cloud model deployments to improve real-time feedback.

Built upon these edge-computing-based machine learning studies, this paper proposes a new edge computing based solution for stress prediction and addresses the unique challenges brought by developing an accurate and fast model for predicting stress in real time from real users' biometric inputs.

III. DATA COLLECTION

A. Overview

This project collected heart rate and stress rating data from Fitbit devices worn by 15 police cadets during four continuous months in a rigorous training academy. Stress ratings were collected directly into the Fitbit when prompted based on increases in the user's heart rate. Stress level options ranged from one to five, with one as "not at all stressed" and five as "extremely stressed". We recorded a zero when the cadet did not respond to the prompt. This scale was chosen through a widely adopted psychological instrument called the Percieved Stress Scale (PSS) [18] which has proven to be an effective measurement for stress perception.

There were two types of data collection pipelines for heart rates and stress responses, respectively. On one hand, heart rate data was transmitted to the Fitbit's remote server by itself (once a Fitbit gets synchronized to the smartphone), and could be retrieved via Fitbit Web API with the credential information of the user account. On the other hand, stress levels needed the private database as it was a metric measured and collected by the custom application. These were stored in key-value format. After uploading and processing of the data, we created a data archive containing all cadets' heart rates and stress responses to begin creating the machine learning models.

B. Fitbit Application

Fitbit's network architecture includes the Fitbit watch and the companion application, a supplementary runtime environment built to extend application capabilities. Fitbit needs to be paired with the companion (in this case, a smartphone) through Bluetooth connection in order to transfer Fitbit-collected data. Because Fitbit does not have an Internet connection by itself, it depends on the companion for any other operations, like fetching information or storing JSON data, except for recording bio-signals. Upon syncing to the companion, Fitbit sends all the biometric signals recorded to the remote server. Fitbit also provides an application on the smartphone to display statistical data of the recording. While the Fitbit focuses mainly on recording biometric measurements, the companion can do more complex operations, such as fetching data or importing external packages. The Fitbit and companion together can build a more elaborate application.

Fitbit supplies a Software Development Kit (SDK) and various Application Programming Interfaces (API) to developers for custom application development. As this initial study intended to catch stress occurrence mainly by heart

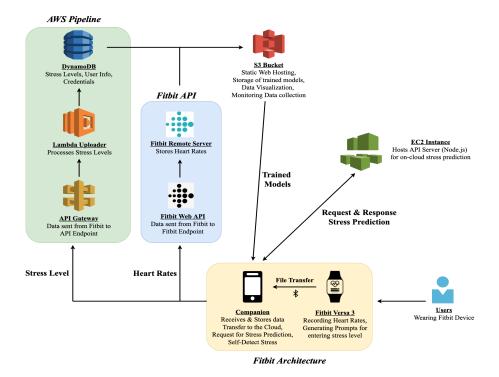


Fig. 1. The Architecture of the Stress Management

rate fluctuation, we built a Fitbit application to detect patterns indicating stress by prompting stress level buttons as shown in Figure 2. The prompting generation mechanism is when the heart rate goes above the resting heart rate by 35 percent for two minutes, the clock's face changes to the prompted face with five buttons ranging with five levels, "No Stress", "A Little", "Moderate", "A Lot", and "Extremely". Users enter their stress levels subjectively. A 30-minute pause period was programmed to occur between prompts to avoid continuous prompting. If a user is unable to respond to a given prompt after 7 minutes, the value is stored as a "missed" prompt in our database.

Through this custom application, once the user enters the stress level, Fitbit keeps the stress input data in the file storage in the CBOR format and sends it to the companion as soon as the Bluetooth connection is established. The companion receives the stress file and concatenates it to the existing data in a key-value format, so the data is temporarily preserved in the companion.

The companion is responsible for sending the key-value dataset to the cloud. When the companion receives the file from Fitbit, it converts the data into a JSON object using the Fitbit File API. Then it sends the dataset to the cloud database using Fitbit Fetch API. The Fetch API allows developers to make GET and POST requests to HTTPS endpoints. To distinguish each cadet's data from the others, the Fitbit-provided unique device identification (ID) number was used as the partition key in the database.

