Vision-Based Recognition of Human Motion Intent During Staircase Approaching

Md R. Islam, Md R. Haque, Masudul H. Imtiaz, Xiangrong Shen, and Edward Sazonov, Senior Member, IEEE

Abstract— Walking in real-world environments involves constant decision-making, e.g., when approaching a staircase, an individual decides whether to engage (climbing the stairs) or avoid. For the control of assistive robots (e.g., robotic lower-limb prostheses), recognizing such motion intent is an important but challenging task, primarily due to the lack of available information. This paper presents a novel vision-based method to recognize an individual's motion intent when approaching a staircase before the potential transition of motion mode (walking to stair climbing) occurs. Leveraging the egocentric images from a head-mounted camera, the authors trained a YOLOv5 object detection model to detect staircases. Subsequently, an AdaBoost and gradient boost (GB) classifier was developed to recognize the individual's intention of engaging or avoiding the upcoming stairway. This novel method has been demonstrated to provide reliable (97.69%) recognition at least 2 steps before the potential mode transition, which is expected to provide ample time for the controller mode transition in an assistive robot in real-world use.

Keywords—staircase detection, intent recognition, YOLOv5, neural network

I. INTRODUCTION

ITH the rapid aging of the population, mobility impairment is becoming an increasingly challenging health problem in the United States [1]. People may suffer from the impaired ability of ambulation in daily living due to a range of reasons, including limb loss [2], age-related muscle strength decline [3], and neuromuscular pathologies (e.g., stroke) [4]. Motivated by this challenging problem, a variety of wearable robots have been developed to restore the lost lower-limb functions (for amputees) (e.g., [5]) and provide motion assistance to supplement the users' lower-limb joint efforts [6].

As a wearable robot is directly coupled with the user' limbs and joints, providing coordinated motion or motion assistance based on the user's motion intent is extremely important. However, recognizing the user's motion intent in complex realworld environments is difficult. The majority of existing methods rely on mechanical sensor signals (joint angle/velocity,

This work was supported by the National Science Foundation under award 1734501. (Corresponding author: Md Rafi Islam.)

This paper has supplementary downloadable material available at https://alabama.box.com/s/i6oc6tb3pc5k340c4hvpq36qvp2vs1ui provided by the authors. This includes one multimedia mp4 format movie clip, which shows stair climb and avoid experiment sessions. This material is 35 MB in size.

Md R. Islam is with the Department of Electrical and Computer Engineering, The University of Alabama, Tuscaloosa, AL 35487 (e-mail: mislam24@crimson.ua.edu).

foot pressure, etc.) [7] or muscle activation signals (measured through electromyography) [8] to deduce the intended motion. Such deductive methods tend to suffer from multiple significant issues, such as low accuracy and long delay, as their inputs have been limited to the physical and/or physiological signals extracted from the user himself/herself. Without access to the information on the environment, these intent recognition methods may only react to the user's actions (which, in turn, are reactions to the upcoming environmental features such as staircases) and thus unable to predict the user's intended motion to obtain smooth mode transitions in locomotion.

Motivated by this problem, multiple researchers investigated the use of vision-based environment sensing for wearable robot control. Laschowski et al. developed the ExoNet, an open-source database of high-resolution images of human walking environments [9]. Using such imagery information, environment recognition systems have been developed, which may serve the purpose of wearable robot control (e.g., [10]). On the other hand, how to use environmental information for motion intent recognition still remains an open question. As a typical example, when an individual approaches a staircase, s/he may still choose to avoid the staircase (i.e., not to go upstairs). Recognizing such motion intent in an accurate and timely manner is critical for the control of wearable robots.

In this paper, the authors present a novel vision-based system for human motion intent recognition, using the staircase approaching as a typical scenario. This specific use scenario was chosen due to the ubiquitous presence of stairs in real-world environments as well as the difficulty of mobility-challenged individuals in stair climbing. Note that the assistance provided by wearable robots constitutes a promising solution to overcome such difficulty [11], and the effectiveness of wearable robot assistance can be quantified with instrumented testbeds (e.g.,[12]). The proposed vision-based intent recognition system consists of two primary modules, including a staircase detection with the You Only Look Once (YOLO) v5 model [13] and a

Masudul H. Imtiaz was with the Department of Electrical and Computer Engineering, The University of Alabama, Tuscaloosa, AL 35487 USA, and now with the Department of Electrical and Computer Engineering, Clarkson University, Potsdam, NY 13699 USA (e-mail: mimtiaz@clarkson.edu).

Xiangrong Shen is with the Department of Mechanical Engineering, The University of Alabama, Tuscaloosa, AL 35487 USA (e-mail: xshen@eng.ua.edu).

Edward Sazonov is with the Department of Electrical and Computer Engineering, The University of Alabama, Tuscaloosa, AL 35487 (e-mail: esazonov@eng.ua.edu).

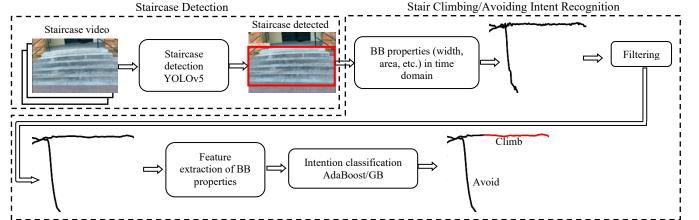


Fig. 1: Overview of the vision-based intent recognizer. The "Staircase Detection" block uses the YOLOv5 model to detect staircases from egocentric videos and generates bounding box-related information. The "Stair Climbing/Avoiding Intent Recognition" block extracts features from the detected staircases' bounding boxes and classifies user intention of either climbing or avoiding the staircase using an AdaBoost/Gradient Boost (GB) algorithm.

bounding box-based intent classification algorithm constructed with the AdaBoost and gradient boost (GB) methods. Fig. 1 illustrates the overview of the proposed method.

As the primary contribution, the vision-based intent recognizer in this paper, to the best of our knowledge, is the first work that explicitly identifies the user's intention when interacting with the environment and its key elements (such as a staircase). In existing works (such as [14]), after the vision system detects an important environmental element (stairs as a typical example), the user is expected to engage the element by default (i.e., ascending or descending stairs). However, in real life, people may still choose to avoid the element (i.e., not to ascend or descend the stairs), and the potential misclassification of the user intent may significantly affect the assistive devices' control performance and increase the risk of fall. To address this problem, the intent recognizer in this paper explicitly deduces the user's motion intent when approaching a significant environmental element (such as a stairway) using egocentric images obtained from the vision system. The egocentric images, compared with images from chest or waist-mounted cameras, better represent the user's focus of attention and thus serve as a better indicator of his/her motion intent. The proposed intent recognizer extracts bounding box features from the egocentric images and classifies the user intent using an AdaBoost /Gradient Boost classifier. The method is simple to implement and can be easily adapted to other environmental elements (e.g., doors and chairs). Further, as the proposed intent recognizer uses the images from the vision system as the sole input, it can work in conjunction with all types of assistive devices (prostheses, exoskeletons, etc.) and motion controllers (impedance control, torque control, etc.). As such, it may become an important building block of wearable robot control systems to improve the robots' performance and functionality in real-world use.

