

Thesis Submitted for Master of Science in Computer Science

Department of Computer Science

Rochester Institute of Technology

Generalized Model of Cognitive Workload

Taylor Carpenter

tjc1575@rit.edu

Chair: Dr. Zack Butler

zjb@cs.rit.edu

Reader: Dr. Esa Rantanen

emrgsh@rit.edu

Observer: Sean Strout

sps@cs.rit.edu

Rochester, NY 14623 USA

July 28, 2015

A Thesis Submitted
in
Partial Fulfillment of the
Requirements for the Degree of
Master of Science in Computer Science

Department of Computer Science
B. Thomas Golisano College of Computing and Information Sciences
Rochester Institute of Technology
Rochester, NY 14623 USA

Abstract

Author remarks are made in italics. Bolded sections signify passages that I am particularly concerned about.

Miscellaneous TODO:

- Change “Study” to “Participant Study” in appropriate places.
- Change “System” to “Classification System” in appropriate places.
- Check tenses
- Add captions to results graphs

Contents

1	Introduction	5
2	Background	5
2.1	Terminology	5
2.1.1	Cognitive Workload	5
2.1.2	Participant Study	6
2.1.3	Classification System	6
2.1.4	Model / Classifier	6
2.2	Cognitive Workload	7
2.3	Psychophysiological Measures	7
2.4	Random Forest	7
2.5	Artificial Neural Network	7
3	Related Work	7
4	Problem Statement	8
4.1	Applications	9
4.1.1	Adaptive Automation	9
4.1.2	Research	9
5	Participant Study	10
5.1	Multi-Attribute Task Battery (MATB)	11
5.2	RanTask	14
5.3	NASA Task Load Index (TLX)	15
5.4	Structure	16
6	System	17
6.1	Sensors	17
6.1.1	Electroencephalogram (EEG)	17
6.1.2	Heart Rate (HR)	19
6.2	Data Preprocessing	19
6.2.1	EEG	20
6.2.2	Heart Rate	21
6.3	Data Processing and Features	22
6.3.1	EEG	23

6.3.2	Heart Rate	24
6.4	Configurations	25
6.4.1	Same Participant - Same Task (SP-ST)	25
6.4.2	All Participants - Same Task (AP-ST)	26
6.4.3	Same Participant - All Tasks (SP-AT)	26
6.4.4	All Participants - All Tasks (AP-AT)	26
6.4.5	Cross Participant - Same Task (CP-ST)	26
6.4.6	Same Participant - Cross Task (SP-CT)	27
6.4.7	All Participant - Cross Task (AP-CT)	27
6.4.8	Cross Participant - All Tasks (CP-AT)	27
6.4.9	Cross Participant - Cross Task (CP-CT)	28
6.4.10	Overview	28
6.5	Classifiers	28
6.5.1	Random Forest	29
6.5.2	Artificial Neural Network (ANN)	29
7	Results	30
7.1	Participant Study	35
7.1.1	Performance	35
7.1.2	NASA Task Load Index (TLX)	35
7.2	Classifier Models	35
7.2.1	Same Participant - Same Task (SP-ST)	35
7.2.2	All Participants - Same Task (AP-ST)	35
7.2.3	Same Participant - All Tasks (SP-AT)	35
7.2.4	All Participants - All Tasks (AP-AT)	35
7.2.5	Cross Participant - Same Task (CP-ST)	35
7.2.6	Same Participant - Cross Task (SP-CT)	35
7.2.7	All Participant - Cross Task (AP-CT)	35
7.2.8	Cross Participant - All Tasks (CP-AT)	35
7.2.9	Cross Participant - Cross Task (CP-CT)	35
8	Discussion	35
8.1	Participant Study	38
8.2	Classifier Models	38
8.3	Potential Sources of Error	38

9 Future Work & Conclusion	38
---------------------------------------	-----------

Appendices

A NASA Task Load Index Records	43
A.1 Multi-Attribute Task Battery (MATB)	43
A.2 RanTask	43
B Performance Data	43
B.1 Multi-Attribute Task Battery (MATB)	43
B.2 RanTask	43
C Model Performance	43
C.1 Artificial Neural Network	43
C.2 Random Forest	43