Fitbit automatically stores biometric data, through the companion, in its remote server, and we can retrieve this data



Fig. 2. This is the clock face of the custom-built Fitbit application. There is a default face for general use and a prompted screen with a cycle button for stress inputs from users.

using the Fitbit Web APIs. However, it requires specific authorization credentials, including access tokens and user identification. To validate the access tokens, users must go through the authorization code grant process, which depends on OAuth 2.0, the protocol to allow a third-party user to access the resources. OAuth application creates the client ID and secret for use to invoke the access token. Each access token is only valid for 8 hours before it needs to be refreshed again.

C. Web Application

To handle multiple devices and data collection pipelines, we created a web application dashboard. The web application enables monitoring of the data collection for each user and shows the device information on the dashboard, such as battery status and the last synced time to their smartphone. It also provides tools for data retrieval and processing, including

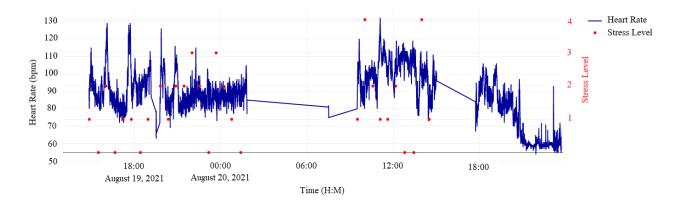


Fig. 3. As seen on our web application, this graph is a visualization of 24 hours of heart rate and stress data from one user. The blue lines represent heart rate in bpm, and the dots represent stress prompts, ranging from no-answer (on the axis line) to 4.

downloading interfaces that convert JSON data to CSV format and visualization page to analyze the heart rate data with stress inputs, as shown in Figure 3.

We implemented this web application using a serverless static web application through AWS. The benefit of static web hosting is that it minimizes the initial cost and eliminates the need for a hosting server, such as an EC2 instance. Instead, the web application runs in the S3 bucket, where HTML, CSS, and JS files are stored. The bucket has a feature to host an application with a designated AWS domain. However, since Fitbit only allows the HTTPS protocol to communicate with the outside, we purchased the private domain for the website and connected it to the S3 bucket using Route53, CloudFront, and Certificate Manager.

We employed DynamoDB, a non-relational key-value NoSQL database, to store all of our collected data. The key-value data is easy to query by the partition key or sort key, like extracting heart rate values for a specific user given start and end dates. Next, it is a non-relational database that allows scaling both vertically and horizontally, which is important to support data generated continuously from many wearable devices.

Most importantly, we utilized AWS Lambda functions to allow our static web application to act as a standard server. AWS Lambda is a serverless computing service that executes code without establishing a server. It is event-driven, executed only when the service is requested. Lambda helps the application build data processing functions by accessing other resources provided by AWS such as DynamoDB and S3. It can handle up to 250MB of code, and the execution time cannot be more than 15 minutes, which is sufficient for our needs. In the case of our application, we built lambda functions for creating and logging in user accounts and uploading and retrieving recorded data to the interface.

We also utilized the API Gateway to build APIs. Our application runs in RESTful APIs, which requests and responds in JSON format with four methods, including CREATE (post), READ (get), UPDATE (put), and DELETE (delete).

AWS's API Gateway provides users in the backend with the endpoint to access data in other services, such as DynamoDB. It also controls authentication to filter out unidentified requests. In our application, we built endpoints and mapped them to Lambda functions so that end-users can access and upload data. More specifically, the companion uploads the stress level data through an API gateway endpoint to DynamoDB.

Note that in this research, we obtained the cadets' permissions to store their data in the cloud and use the security mechanisms provided by AWS to protect the data. In real use scenarios, users may not permit their private data to leave their personal devices. Our system can be extended to employ privacy-preserving machine learning methods such as federated learning to support such scenarios.

IV. MACHINE LEARNING BASED STRESS MANAGEMENT

A. Problem Definition

As stated earlier, this research aimed to build a real-time stress prediction solution using machine learning. Our initial goal was to detect the patterns of heart rate that could indicate stress. Heart rate was collected through Fitbit Versa 3, and the user provided real-time stress perception scores using the device. Our dataset, comprised of these data points, proved to be suitable for supervised learning algorithms, due to its labeled nature. The labels were also categorical variables, so we decided to employ classification algorithms to determine the input instances as "stressed" or "not stressed". Table I shows the distribution of stress levels from all the prompts that the users received. Stress level 0 indicates no response at the time of the stress prompt, while 1-5 indicates the increasing scale of stress.