For the development of this novel method of intent recognition, the research generated a number of technical contributions, including: 1) the establishment of the vision-base motion intent recognition framework comprising an environment-sensing module and an intent recognition module; 2) developing the YOLOv5 staircase detection model; 3) identifying the most significant features of the bounding box that can help in detecting not only the intention of staircase engaging/avoiding but also other related intentions/decisions such as how to avoid crossing pedestrians; 4) lastly, a robust and highly accurate simple classification model for stair climb or avoid intention detection without heavy computational load.

II. STAIRCASE DETECTION

Identifying a staircase to prepare for the subsequent climbing motion is one of the fundamental challenges for the control of mobile robots and unmanned ground vehicles [15]. Lower-limb prostheses and assistive devices use different actuator control algorithms in different walking conditions, such as walking, running, stair climbing, and so on, to provide the most convenience to the user [16], [17]. So, to choose an appropriate control algorithm, it is crucial to detect the object, such as the staircase, and the user's intention to climb or avoid that staircase before the activity occurs.

However, predicting the wrong intention can result in an improper control algorithm and wrong actuator control, compromising the user's safety. So, improving the accuracy of user intention recognition is a very active research field nowadays [18]–[20]. Also, combining environment information with the prosthesis or assistive device user's walking biomechanics or muscle activation improves intention classification accuracy and allows for more robust control systems [16], [21]-[22]. Thus, when a prosthesis or assistive device user approaches a staircase, for robust decision-making on user intention, the prosthesis or assistive device first has to identify the staircase and then classify the user's intention of either climbing or continuing usual ground walking.

Many existing approaches extract stair edges from RGB images to detect staircase [23]–[25]. Some use depth images to include range information and make the staircase detector model more robust [26], [27]. The problem with line extraction-based methods lies in the detected straight lines' reliability and the hypothesis that staircases are always parallel lines. Also, in sunlight, most depth cameras, like Kinect, do not work till the stair is just a few feet apart from the camera [28]. Considering these drawbacks, many researchers have applied deep learning in staircase detection, and these methods have demonstrated better performance [26], [29].

However, deep learning-based methods require a vast amount of data to train the staircase detector model so the model can be generalized properly. Still, to the best of the authors' knowledge, no annotated dataset contains more than 10,000 staircase RGB images. Also, it is necessary to test the intent recognition model on a dataset that simulates a user approaching a staircase and then climbing or avoiding it. Again, to the best of the authors' knowledge, a dataset like that does not exist either. These limitations inspired us to create two datasets. The first consists of 12187 still staircase images (640×480 resolution) annotated in YOLOv5 format taken from an egocentric perspective. As most of existing stairs-related images are taken at lower heights (e.g., using chest-mounted cameras), our new dataset complements existing datasets and may improve the performance of visions systems using head-mounted cameras (e.g., cameras embedded in eyeglasses such as Google Glass). The second dataset consists of 375 egocentric videos taken with a head-mounted camera when the user approaches a staircase. Videos in this comprehensive staircase-approaching dataset cover a variety of approaching directions (approaching a staircase from the left/right side or in approximately the same direction as the staircase) and motion intents (ascending/descending stairs or avoiding the staircase). Different from still images of stairs, this dataset of videos exemplifies the visual perception of a human when approaching a staircase and making the related decision (engage/avoid) in real-world ambulation scenarios. As such, it may be very useful to computer vision and wearable robot researchers in the development and testing of vision-based intent recognizers for wearable robots. Further, it may be used as a typical example of humans' visual perception of the environment when interacting with an important environmental element, and thus serve as a basis for future expanded datasets that cover a wider variety of real-life scenarios when interacting with other types of important environmental elements (such as doors and chairs).

A. Training YOLOv5 for Staircase Detection

Considering the intent recognizer's target application of realtime control of wearable robots, we selected YOLO as the method for detecting stairs, leveraging its very fast speed of processing [30]. We specifically used YOLOv5 in the development of the staircase detector to obtain high accuracy of detection, which is highly important for the control of wearable robots and assistive devices [13].

For training the YOLOv5 model, we annotated the 12,187 staircase still images using the 'imageLabeler' application in MATLAB and converted the annotations to the YOLOv5 image annotation format.

We augmented the annotated images to increase variability and the number of images. By default, YOLOv5 applies different augmentations on the training images, such as changing hue, saturation, value, image rotation, translation, shear, and mosaic [31]. The rotation and sheer augmentation values were $\pm 15^{\circ}$, and default values were used for the rest of the augmentations [31].

In real world scenario, when we approach a staircase with a camera, sometimes the camera goes out of focus and makes the video blur. Also, image brightness and contrast differ considerably. Considering these situations, apart from the YOLOv5 augmentations, we applied blur, exposure, and noise augmentations. Out of 12,187 annotated staircase images, we used 80% or 9750 images for training and the rest for validation. After applying the three mentioned augmentations, the total number of training images became 29,250.

We trained the YOLOv5s, YOLOv5l, and YOLOv5x models with pretrained weights for 500 epochs with a batch size of 8, SGD optimizer, and patience of 20 epochs. The models were pretrained on the MS COCO dataset [32].

B. Testing YOLOv5 Model for Staircase Detection

To detect user intention for climbing or avoiding the staircase, we needed videos that show approaching the staircase and climbing or avoiding the staircase. So, based on how a person will approach a staircase, we divided the videos into three main categories:

- 1) Heading straight to the stair
- 2) Heading from the left side of the stair
- 3) Heading from the right side of the stair

These main categories were divided into five subcategories: one instance was climbing, and the rest were avoiding the staircase. So, 15 videos for each staircase. Finally, 375 videos were taken from 25 staircases inside the University of Alabama to create the testing dataset. All the videos were collected at 30 FPS using a head-mounted camera (Campark X25) oriented in portrait mode. Fig. 2 illustrates the head-mounted camera setup for capturing egocentric videos.

Then we tested the previously trained YOLOv5 models on these staircase videos and stored the bounding box (BB) properties of the detected staircase. As YOLOv5x displayed the best detection capability, we used that model's predicted BB properties for classifying stair climbing or avoiding intention detection.