List of Figures

1	Multi-Attribute Task Battery (MATB) Main View	12
2	General Overview of Complete System Workflow	18
3	Overview of EEG Preprocessing Workflow	21
4	Overview of the heart rate preprocessing workflow.	22
5	Overview of the processing and feature generation workflow.	23
6	31
7	32
8	32
9	33
10	33
11	34
12	34
13	35
14	36
15	36
16	37
17	37

List of Tables

1	Multi-Attribute Task Battery (MATB) Load Condition Parameters .	13
2	RanTask Load Condition Tone Counts	14
3	Configuration Model Counts	28
4	Performance Difference Between Conditions	38
5	NASA-TLX Difference Between Conditions	38
6	Model Accuracies	38

1 Introduction

2 Background

2.1 Terminology

The following subsections layout the definitions that are used for various terminology that appear throughout this research. Some of the terms are commonly used in this area of research but have various meanings, depending on the field and background from which the research originates. Other terms are being assigned unique meanings for the purposes of this research. For the sake of clarity, the major terms are explained explicitly, ideally removing confusion on how the terms are being used.

2.1.1 Cognitive Workload

Many definitions exist for cognitive workload, as it is a theoretical concept that is intuitive to understand but difficult to fully encapsulate within a single statement. For the purposes of this research, the definition of cognitive workload is that proposed by Meshkati: “mental workload [is] a multidimensional construct that reflects the interaction of such elements as task and system demands, operator processing capabilities and effort, subjective performance criteria, operator information processing behavior and strategies, and finally, operator training and prior experience” [1]. Since cognitive workload is a theoretical construct that exists in the mind, there are no methods of measuring it directly. Instead, loading tasks can be used to affect workload while measurements, such as psychophysiological readings, can be examined that allow for inferences to be made about cognitive workload. With respect to this research, descriptions about measuring cognitive workload are in fact referring to the indirect measurement of factors affected by workload and making inferences about the underlying load. Another term commonly used in this line of research is “operator functional state” or OFS. The definition of OFS used in this study is that put forth by Hockey: “the variable capacity of the human operator for effective performance in response to task and environmental demands, and under the constraints imposed by cognitive and physiological processes that control and energize behavior” [2]. In many applications, such as the current research and various other similar studies, cognitive workload and OFS can be seen as the same construct and used interchangeably. For the sake of clarity, cognitive workload is the term that is

used throughout this research, although studies that are referenced herein may use the term OFS.

2.1.2 Participant Study

Participant study, for the purposes of this research, refers to the data collection process that involved the help of participants. The end result of the participant study was the production of psychophysiological data that was then processed and analyzed as part of the remaining research.

2.1.3 Classification System

In this research, system refers to the processes and Python scripts that acted on the data collected from the study. This includes everything from the preprocessing of sensor data to the training and evaluation of classification models.

2.1.4 Model / Classifier

Another term that is used throughout this research is “model”. In this case, the more mathematical definition of model that is used in the context of machine learning is intended, rather than a strictly psychological definition. A model, as used in machine learning, is a system that can take in input data and produce some output as defined by patterns present in historical data. This system may be a black-box, which cannot be defined or explained effectively by a person, or it may be described with well defined rules, depending on the technique used. Since the current research deals with the classification of data, the term classifier is also used to refer to the model. It should be noted that “classifier”, which is the trained model, is distinct from “classification method”, which is the particular format of the model, such as artificial neural network.

2.2 Cognitive Workload

2.3 Psychophysiological Measures

2.4 Random Forest

2.5 Artificial Neural Network

3 Related Work

Classifying cognitive workload has been the subject of many studies. The majority of studies have dealt with creating an individual model unique to each participant, with the study consisting of only one task. These studies commonly use EEG and heart rate as psychophysiological features [3–6]. While EEG data is consistently used in cognitive workload studies, the subsequent features generated from the data vary. Some of studies, such as the work performed by Wilson et al. [5], use log power spectrum information from the EEG measurements while other studies, such as that by Zhang et al. [4], use combinations of EEG power spectrum information, called task load indices. In addition to EEG and heart rate, some studies include blink rate [5, 7] and respiration rate [8] as features. Other studies used blinking and eye movement measures to correct artifacts in EEG data [3]. The general layout for each of the studies was roughly the same; participants were attached to a variety of sensors and then asked to perform a loading task at varying difficulty levels to induce different cognitive workloads. The two most common loading tasks used were the Multi-Attribute Task Battery (MATB) [9] and AutoCAMS [10]. These systems are commonly used due to their ability to systematically vary the difficulty settings.