B. Data Processing

As mentioned above, our work is appropriate for classification models of supervised learning given its labeled nature. We began by extracting heart rate segments from the original, continuous heart rate data for each user. To capture heart rate leading up to the stress prompt (programmed to generate when the current heart rate exceeds 35% above the baseline

Stress level	Numerical value	Number of responses
No Response	0	9578
Not Stressed	1	2935
A Bit Stressed	2	355
Moderate	3	89
A Lot	4	19
Extremely Stressed	5	5

TABLE I DISTRIBUTION OF STRESS RESPONSES

for 2 minutes), we extracted the 2 minutes of data prior to the prompt. Each heart rate segment was then mapped to corresponding stress levels to label each of them as stressed or unstressed. We converted the stress level to "0" for unstressed and "1" for stressed for binary classification. Since we let users indicate their stress in five different levels, we can map them into "unstressed" and "stressed" in different ways: 1) To include dat with stress levels from 2 to 5 as "stressed", and 2) To include data with stress levels from 3 to 5 as "stressed". The reason for these two options is because of the ambiguity of the stress level 2. Therefore, we processed the datasets in two ways separately according to how to map the stress levels.

The pre-processed raw dataset includes the prompting threshold, resting heart rates, and 40 samples of heart rates. The granularity of heart rates recorded by Fitbit is between 5 to 10 seconds, so the 40 samples of heart rates represent about 2 to 3 minutes of data prior to the prompt. Since the resting heart rates are calculated by Fitbit automatically, we did not have to go through feature extraction steps to get those values.

The next step is feature extraction. Since we set the prompting algorithm based on how heart rate fluctuates compared to the resting heart rate, we extracted the features representing distribution and fluctuation-related aspects of heart rate. There were seven features extracted from the heart rate. The primary five were 1) mean, 2) standard deviation, 3) minimum and 4) maximum value of the heart rate, and 5) resting heart rate. The additional two features were 6) the difference between mean heart rate from the resting heart rate by percentage (DiffRest) [9] and 7) the root mean square of successive difference between normal heartbeats (RMSSD) which Fitbit uses for calculating heart rate variability from heartbeats.

One issue we faced was the unbalanced nature of the dataset. Over 85% of the dataset contained unstressed instances. To avoid a learning bias across our models, we explored two approaches for dataset balancing: 1) To under-sample the majority class by randomly choosing as many instances as the number of the minority, and 2) To over-sample the minority instances using Synthetic Minority Over-Sampling Technique (SMOTE) [19]. The former method actually reduced the size of the data, proving to be ineffective for training sessions; however, it presented the ability to learn both stressed and unstressed equally. The latter increased the dataset size and the effectiveness of the learning stage, but created similar but synthetic sampled instances. Lastly, to ensure all attributes had equal, unbiased influence, we implemented standard scaling on all the dataset columns to make them have the same

distribution with 0 means and the unit standard deviation.

C. Training Models

Approach	"Stressed" levels	"Not Stressed" levels	Resampling method
1	2 to 5	1	Under-sampling
2	2 to 5	1	SMOTE
3	3 to 5	1 to 2	Under-sampling
4	3 to 5	1 to 2	SMOTE

TABLE II FOUR APPROACHES FOR TRAINING MODELS

We considered four approaches for training models as specified in Table II. We trained models separately by how we defined "stressed" (level 2 to 5 or 3 to 5) and how we balanced the dataset (under-sampling or over-sampling). Regarding classification algorithms, we considered four classic models and TensorFlow's feed-forward network. The set of classic models include Decision Tree, Random Forest, Adaboost, and XGBoost. The first three models were imported from Scikit-learn (Sklearn) [20], which provide not only reliable classification models but also useful built-in functions for the machine learning process, such as *train_test_split* or *standard scalaer*. Before running the dataset for training, we split the dataset into a training set (80%) and a testing set (20%) to evaluate models' capability to handle unseen data.

The TensorFlow feed-forward network was developed with six hidden layers activated by the ReLu function. The output layer was set with a Sigmoid function returning a value ranging from 0 to 1. If the final return value from an input instance is lower or equal to 0.5, it is unstressed, which is labeled as 0. Otherwise, it is stressed labeled as 1. Since we used binary classification, we used binary cross-entropy for the loss function.