III. HUMAN INTENT RECOGNITION

For intention detection, a common practice is to use electromyography (EMG) collected from the residual limb and inertial measurement units (IMUs) [33]-[36]. However, these methods are highly user-dependent, and the signals are typically delayed, resulting in delayed prediction [34]. On the other hand, vision-based intention detection methods are mostly userindependent [37]. However, all these methods typically depend on online adaptation and dataset expansion for accurate intention detection [20]. Also, they are trained only on a subset of possible sensor data related to all possible prosthesis configurations [20]. Thus, even with the large datasets and complex deep learning-based classifiers, it is not guaranteed that the existing models would detect intentions in all possible realworld situations. These shortcomings motivated us to develop a simple linear classifier that does not depend on the user; instead, the classifier uses the properties of the bounding box of the detected object and classifies the user's intention of either engaging or avoiding it. We hypothesize that when a person approaches an object in the environment (in our case, staircases), the properties, i.e., width, centroid coordinates, and area of the bounding box increase to some extent. Suppose the person engages that object (climbing the staircase). In that case, the

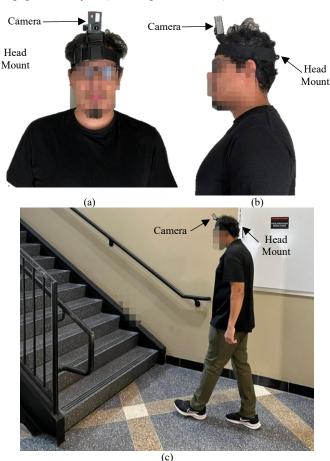


Fig. 2. Camera setup for egocentric video capture: (a) front view, (b) side view, (c) 2-point perspective view.

bounding box properties go close to maximum, and if that person avoids the staircase (continues ground walking), the property of the bounding box goes to minimum and vanishes at some point. Thus, by extracting features from those bounding box properties in the time domain, we should be able to detect the user intention using a linear classifier.

A. Bounding Box Property Extraction

When the YOLOv5 model detects the staircase in the test videos, it provides four features of the bounding box surrounding the staircase. The four features are, normalized bounding box centroid coordinates along the horizontal and vertical axis and normalized width and height of the bounding box. By observing the bounding boxes in the detected videos, we realized the area of the bounding box could be another essential property to detect the user's intention. So, from the normalized height and width of the bounding box, we also calculated the area of the bounding box.

In some frames of a staircase video, the YOLO model detected more than one bounding box. In cases like that, we took the highest height, width, and average centroid value. For the area, we calculated the total area using the following equation:

$$A_{RI} = W_{RI} * H_{RI} \tag{2}$$

$$Total\ area = A_{RI} + A_{R2} - A_I \tag{3}$$

Here, A_{RI} and A_{R2} represent the area of the 1st and 2nd rectangles, W_{RI} is the width of the rectangles, H_{RI} is the height of the 1st rectangle, and A_I is the intercepting/overlapping area between rectangles.

$$W_{RI} = abs (l1.x - r1.x)$$
 (4)

$$H_{RI} = abs (l1.y - r1.y)$$
 (5)

Here, ll.x is the x-axis coordinate on the left side, and rl.x is the x-axis coordinate on the right side. ll.y is the y-axis coordinate on the left side, and rl.y is the y-axis coordinate on the right side.

Similarly, the 2nd rectangle area was calculated. Next, we calculate the intercepting/overlapping area between rectangles using the following equations,

$$W_I = min(r1.x, r2.x) - max(l1.x, l2.x)$$
 (6)

$$H_I = min(r1.y, r2.y) - max(l1.y, l2.y)$$
 (7)

$$A_I = W_I * H_I \tag{8}$$

If W_I or H_I is negative, then the two rectangles do not intersect. In that case, the A_I is 0. With these data, we calculated the total area using (3).

Due to the frame-to-frame fluctuation of detected the bounding box, we observed high-frequency noise in the bounding box property values when plotted against frames. We applied a fifth-order moving average filter to remove the high-frequency noise without compromising the original bounding box properties (Fig. 3).

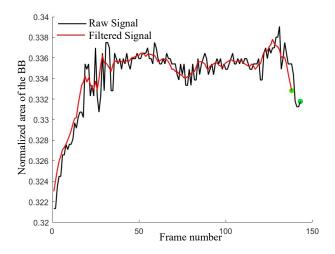


Fig. 3. Raw and filtered normalized area of the bounding box Vs. frame number

B. Feature Extraction from bounding box Properties

We ran a sliding window on the bounding box property plots to get the feature values. In this experiment, we used five different widths of sliding windows and compared them to get the best result. The widths of the windows were 5, 10, 15, 20, and 30 frames. We applied overlapping windows with a stride value of 2 during feature extraction.

We extracted simple and computationally less intensive time domain features from the bounding box properties in this paper, i.e., mean and the difference between the maximum and minimum peak of the plots in each sliding window. We extracted the features and plotted box plots to analyze data variability for the two classes. Fig. 4 illustrates the box plot of the features for 15 frame sliding window.

From Fig. 4, we observe that the most variability between classes was in the mean area, width, and centroid Y coordinate feature. Thus, these features were used to train and test the classifiers.

C. Intention Recognition Classifiers

Our goal was to prove that computationally inexpensive models can robustly classify user intention from bounding box properties which is user independent. Also, from Fig. 4, we observe that the stair climbing and avoid data have considerable differences that simple classification models can address. So, we trained and compared the results of AdaBoost and gradient-boosted (GB) tree classifiers as they are renowned for their faster training speed, low memory usage, and higher efficiency [38]-[39].

D. Dataset for Intention Recognition Classifiers

We annotated the first 35% of each bounding box property data as 'walk' and the last 35% of the signals as either 'walk' or 'climb.' We used the leave-one-out approach for training the classifier. So all the bounding box property data from one

staircase was separated for testing, and the rest of the bounding box property data from the other 24 staircases were used for training the classifiers. After training, the classifiers were tested on the left-out staircase data. This process was repeated 25 times to get results for all 25 test videos.

E. Training Parameters of the Classifiers

1) AdaBoost Classifier:

The AdaBoost algorithm uses poor learners and adaptively adjusts them by maintaining a collection of weights during the training [38]. We performed a grid search for the optimum parameters of the AdaBoost algorithm. In the grid search, the different number of learners were 5, 10, 20, 30, 50, 100, 200, learning rates were 0.025, 0.05, 0.1, 0.2, 0.3, and the maximum number of splits were 2, 5, 10, 20. The optimum hyperparameters for the AdaBoost classifier were this: number of learners 50, learning rate 0.05, and maximum number of split 5.