A variety of different classifiers have been used in studies in an attempt to model cognitive workload. The study covered by Yang et al. [6] uses a system of fuzzy inference rules to perform realtime classification of cognitive workload for the purposes of adaptive automation triggering. An SVM Regressor was used by Ke et al. [11] to produce a cognitive workload index, a particular number corresponding to cognitive workload level. Another study [3] explored the use of an Adaptive-Network-based Fuzzy Inference System that was trained using differential evolution and ant colony search as a means of predicting levels of cognitive workload. Fuzzy C-Means clustering was used in a study [4] as a means of classification through clustering. The study completed by Wilson et al. [5] used an artificial neural network. All of these studies have shown adequate results but have failed to address the larger scope of

generalizability. One study [12] used a bayesian model to explore the ability of a model to be used on multiple subjects, learning from and tested on a dataset consisting of data from multiple participants. What the study lacked, however, was the ability to handle novel data; that is, classify on a participant that was previously unseen. A different study [11] falls in a similar category, however it deals with multiple tasks rather than multiple participants. The study involved testing on data from an unseen task, showing the possibility for full generalization. An additional study [13] demonstrated how cognitive workload metrics for a single individual can vary substantially across multiple days. This adds an additional challenge to generalizability as the model not only has to account for different subjects and tasks, but also the variation within subjects across multiple days.

4 Problem Statement

While many studies have focused on individualized models of cognitive workload using machine learning, few studies have produced results that are of use outside of a controlled lab setting. In a practical application, the training of a model to each individual for every task they may perform would be far too time consuming and costly. This research investigates the effectiveness of a generalized model of cognitive workload based on psychophysiological measures. In order for the model to be generalized, it should be robust against differences between individuals and between the tasks being performed. It should also be effective at handling novel individuals and tasks. Such a model may or may not be reasonably developed due to the complexity of the human mind. Ideally the resulting model could be used outside of a controlled setting with automatic data processing and realtime classification.

A driving factor behind this work is usability in a real-world setting. As such, design decisions were made that favor low-cost, easy-to-use sensors with little to no empirical cleaning of data. The models were trained using the machine learning techniques of artificial neural networks and random forests and tested in a variety of configurations to determine their effectiveness. As is common in this type of study, data for experimentation was collected with sensors from subjects as they participated in cognitive loading tasks of varying difficulties. While no methods used would inhibit realtime analysis, the system focused on offline data. In the end, this research hopes to address the following hypothesis: a generalized model

of cognitive workload that is more effective than random and that can be trained using methods suitable for real-time classification without manual data processing.

4.1 Applications

A fully generalized model of cognitive workload would be of use in a variety of applications. In some instances it would simply improve on areas where cognitive workload models are already used but in other instances it would open the door to new research.

4.1.1 Adaptive Automation

One application in which cognitive workload models are already used is adaptive automation. Adaptive automation, as described by Byrne and Parasuraman, is automation in which the “assignment of tasks between the human operator and automation is dynamically adjusted based on task demands, user capabilities, and total system requirements to promote optimal system performance” [14]. It has been found in previous studies that static automation, automation that is constant and which puts the operator in a state of monitoring for failures, can result in a deterioration of performance due to the strain of monitoring [15, 16]. In adaptive automation, through the use of a cognitive workload model, the state of the operator can be monitored and the automation system adjusted accordingly to maintain the desired balance between manual and automatic control. A concept of an adaptive automation system using participant specific cognitive workload models has been done by Wilson, showing improvement in the overall task performance [17]. The addition of a generalized model of cognitive workload would greatly improve the flexibility of an adaptive automation system as it would eliminate the need for training on each individual user.

4.1.2 Research

While some tasks, such as the Multi-Attribute Task Battery (MATB) [9] and AutoCAMS [10] were specifically developed to have control over the level of the load placed on the user, many other tasks can be much less predictable. Therefore, while it is possible to train individualized models of cognitive workload effectively on specialized laboratory tasks, it is unlikely that the same level of training could occur with more complicated tasks that are reflective of other real-world situations. In

many cases, knowledge gained while studying one situation or task is used to make predictions and assumptions about another domain. It has been noted, however, that special care must be taken when dealing with analogous systems as the situations could be more different than they seem due to underlying complexity [18]. If a generalized model of cognitive workload were created that was effective independent of task, it could be trained on the specialized loading tasks, such as MATB, and then used in other domains to better research the workload of tasks, without the need for analogous systems and the danger of faulty comparisons.