D. Evaluation

Table III shows the accuracy and F1 scores for each approach. Approach 1 gave us the worst result. The models could not distinguish between unstressed and stressed by heart rate features. This might be because of either insufficient number of instances or ambiguity of features from stress level 2. We could gain better results from Approach 2, which oversampled the minority class. Approach 3 showed more reliable results than Approach 1 by not using over-sampling with synthetic samples. The dataset size was reduced even more than Approach 1, as it excluded stress level 2 and re-sampled the unstressed data as many as the number of stress levels ranging from 3 to 5. Although it does not have a sufficient dataset, it showed pretty good accuracy and F1 scores by properly inferencing on the testing set. Lastly, Approach 4 gave us the best results. It excluded stress level 2 and oversampled the minority class. The best accuracy it reached is 96.98% from XGBoost and 95.98% from the Tensorflow neural network.

	Model	Accuracy	F1 Score
Approach 1	Decision Tree	61.23%	60.51%
	KNN	60.14%	58.33%
	Random Forest	64.85%	65.72%
	AdaBoost	60.86%	58.33%
	XGBoost	70.28%	70.50%
	Neural Networks	56.52%	56.52%
Approach 2	Decision Tree	85.29%	85.63%
	KNN	81.08%	83.18%
	Random Forest	87.63%	88.14%
	AdaBoost	71.22%	73.37%
	XGBoost	88.60%	88.75%
	Neural Networks	87.77%	88.38%
Approach 3	Decision Tree	77.27%	76.92%
	KNN	77.27%	72.72%
	Random Forest	81.72%	81.83%
	AdaBoost	72.12%	71.64%
	XGBoost	78.78%	78.12%
	Neural Networks	86.36%	88.46%
Approach 4	Decision Tree	91.79%	91.73%
	KNN	88.34%	88.73%
	Random Forest	94.41%	94.34%
	AdaBoost	79.25%	80.15%
	XGBoost	96.98%	96.95%
	Neural Networks	95.98%	96.02%

TABLE III
EVALUATIONS BY ACCURACY AND F1-SCORE

V. EDGE COMPUTING FOR STRESS PREDICTION

A. In-Cloud vs. On-Edge Stress Prediction

Edge computing enables operation of classification models on edge device such as smart wearables and smartphones, alleviating the burdens of reaching the server in cloud. Before we established the edge-based stress prediction, as a baseline, we implemented in-cloud stress prediction by running pretrained models on EC2 instances with the input from the companion. On-edge stress prediction, in comparison, uses the companion, instead of the cloud server, to load models and predict stress in order to save the communication time between the companion and server.

Our Fitbit app always maintains the latest 2 to 3 minutes of heart rates using the queue data structure. When it needs to request stress prediction, it generates an input instance by calculating the seven features specified in Section B. It sends a request with a feature vector to either the companion or cloud after extracting the features.

For on-edge stress prediction, the companion loads the pretrained models and runs the inference. Unlike typical Android applications, the Fitbit companion cannot load files from the device's storage, and thus it has to retrieve the pre-trained models from the cloud. Specifically, the companion uses the Fitbit Fetch API to retrieve the models from the AS S3 buckets via HTTPS. There is an initial overhead when the companion loads the models from cloud into its memory when it performs the first inference; the following inferences can directly use the models that are already in memory (until the app restarts). Because the smartphone has limited computing power and memory capacity, we considered model optimization techniques such as quantization for reducing the complexity and size of the models. Since the Fitbit app and the companion run in JavaScript, we had to convert the models, originally written in Python for TensorFlow, to JavaScript versions, utilizing the TensorFlow library. Among the four approaches with which we experimented, we chose Approach 4's model, which achieved the highest accuracies and F1 scores.

For in-cloud stress prediction, we built an API server that handles requests from Fitbit, runs model inference, and returns the prediction result. Fitbit communicates with the server using the Fitbit Fetch API to upload, retrieve, and manipulate data from our storage tables. The companion is merely a bridge redirecting requests and responses. When requesting stress prediction, Fitbit uses the *POST* method of the Fetch API to send the inference input which contains the heart rate features. Upon receiving the request, the server proceeds with the stress prediction method using pre-trained models fetched from the S3 bucket. Since the cloud has sufficient resources, far more superior than the companion, the server-side models have not gone through model optimization processes.