2) Gradient Boosting Classifier (GB):

In the GB algorithm, the decision procedure combines the outcome of many weak models to provide a more accurate estimation of the response variable. The principle of this algorithm is to update the new base models in a way that correlates with the negative gradient of the loss function, which represents the whole ensemble [39]. Again, we performed a grid search to find the optimum hyperparameters. The sub-sampling factors tested in the grid search were 0.1, 0.15, 0.5, 0.75, 1; different learning rates were 0.025, 0.05, 0.1, 0.2, 0.3, and different maximum tree depths were 2, 3, 5, 7, 10. The best result was found for the sub-sampling factor of 0.15, the learning rate of 0.25, and the maximum tree depth of 2.

F. Majority Voting for Final Decision

Till now, the proposed approach classified the user intention only based on the features of the bounding box property signal inside the sliding window. Due to noise or other factors like, incorrect object detection, the bounding box features can change rapidly and as a result the classifier can repeatedly change its decision. However, to supply a robust control signal to the prosthesis controller, the classifier should not change its decision too often. So, to improve the robustness, we introduced majority voting, i.e., the consecutive classifier decisions voted either climb or walk to decide the final intention. In this study, we experimented with four sets of number of classifier predictions to decide the final intention class. The number of votes (predictions) in those four sets are 5, 10, 15, or 30. The class that got the most votes were chosen as final class. If the classes received equal number of votes, then the last class that was selected from majority voting was chosen as final intention.

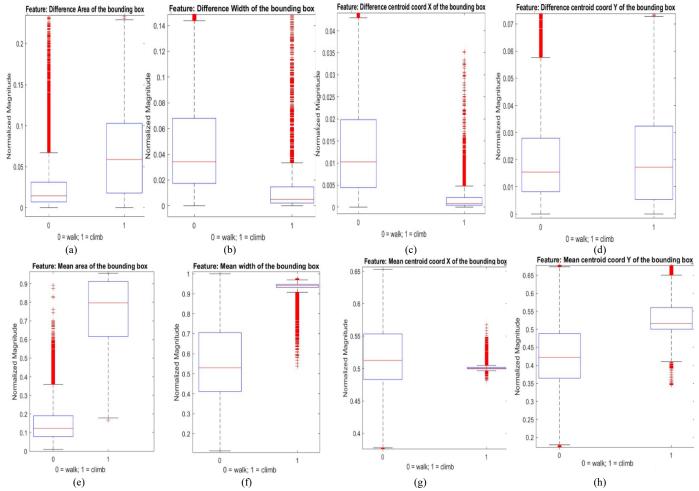


Fig. 4. Box plot of features: (a) Max-min peak difference of the bounding box area; (b) Max-min peak difference of the bounding box width; (c) Max-min peak difference of the bounding box X coordinate; (d) Max-min peak difference of the bounding box Y coordinate; (e) Mean of the bounding box area; (f) Mean of the bounding box width; (g) Mean of the bounding box X coordinate; (h) Mean of the bounding box Y coordinate.

IV. RESULT

First, we compared different YOLOv5 model performances to determine the best model and used the model for further analysis. Next, we determined three performance metrics representing real-world requirements from an intention detector, like detection time and robustness. Then we compared the classifier's performance for different lengths of feature windows. Finally, we compared different voting lengths to get the best possible outcome from the classifiers.

A. Performance evaluation of YOLOv5 Models

The YOLOv5x is the largest model and had the best performance on the validation set. The following best was YOLOv5l, and the last was YOLOv5s, the smallest model (Table I). Thus, we applied the YOLOv5x on the test dataset of

375 videos of approaching, climbing, or avoiding a staircase and stored the bounding box properties.

B. Metrics for Quantifying Intention Classifier Performance

We addressed three aspects of intention detection important for the control of an assistive or prosthetic device. First, the intention must be detected sometime before the actual intended action occurs so that the operating device has sufficient time to change its mode of operation according to the intended action.

TABLE I.
PERFORMANCE OF YOLOV5 MODELS

Parameters	YOLOv5s	YOLOv5l	YOLOv5x
mAP_0.5 (%)	65.3	78.9	83.4
mAP_0.5:0.95 (%)	28.6	48.2	51.8
Precision	0.71	0.82	0.86
Recall	0.65	0.73	0.74

Here, mAP is mean average precision.

Second, the intention detection must be precise, and misclassification should be as little as possible. Lastly, the classifier should not change its class too often after classifying an intended action. If the classifier changes the detected class multiple times, the prosthesis or assistive device would also switch the mode of operation, and the user would feel discomfort and fall in a worst-case scenario. These three metrics are directly related to user safety which is one of the prime concerns in a prosthesis design. That is why we addressed misclassification, mean prediction time before stepping on the staircase, and percentage of change of class while approaching a staircase as performance metrics of our classifier.

To calculate the time between intention classification and action, first, we took the time of classification. Second, we recorded the time of the first step on the stairs. We defined the time between these two steps as the prediction time in advance of the action. Fig. 5 illustrates those two steps.

Fig. 5(a) illustrates that the intention was detected on frame 241, i.e., 8.03 seconds (30 FPS), and Fig. 5(b) illustrates the step on the staircase is at 289th frame, i.e., 9.64 seconds. So, the intention was classified about 1.61 seconds before the intended action occurred.

For the misclassification calculation, if the final prediction by the classifier does not match the actual class, the outcome was classified as misclassification. As we tested using the leave-oneout approach, we added all the misclassifications, divided the sumby the total number of test cases, and got the mean misclassification by the classifiers.

Finally, we calculated the percentage of the classifier's change of decisions after predicting the class for the first time. To calculate it, we added the total number of times the decision was changed and then took the percentage to obtain the final percentage of change of class for climbing or avoiding a staircase.

Table II shows the comparison of classifiers' performance for different lengths of feature windows. The GB model performs best in all three-performance metrics with only 0.45% misclassification. Although AdaBoost has a comparable outcome with GB on misclassification rate and percentage of change of class in each prediction, GB outperforms AdaBoostin all cases.

Next, the feature collected with 15 frame window performs better in misclassification rate mean percentage of change of class performance metrics. The only metric the 15 frame signal window lags a bit in performance is the mean prediction time before stepping on the staircase, where features with 30 frames in each window perform better with the GB classifier.

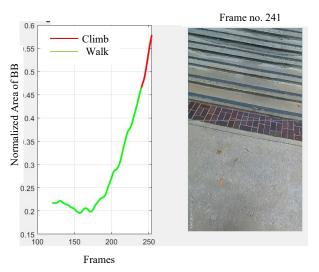
So, after analyzing the results from Table II, we can say that the GB is the best classifier in intention detection in this study, and the feature window with 15 frames displays the best possible outcome.