5 Participant Study

A study was conducted in order to generate and collect the data needed to train the appropriate models. While some previous research [12] has relied on the use of established and validated data collected from other experiments [19], the amount of data required to adequately train artificial neural networks was not available. Consequently, a new study was required to generate the large amount of data that was necessary. The participant study consisted of two tasks, each with a leading baseline trial and three load conditions. The load condition trials were ten minutes long and the baseline trials were five minutes long. The initial design of the study specified fifteen minute long trials but feedback from the pilot group suggested that at fifteen minutes, task fatigue started affecting performance.

The participant study began with a total of ten college-aged participants, similar to existing works [3,4,19–21], however, due to complications, only seven participants completed the entire process. Two of the removed participants were excluded due to misfitting of the EEG headset as well as discrepancies between the task difficulties reported versus the average difficulties reported by the other participants. The last removed participant was excluded due to personal issues preventing the completion of the tasks.

Of the seven participants that completed the participant study, the ages ranged from 18 to 23 years old, with an average age of 21.1 years. The gender ratio was balanced with three males and four females. Participants lacked any uncorrected hearing or visual impairments that would affect performance on the tasks of the study. The following subsections describe each of the tasks as well as the overall layout of the study in its entirety.

5.1 Multi-Attribute Task Battery (MATB)

The first loading task presented to participants was the Multi-Attribute Task Battery (MATB), a task loading system originally developed by NASA [9]. MATB is a system that has been used in a variety of studies as a means of task loading for the collection of psychophysiological data [5, 12, 22]. The system consists of a multitude of subtasks arranged in such a way as to simulate the tasks a pilot might be required to perform during flight. The particular version of the system used in this study was AF-MATB, an updated version of the original software [23]. This version is very similar to the original with the majority of the changes affecting how it runs on modern operating systems and adding different options for subtask automation. All of the subtasks of MATB were used in the study, while the scheduling view that shows upcoming events was disabled. Input to the system was in the form of joystick and either keyboard or mouse. Participants were free to decide whether to use the keyboard, mouse, or a combination of the two with the right hand while the joystick was controlled with the left hand. In Figure 1, the main view of MATB with which participants interacted is shown. This includes resource management, systems monitoring, communications, and tracking.

The resource management subtask involved maintaining the amount of fuel that was present in two main tanks within a small acceptable range. Secondary tanks were present and could be used to temporarily store fuel. The participant was required to toggle pumps, controlling the flow of fuel between the tanks. Throughout the trial, pumps would temporarily break, or silently shut-off, preventing the flow of fuel through the pump. The number and frequency of these occurrences depended on the difficulty condition of the trial.

The systems monitoring subtask were centered around two buttons and four gauges. One button was to be kept on while the other button was to be kept off. Throughout the simulation, the buttons changed states and the participant was required to click on the button indicator, or press a corresponding key, to revert the button to the correct state. Colors were used to indicate the state of the buttons. Background color indicated that the button was off, while either red, for the “off” button, or green, for the “on” button, was used to indicate on. The gauges presented had sliders that would vary in vertical position throughout the simulation. Occasionally, depending on the difficulty condition, a slider would move outside the acceptable range for the gauge and require interaction from the participant, either