Model	Type	Value
Original Model	Accuracy	85.98%
	Topology Size	0.005 MB
	Weight Size	3.0 MB
Quantized 8-bit Model	Accuracy	84.10%
	Topology Size	0.0067 MB
	Weight Size	0.774 MB

TABLE IV
ACCURACY AND THE SIZE OF MODELS

B. Evaluation

Since the companion has limited computing resources, we compressed the model to optimize for the edge device. We applied quantization to reduce the precision of model parameters to 8-bits from the original model. There are two approaches to quantize, post quantization and quantization-aware training. In this work, we used the quantization-aware training because it generally creates more accurate quantized models by considering the error introduced by quantization during the training process. Table IV shows the accuracy and the size of models. The topology size slightly increases from the original to the quantized model, but decreases by 74.2% for the overall weight size. Most significantly, the accuracy across both models only varies by 2%.

Figure 4 illustrates the stress prediction latencies and their breakdowns. We compared three scenarios: original model on-edge (Companion), quantized 8-bit model on-edge (Companion-8bit), and original model in-cloud (Server-orig). On one hand, the overall latency difference for both on-edge scenarios prove to be negligible, with small variation across data transfer steps. Nonetheless, the quantized model still saves X% of memory usage, which is still valuable for resource-limited smartphones. On the other hand, it can be clearly observed that the in-cloud deployment takes an extra 200ms on average, mainly contributed by the additional companion to server data transfer overhead. This result confirms that on-edge deployment of stress prediction produces lower latency than in-cloud deployment.

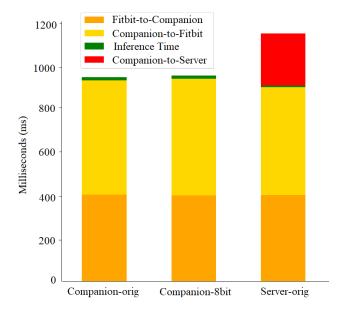


Fig. 4. This graph shows the latency during each phase for in-cloud and onedge stress prediction. For 'server-orig', it is clear that companion-to-server adds to the overall latency compared with companion-based models.

VI. CONCLUSION

In this paper, we studied the feasibility of edge machine learning based stress prediction using wearable devices with police cadets during a training academy. We focused on heart rate as an important starting point, and verified that heart rate recorded by commercial wearable devices can be effectively used for stress prediction. Unlike other equipment which records ECG data in milliseconds, Fitbit reads heartbeats in 5 to 10 second granularity, with the option of extracting in 1-second segments. However, with features extracted from the heartbeats and stress levels provided by participants, we found that Fitbit-based heartbeat data also can be a validated data type indicating stressful or non-stressful circumstances, even when measured under real-world conditions.

Next, we employed five classic classification models and one neural network for binary classification. After segmenting heart rate data into 2 to 3-minute windows and extracting five statistical features and two features representing heart rate variability, we could apply the dataset to machine learning algorithms. We also resolved the imbalance of the dataset by either under-sampling or over-sampling. The best accuracies were 96.98% from XGBoost and 95.98% from a 6-layer feedforward neural network using an approach that oversampled the minority instances and only used levels 3 to 5 as stressed, instead of including the 2.

Lastly, we optimized our neural network model for edge deployment (on the Fitbit companion) using quantization which reduced the model size by 74.2% and maintained the accuracy as before the quantization. On-edge deployment of the quantized model to the companion saved approximately 200 milliseconds, proving to be promising as we move forward.

VII. ACKNOWLEDGEMENT

We thank the anonymous reviewers for their feedback. This research is sponsored by U.S. National Science Foundation awards OAC-2126291 and CNS-1955593. We also thank Phoenix Regional Police Academy for their support of this study.