As described above, we applied majority voting to improve the performance even further. Table III contains the outcome of majority voting with the different number of votes used for deciding the class. The table shows that although the mean

prediction time before stepping on the staircase and percentage of change of class per staircase video metrics is improved with increased votes, the misclassification rate increases slightly. Plotting the GB classifier's output in one figure on all 25 test staircase videos will obscure it. So, we illustrated the GB classifier's outcome on one test video in Fig. 6, where the feature window had 15 frames, and 15 votes were used to decide the class.

V. DISCUSSION

Our objective was to develop a robust intent recognization method that can be used to provide reliable detection of the user's motion intent for control of wearable robots/assistive devices. Intention recognition from environmental features still requires computationally expensive deep learning models [20], [37], [40] We hypothesized that if an object, toward which a.



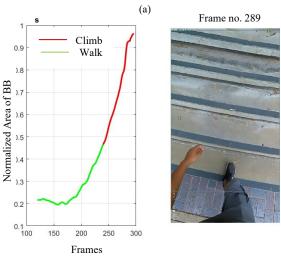


Fig. 5. Illustration of calculating intention classification time before stair climb. (a) Frame at which intention was detected. (b) Frame at which first stepped on the staircase.

TABLE II.

PERFORMANCE COMPARISON OF THE CLASSIFIER FOR DIFFERENT FEATURE WINDOWS

#of frames in feature windows	Percentage of Misclassification (%)		Mean prediction time before stepping on the staircase (s)		Percentage of change of class per staircase video	
	AdaBoost	Gradient Boost	AdaBoost	Gradient Boost	AdaBoost	Gradient Boost
5	2.59	1.32	0.68	1.24	23.24	17.38
10	1.86	0.56	0.95	1.51	18.08	11.26
15	1.48	0.45	1.1	1.73	15.29	8.03
20	1.32	0.47	1.33	1.77	14.83	13.04
30	1.45	0.70	1.52	1.94	14.56	13.07

person is moving, can be detected properly with a bounding box, the magnitude of the properties of that bounding box will increase when the person approaches closer to that object as the object would become more prominent to the detector. Then the features of the bounding box properties can be used to robustly detect the user's intention

Vision-based object sensing does not suffer from drawbacks like low depth sensor range. YOLOv5x is the largest among the YOLOv5 models, contributing to its superior performance on the validation and test datasets (Table I) [41]. However, the YOLOv5l and YOLOv5s also show promising results, and they can be used to detect staircases on devices with less computational power.

The staircase engaging/avoiding classifier was developed on the features of bounding box properties. It can be argued that the objectdetector might be unable to detect the staircase for some frames. Also, there can be possible misdetection of objects. Noises like this can be mitigated through the low-pass filter used before the feature extraction step. Also, this method of intention detection does not depend on the user gait cycle of other physiological data; instead, the proposed approach uses staircase video collected from a camera. So, we can argue that our approach should be able to detect any user's intention from any video that shows approaching a staircase. Thus, this makes the proposed approach user-independent.

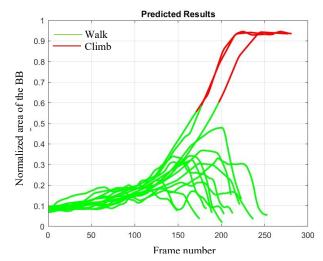


Fig. 6. GB classifiers output on a staircase video

The GB outperforms the AdaBoost in all the performance evaluation metrics. The use of loss functions instead of penalizing misclassification and the sensitivity to the outliers of the AdaBoost algorithm might have contributed to the better performance of the GB algorithm. With 15 decision majority voting, the GB model can predict staircase climbing intention about 1.01 seconds before the staircase climbing on average. Its misclassification rate is only 1.37%, which shows that the proposed method not only predicts outcomes about 1 second in advance but also with high accuracy. Also, its tendency to change its prediction is only 2.31% which represents the robustness of the model and proves that this model fulfills the essential requirements in an intention predictor, i.e., robust, accurate, and predicts well in advance of the action. In the future, we will research creating an ensemble model to improve the classifier performance even more. If the application of the proposed method requires more accuracy, the user can easily switch to the no-voting option and get accuracy above 99%.

On average, the mean prediction time before stepping on the staircase is higher for 30 frames feature window. That is because the feature window with 30 frames has more data and more significant features to distinguish between climbing and avoiding intention. Thus, the classifier can distinguish climbing intention early with 30 frame window, resulting in more time to decide the intention to stair climb before the actual climbing happens.

On the other hand, the mean prediction time before stepping on the staircase decreases with an increased number of voting in Table III. This is because an increased number of voting means the method has to wait more time to get a majority vote to decide a class and come to a conclusion. So, more time is lost in getting the decision from the voting and thus less time before

 $\label{thm:condition} TABLE~III.$ Results after Applying Voting on GB Classifier

#of decisions used for voting	Percentage of Misclassification (%)	Mean prediction time before stepping on the staircase (s)	Percentage of change of class per staircase video
No voting	0.45	1.731	8.03
5	0.91	1.383	6.88
10	0.91	1.269	4.82
15	1.37	1.012	2.31
30	1.82	0.662	0.92

getting the final decision of intention detection. This phenomenon can also explain the decreasing accuracy with increasing voting. In the proposed approach, the accuracy is calculated based on the last frame, whether it is a climb or avoid. Here, misclassified classes play a more significant role when we use 30 votes to decide a class. So, suppose at the end part of a staircase video, the classifier misclassified a class. In that case, there is simply not enough data point left for the classifier to correct the classification as 30 vote majority voting requires at least 16 votes for deciding a class.

Apart from this, our data suggests that we can make a robust predictor model with simple classifiers and minimum features of bounding box properties. To the best of our knowledge, this is the first time bounding box properties were used to classify stair climbing. or avoid intention. In the future, we will research how we can accurately classify intention with even fewer properties of the bounding box while increasing the robustness.

For the future application of the proposed intent recognizer, it can be integrated into a typical hierarchical control system of a wearable robot to identify the user's motion intent, which will enable the robot's lower-level controller to calculate the specific control commands for its powered actuators (e.g., desired assistive torque or joint angle trajectory). As the intent recognizer does not rely on external sensor signals to function, and thus can be used in conjunction with all types of wearable robots such as robotic lower-limb prosthesis and assistive ankle/knee orthoses/exoskeletons. Currently the intent recognizer only identifies the user's motion intent in engaging or avoiding stairs, but the algorithm can be adapted to the recognition of the user's intent when countering other important environmental elements such as doors and chairs. Another limitation is its hardware requirement. YOLO, as the most popular object detection algorithm, provides multiple advantages such as fast processing and high accuracy of recognition. On the other hand, YOLO also requires highperformance hardware, which limits its use in real-time embedded systems. Multiple YOLO models have been developed in recent years to enable the implementation in embedded system (e.g., the Fast YOLO [42] and the Efficient YOLO [43]). YOLOv5 has been successfully deployed on standalone devices such as the Raspberry Pi 4 [44], and the integration of a Coral USB accelerator further enhances the detection speed on the Raspberry Pi 4 [45]. In addition, YOLOv5 can also be executed in compact portable devices like the Jetson Nano, equipped with a high-resolution camera such as the Waveshare IMX477 CSI [46]. Based on such recent trend, we envision that, with the fast technological advances in high-performance embedded microprocessor and computer vision, the proposed intent recognizer can be implemented in wearable robot control systems in the near future and significantly improve the robots' performance in navigating real-world environments.

For the future work, we plan to incorporate the proposed intent recognizer into complete control systems for wearable robots and characterize its performance in human testing. Note that, as most assistive device users are expected to have normal vision (or nearly normal vision after correction), the way they perceive the environment during locomotion should remain largely unchanged, and thus we expect the proposed intent recognizer to be appliable to new users without needing much additional training.

VI. CONCLUSION

We have created a database of annotated staircase images and trained three YOLOv5 models to detect staircases. We also created a database of 375 videos of 25 staircases, simulating a person walking toward a staircase and avoiding or climbing it. We validated that the properties of the bounding box from a detected staircase can be used to create a simple but robust classifier to classify stair climbing or avoiding intentions. Our data shows that the intention detection classifier (GB) is highly accurate (97.69%), and on average, the classifier can detect intention about 1 second before the stair climb. These results advocate that the proposed method can be utilized to send a robust control signal for the prosthesis or assistive device. Our approach also discovered the essential features of bounding box properties for intention detection, leading to future research on intention detection with different types of objects. We also believe our created datasets would help create and test new staircase detection and intention classification models in the future.

VII. REFERENCES

- J. Gill and M. J. Moore, "The State of Aging and Health in America 2013," 2013, Accessed: Nov. 03, 2022. [Online]. Available: www.cdc.gov/aging.
- [2] K. Ziegler-Graham, E. J. MacKenzie, P. L. Ephraim, T. G. Travison, and R. Brookmeyer, "Estimating the prevalence of limb loss in the United States: 2005 to 2050," *Arch Phys Med Rehabil*, vol. 89, no. 3, pp. 422–429, Mar. 2008, doi: 10.1016/J.APMR.2007.11.005.
- [3] J. L. Hodgson and E. R. Buskirk, "Physical fitness and age, with emphasis on cardiovascular function in the elderly," *J Am Geriatr Soc*, vol. 25, no. 9, pp. 385–392, 1977, doi: 10.1111/J.1532-5415.1977.TB00671.X.
- [4] F. M. Collen, D. T. Wade, and C. M. Bradshaw, "Mobility after stroke: reliability of measures of impairment and disability," *Int Disabil Stud*, vol. 12, no. 1, pp. 6–9, 1990, doi: 10.3109/03790799009166594.
- [5] B. E. Lawson, J. Mitchell, D. Truex, A. Shultz, E. Ledoux, and M. Goldfarb, "A robotic leg prosthesis: Design, control, and implementation," *IEEE Robot Autom Mag*, vol. 21, no. 4, pp. 70–81, Dec. 2014, doi: 10.1109/MRA.2014.2360303.
- [6] H. Lee, P. W. Ferguson, and J. Rosen, "Lower limb exoskeleton systems-overview," Wearable Robotics: Systems and Applications, pp. 207–229, Jan. 2019, doi: 10.1016/B978-0-12-814659-0.00011-4.
- [7] H. A. Varol, F. Sup, and M. Goldfarb, "Multiclass Real-Time Intent Recognition of a Powered Lower Limb Prosthesis," *IEEE Trans Biomed Eng*, vol. 57, no. 3, p. 542, Mar. 2010, doi: 10.1109/TBME.2009.2034734.
- [8] H. Huang, T. A. Kuiken, and R. D. Lipschutz, "A strategy for identifying locomotion modes using surface electromyography," *IEEE Trans Biomed Eng*, vol. 56, no. 1, pp. 65–73, Jan. 2009, doi: 10.1109/TBME.2008.2003293.
- [9] B. Laschowski, W. McNally, A. Wong, and J. McPhee, "ExoNet Database: Wearable Camera Images of Human Locomotion Environments," Front Robot AI, vol. 7, p. 188, Dec. 2020, doi: 10.3389/FROBT.2020.562061/BIBTEX.