REFERENCES

- [1] R. S. Lazarus and S. Folkman, *Stress, appraisal, and coping*. Springer publishing company, 1984.
- [2] N. Atsan, "Decision-making under stress and its implications for managerial decision-making: A review of literature," *International Journal of Business and Social Research*, vol. 6, p. 38, 04 2016.
- [3] S. P. Griffin and T. J. Bernard, "Angry aggression among police officers," Police quarterly, vol. 6, no. 1, pp. 3–21, 2003.
- [4] —, "Angry aggression among police officers," *Police Quarterly*, vol. 6, no. 1, pp. 3–21, 2003.
- [5] K. Hoemann, Z. Khan, M. J. Feldman, C. Nielson, M. Devlin, J. Dy, L. F. Barrett, J. B. Wormwood, and K. S. Quigley, "Context-aware experience sampling reveals the scale of variation in affective experience," *Scientific reports*, vol. 10, no. 1, pp. 1–16, 2020.
- [6] S. Pieper and J. F. Brosschot, "Prolonged stress-related cardiovascular activation: Is there any?" *Annals of Behavioral Medicine*, vol. 30, no. 2, pp. 91–103, 2005.
- [7] S. Beniczky, P. Karoly, E. Nurse, P. Ryvlin, and M. Cook, "Machine learning and wearable devices of the future," *Epilepsia*, vol. 62, pp. S116–S124, 2021.
- [8] Y. Rykov, T.-Q. Thach, I. Bojic, G. Christopoulos, J. Car et al., "Digital biomarkers for depression screening with wearable devices: cross-sectional study with machine learning modeling," *JMIR mHealth and uHealth*, vol. 9, no. 10, p. e24872, 2021.
- [9] W. Lawanont, P. Mongkolnam, C. Nukoolkit, and M. Inoue, Daily Stress Recognition System Using Activity Tracker and Smartphone Based on Physical Activity and Heart Rate Data, 01 2019, pp. 11–21.
- [10] R. Dai, C. Lu, L. Yun, E. Lenze, M. Avidan, and T. Kannampallil, "Comparing stress prediction models using smartwatch physiological signals and participant self-reports," *Computer Methods and Programs* in *Biomedicine*, vol. 208, p. 106207, 2021.
- [11] Y. S. Can, N. Chalabianloo, D. Ekiz, and C. Ersoy, "Continuous stress detection using wearable sensors in real life: Algorithmic programming contest case study," *Sensors*, vol. 19, no. 8, p. 1849, 2019.
- [12] A. Tiwari, R. Cassani, J.-F. Gagnon, D. Lafond, S. Tremblay, and T. H. Falk, "Prediction of stress and mental workload during police academy training using ultra-short-term heart rate variability and breathing analysis," in 2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC). IEEE, 2020, pp. 4530–4533.
- [13] M. L. Erickson, W. Wang, J. Counts, L. M. Redman, D. Parker, J. L. Huebner, J. Dunn, and W. E. Kraus, "Field-based assessments of behavioral patterns during shiftwork in police academy trainees using wearable technology," *Journal of Biological Rhythms*, vol. 37, no. 3, pp. 260–271, 2022.
- [14] Y. Chen, S. Biookaghazadeh, and M. Zhao, "Exploring the capabilities of mobile devices supporting deep learning," in *Proceedings of the 27th International Symposium on High-Performance Parallel and Distributed Computing*, ser. HPDC '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 17–18. [Online]. Available: https://doi.org/10.1145/3220192.3220460
- [15] T. R. Mauldin, M. E. Canby, V. Metsis, A. H. H. Ngu, and C. C. Rivera, "Smartfall: A smartwatch-based fall detection system using deep learning," *Sensors*, vol. 18, no. 10, 2018. [Online]. Available: https://www.mdpi.com/1424-8220/18/10/3363
- [16] S. S. Ogden and T. Guo, "MODI: Mobile deep inference made efficient by edge computing," in *USENIX Workshop on Hot Topics in Edge Computing (HotEdge 18)*. Boston, MA: USENIX Association, Jul. 2018. [Online]. Available: https://www.usenix.org/conference/hotedge18/presentation/ogden
- [17] T. Guo, "Cloud-based or on-device: An empirical study of mobile deep inference," in 2018 IEEE International Conference on Cloud Engineering (IC2E). IEEE, 2018, pp. 184–190.

- [18] S. Cohen, T. Kamarck, R. Mermelstein et al., "Perceived stress scale," Measuring stress: A guide for health and social scientists, vol. 10, no. 2,
- pp. 1–2, 1994. [19] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "Smote: synthetic minority over-sampling technique," *Journal of artificial intelligence research*, vol. 16, pp. 321–357, 2002.
 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg *et al.*,
- "Scikit-learn: Machine learning in python," *Journal of machine learning research*, vol. 12, no. Oct, pp. 2825–2830, 2011.