- [10] B. Laschowski, W. McNally, A. Wong, and J. McPhee, "Preliminary Design of an Environment Recognition System for Controlling Robotic Lower-Limb Prostheses and Exoskeletons," *IEEE Int Conf Rehabil Robot*, vol. 2019, pp. 868–873, Jun. 2019, doi: 10.1109/ICORR.2019.8779540.
- [11] J. Jang, K. Kim, J. Lee, B. Lim, and Y. Shim, "Assistance strategy for stair ascent with a robotic hip exoskeleton," *IEEE International Conference on Intelligent Robots and Systems*, vol. 2016-November, pp. 5658–5663, Nov. 2016, doi: 10.1109/IROS.2016.7759832.
- [12] Maugliani, N., Caimmi, M., Malosio, M., Airoldi, F., Borro, D., Rosquete, D., Sergio, A., Giusino, D., Fraboni, F., Ranieri, G. and Pietrantoni, L., 2022. Lower-limbs exoskeletons benchmark exploiting a stairs-based testbed: the STEPbySTEP Project. In Wearable Robotics: Challenges and Trends: Proceedings of the 5th International Symposium on Wearable Robotics, WeRob2020, and of WearRAcon Europe 2020, October 13–16, 2020 (pp. 603-608). Springer International Publishing.
- [13] "GitHub ultralytics/yolov5: YOLOv5 in PyTorch > ONNX > CoreML > TFLite." https://github.com/ultralytics/yolov5 (accessed Aug. 31, 2022).
- [14] Tricomi, E., Mossini, M., Missiroli, F., Lotti, N., Xiloyannis, M., Roveda, L. and Masia, L., 2022. Environment-based Assistance Modulation for a Hip Exosuit via Computer Vision. arXiv preprint arXiv:2211.15346.
- [15] K. Lee, V. Kalyanram, C. Zhengl, S. Sane, and K. Lee, "Vision-based Ascending Staircase Detection with Interpretable Classification Model for Stair Climbing Robots," 2022 International Conference on Robotics and Automation (ICRA), pp. 6564–6570, May 2022, doi: 10.1109/ICRA46639.2022.9812456.
- [16] M. R. Tucker et al., "Control strategies for active lower extremity prosthetics and orthotics: A review," J Neuroeng Rehabil, vol. 12, no. 1, pp. 1–30, Jan. 2015, doi: 10.1186/1743-0003-12-1/FIGURES/2.
- [17] A. J. Young and D. P. Ferris, "State of the art and future directions for lower limb robotic exoskeletons," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 25, no. 2, pp. 171–182, Feb. 2017, doi: 10.1109/TNSRE.2016.2521160.
- [18] R. Stolyarov, G. Burnett, and H. Herr, "Translational Motion Tracking of Leg Joints for Enhanced Prediction of Walking Tasks," *IEEE Trans Biomed Eng*, vol. 65, no. 4, pp. 763–769, Apr. 2018, doi: 10.1109/TBME.2017.2718528.
- [19] M. Liu, F. Zhang, and H. H. Huang, "An Adaptive Classification Strategy for Reliable Locomotion Mode Recognition," *Sensors* (Basel), vol. 17, no. 9, Sep. 2017, doi: 10.3390/S17092020.
- [20] B. Hu, A. M. Simon, and L. Hargrove, "Deep Generative Models with Data Augmentation to Learn Robust Representations of Movement Intention for Powered Leg Prostheses," *IEEE Trans Med Robot Bionics*, vol. 1, no. 4, pp. 267–278, Nov. 2019, doi: 10.1109/TMRB.2019.2952148.
- [21] M. Liu, D. Wang, and H. Helen Huang, "Development of an Environment-Aware Locomotion Mode Recognition System for Powered Lower Limb Prostheses," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 24, no. 4, pp. 434–443, Apr. 2016, doi: 10.1109/TNSRE.2015.2420539.
- [22] B. Laschowski, W. McNally, A. Wong, and J. McPhee, "Computer Vision and Deep Learning for Environment-Adaptive Control of Robotic Lower-Limb Exoskeletons," *Proceedings of the Annual International Conference of the IEEE Engineering in Medicine and Biology Society, EMBS*, vol. 2021-January, pp. 4631–4635, 2021, doi: 10.1109/EMBC46164.2021.9630064.
- [23] J. A. Hesch, G. L. Mariottini, and S. I. Roumeliotis, "Descending-stair detection, approach, and traversal with an autonomous tracked vehicle," *IEEE/RSJ 2010 International Conference on Intelligent Robots and Systems, IROS 2010 Conference Proceedings*, pp. 5525–5531, 2010, doi: 10.1109/IROS.2010.5649411.
- [24] S. Carbonara and C. Guaragnella, "Efficient stairs detection algorithm Assisted navigation for vision impaired people," INISTA 2014 - IEEE International Symposium on Innovations in Intelligent Systems and Applications, Proceedings, pp. 313–318, 2014, doi: 10.1109/INISTA.2014.6873637.
- [25] Y. Cong, X. Li, J. Liu, and Y. Tang, "A stairway detection algorithm based on vision for UGV stair climbing," *Proceedings of 2008 IEEE*

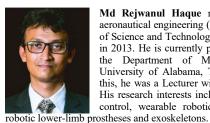
- International Conference on Networking, Sensing and Control, ICNSC, pp. 1806–1811, 2008, doi: 10.1109/ICNSC.2008.4525517.
- [26] K. Lee, V. Kalyanram, C. Zhengl, S. Sane, and K. Lee, "Vision-based Ascending Staircase Detection with Interpretable Classification Model for Stair Climbing Robots," 2022 International Conference on Robotics and Automation (ICRA), pp. 6564–6570, May 2022, doi: 10.1109/ICRA46639.2022.9812456.
- [27] S. Murakami, M. Shimakawa, K. Kivota, and T. Kato, "Study on stairs detection using RGB-depth images," 2014 Joint 7th International Conference on Soft Computing and Intelligent Systems, SCIS 2014 and 15th International Symposium on Advanced Intelligent Systems, ISIS 2014, pp. 1186–1191, Feb. 2014, doi: 10.1109/SCIS-ISIS.2014.7044705.
- [28] T. Westfechtel *et al.*, "Robust stairway-detection and localization method for mobile robots using a graph-based model and competing initializations:," *The International Journal of Robotics Research*, vol. 37, no. 12, pp. 1463–1483, Sep. 2018, doi: 10.1177/0278364918798039.
- [29] D. E. Diamantis, D. C. C. Koutsiou, and D. K. Iakovidis, "Staircase detection using a lightweight look-behind fully convolutional neural network," *Communications in Computer and Information Science*, vol. 1000, pp. 522–532, 2019, doi: 10.1007/978-3-030-20257-6 45/TABLES/3.
- [30] Redmon, J., Divvala, S., Girshick, R. and Farhadi, A., 2016. You only look once: Unified, real-time object detection. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 779-788).
- [31] G. Jocher and J. Borovec, "Yolov5/hyp.scratch-low.yaml at B94b59e199047aa8bf2cdd4401ae9f5f42b929e6 ultralytics/yolov5," GitHub. [Online]. Available: https://github.com/ultralytics/yolov5/blob/b94b59e199047aa8bf2cdd4401ae9f5f42b929e6/data/hyps/hyp.scratch-low.yaml#L6-L34. [Accessed: 17-Mar-2023].
- [32] "Training the YOLOv5 Object Detector on a Custom Dataset PyImageSearch." https://pyimagesearch.com/2022/06/20/training-the-yolov5-object-detector-on-a-custom-dataset/ (accessed Nov. 03, 2022).
- [33] M. Goršič *et al.*, "Online phase detection using wearable sensors for walking with a robotic prosthesis," *Sensors (Basel)*, vol. 14, no. 2, pp. 2776–2794, Feb. 2014, doi: 10.3390/S140202776.
- [34] H. F. Maqbool, M. A. B. Husman, M. I. Awad, A. Abouhossein, N. Iqbal, and A. A. Dehghani-Sanij, "A Real-Time Gait Event Detection for Lower Limb Prosthesis Control and Evaluation," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 25, no. 9, pp. 1500–1509, Sep. 2017, doi: 10.1109/TNSRE.2016.2636367.
- [35] T. R. Clites et al., "Proprioception from a neurally controlled lower-extremity prosthesis," Sci Transl Med, vol. 10, no. 443, May 2018, doi: 10.1126/SCITRANSLMED.AAP8373/SUPPL_FILE/AAP8373_TA BLE S1.ZIP.
- [36] A. Doulah, X. Shen, and E. Sazonov, "Early Detection of the Initiation of Sit-to-Stand Posture Transitions Using Orthosis-Mounted Sensors," *Sensors 2017, Vol. 17, Page 2712*, vol. 17, no. 12, p. 2712, Nov. 2017, doi: 10.3390/S17122712.
- [37] K. Zhang et al., "Environmental Features Recognition for Lower Limb Prostheses Toward Predictive Walking," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 27, no. 3, pp. 465–476, Mar. 2019, doi: 10.1109/TNSRE.2019.2895221.
- [38] R. Wang, "AdaBoost for Feature Selection, Classification and Its Relation with SVM, A Review," *Phys Procedia*, vol. 25, pp. 800– 807, Jan. 2012, doi: 10.1016/J.PHPRO.2012.03.160.
- [39] A. Natekin and A. Knoll, "Gradient boosting machines, a tutorial," Front Neurorobot, vol. 7, no. DEC, 2013, doi: 10.3389/FNBOT.2013.00021.
- [40] B. Zhong, R. L. da Silva, M. Li, H. Huang, and E. Lobaton, "Environmental Context Prediction for Lower Limb Prostheses with Uncertainty Quantification," *IEEE Transactions on Automation Science and Engineering*, vol. 18, no. 2, pp. 458–470, Apr. 2021, doi: 10.1109/TASE.2020.2993399.

- [41] "Train Custom Data YOLOv5 Documentation." https://github.com/ultralytics/yolov5/wiki/Train-Custom-Data (accessed Sep. 08, 2022).
- [42] Shafiee, M.J., Chywl, B., Li, F. and Wong, A., 2017. Fast YOLO: A fast you only look once system for real-time embedded object detection in video. arXiv preprint arXiv:1709.05943.
- [43] Wang, Z., Zhang, J., Zhao, Z. and Su, F., 2020, July. Efficient yolo: A lightweight model for embedded deep learning object detection. In 2020 IEEE International Conference on Multimedia & Expo Workshops (ICMEW) (pp. 1-6). IEEE
- [44] J. Johnston, "Tutorial: Running yolov5 machine learning detection on a raspberry pi 4," Medium, 08-Apr-2021. [Online]. Available: https://jordan-johnston271.medium.com/tutorial-running-yolov5machine-learning-detection-on-a-raspberry-pi-4-3938add0f719. [Accessed: 25-Feb-2023].
- [45] S. Virahonda, "Deploying YOLOv5 Model on Raspberry Pi with Coral USB Accelerator" [Online]. Available: https://www.codeproject.com/Articles/5293079/Deploying-YOLOv5-Model-on-Raspberry-Pi-with-Coral [Accessed: 24-Feb-2023]
- [46] A. Heydarian, "Yolov5 Object Detection on NVIDIA Jetson Nano" [Online]. Available: https://towardsdatascience.com/yolov5-object-detection-on-nvidia-jetson-nano-148cfa21a024 [Accessed: 24-Feb-2023]



Md. Rafi Islam received the B.Sc. degree in electrical and electronic engineering (EEE) from the Chittagong University of Engineering and Technology (CUET), Chittagong, Bangladesh, in 2017. He is currently pursuing a Ph.D. degree with the Department of Electrical and Computer Engineering, The University of Alabama, Tuscaloosa, AL, USA. Before this, he was a Lecturer with the Department of EEE, Premier University, Chittagong. His research interests include embedded systems, machine learning, wearable

robotics, and intent recognition.

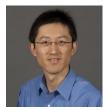


Md Rejwanul Haque received the B.Sc. degree in aeronautical engineering (AE) from the Military Institute of Science and Technology (MIST), Dhaka, Bangladesh, in 2013. He is currently pursuing the Ph.D. degree with the Department of Mechanical Engineering, The University of Alabama, Tuscaloosa, AL, USA. Before this, he was a Lecturer with the AE Department, MIST. His research interests include advanced sensing, motion control, wearable robotics, and intent recognition of



Masudul H Imtiaz is currently an assistant professor with the Department of Electrical and Computer Engineering, Clarkson University, Potsdam, NY, USA, and head of the AI Vision, Health, Biometrics, and Applied Computing (AVHBAC) lab. Dr. Imtiaz received bachelor's and master's degrees in applied physics, electronics, and communication engineering from the

University of Dhaka, Bangladesh, and a Ph.D. degree from the University of Alabama in the summer of 2019. He was also a Postdoctoral Fellow with the Department of Electrical and Computer Engineering at the University of Alabama. His research interests include the development of wearable systems, mHealth, deep learning on wearables, biomedical signal processing, and computational intelligence for preventive, diagnostic, and assistive technology, with a special focus on health monitoring and rehabilitation. Dr. Imtiaz is also developing novel biometric technologies ensuring high accuracy in enrollment and person verification.



Xiangrong Shen received his Ph.D. degree in mechanical engineering from Vanderbilt University, Nashville, TN, USA, in 2006. He subsequently completed a two-year post-doctoral training in rehabilitation robotics at Vanderbilt University. He is a Professor with the Department of Mechanical Engineering, The University of Alabama, Tuscaloosa, AL, USA, where he has been a Faculty Member since 2008. His research interests include assistive and

rehabilitation robotics, and the specific topics include robotic lower-limb prostheses, portable power-assist lower-limb orthoses, and legged/wheeled robotic platforms for assistive and therapeutic purposes. His research has been supported by multiple NSF and NIH grants, including the NSF CAREER Award in 2014. Dr. Shen served as an Associate Editor for the journal of Control Engineering Practice from 2008 to 2014. He also serves as the Chair for the ME Dynamic Systems and Control (DSC) Group and a Fellow for the Alabama Life Research Institute (ALRI).



Edward Sazonov (Senior Member, IEEE) received the Diploma of Systems Engineer from the Khabarovsk State University of Technology, Khabarovsk, Russia, in 1993, and the Ph.D. degree in computer engineering from West Virginia University, Morgantown, WV, USA, in 2002. He is currently a James R. Cudworth Endowed Professor with the Department of Electrical and Computer Engineering, The University of Alabama, Tuscaloosa, AL, USA, and the Head of the Computer Laboratory of Ambient and

Wearable Systems (http://claws.eng.ua.edu), The University of Alabama. His research interests include wearable devices, sensor-based behavioral informatics, and methods of biomedical signal processing and pattern recognition. Devices developed in his laboratory include a wearable sensor for objective detection and characterization of food intake automatic ingestion monitor (AIM); a highly accurate physical activity and gait monitor integrated into a shoe insole (SmartStep, winner of Bluetooth Innovation WorldCup 2009); a wearable sensor system for monitoring of cigarette smoking [Personal Automatic Cigarette Tracker (PACT)]; and others. Dr. Sazonov's research in his lab was recognized by several awards, including best paper awards and the President's Research Award at the University of Alabama. He served as a Fulbright Distinguished Chair for the University of Newcastle, Callaghan, NSW, Australia. His research has been supported by the National Institutes of Health, the National Science Foundation, and the National Academies of Science, as well as by state agencies, private industry, and foundations. He serves as a Specialty Chief Editor for Wearable Electronics and Frontiers in Electronics and an associate editor for several IEEE journals.