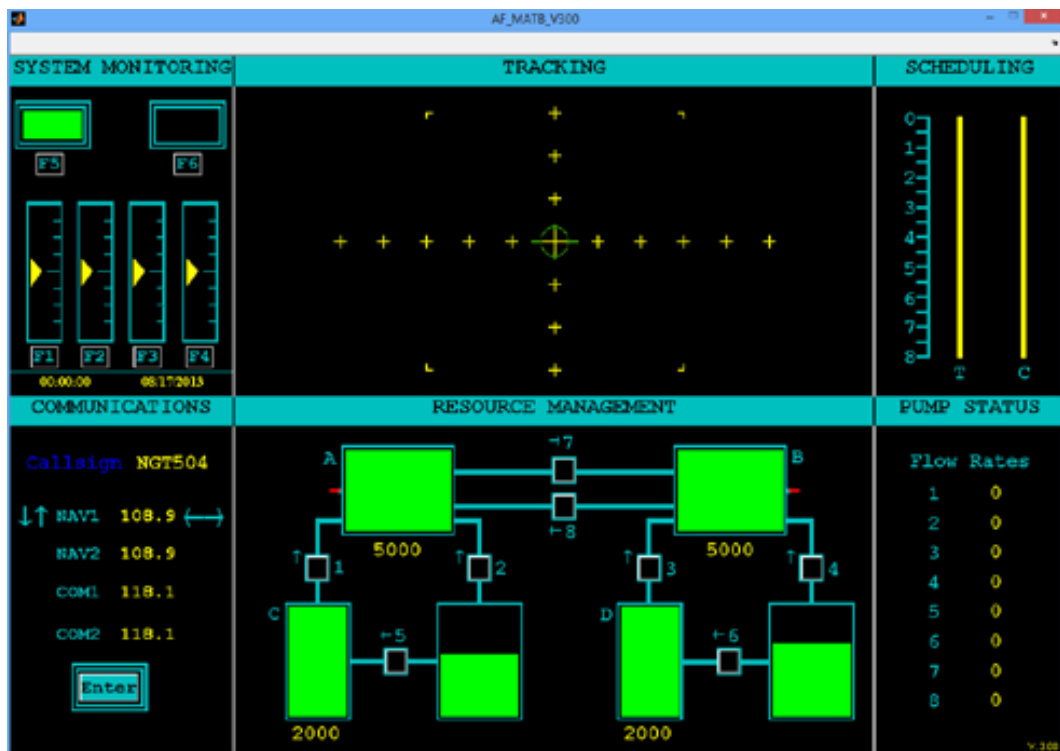


Figure 1: Main view of the MATB simulation software.

Table 1: Parameters used in MATB script generator specifying subtask difficulty and event occurrences per 10-minute trial

Load Condition	Tracking	Communication		System		Resources	
		Target	Distractor	Lights	Gauges	Failures	Shut-offs
Low	Low	6	2	48	42	4	2
Moderate	Moderate	18	10	96	102	20	10
High	High	30	12	110	120	40	20

through a mouse click or key press, to reset the gauge. If an out-of-state system component was not addressed within the response time period of ten seconds, the component was automatically reset to its defaults, within-range state.

The communication subtask required the use of sound. The participant was assigned a particular aircraft callsign. For the duration of the trial, the participant was required to listen to their callsign and respond to the instructions. The communication window included four different communication channels, each set to a particular frequency. Throughout the simulation, recorded audio tracks would play, addressing a particular callsign to change a channel to a desired frequency. If the callsign matched that of the participant, they were to change the specified communication channel to the appropriate frequency using either the mouse or the arrow keys. If the callsign did not match that of the participant, the request was to be ignored. The number and rate of true- and false-alarm requests depended on the difficulty condition of the trial.

The tracking subtask was the only one that required the use of the joystick. The participant simulated steering a plane by keeping a reticle within the acceptable parameters of the tracking window. The amount of drift affecting the reticle and the amount of control the joystick emitted onto the reticle varied throughout the simulation in accordance with the difficulty condition of the trial.

The baseline condition for MATB involved the participant watching the screen while the simulation ran with the low condition parameters and automation of all available tasks enabled. This allowed for the visual and auditory stimulation without the cognitive workload. **Presented in Table 1 are the parameter values that were entered into the script generator of MATB to produce the trials for each load condition. Each parameter, with the exception of tracking difficulty, represents a type of event that can occur during the simulation. The value for a parameter specifies how many times the event should**

Table 2: RanTask Desired Tone Counts

Load Condition	Tone Frequency		
	Low	Moderate	High
Low	0	5	0
Moderate	2	0	3
High	2	2	3

occur over the length of the 10-minute trial. Two scripts were generated for each difficulty condition to ensure that each trial was different and the participants could not memorize any patterns. The parameter values used are based off those in a previous study [24]. The majority of the parameter values are scaled from the referenced values to account for the difference in trial lengths, however the parameters for “lights” and “gauges” in the high condition were further modified due to errors in scheduling from the script generator.

5.2 RanTask

The second loading task used in this study is a custom tone counting task, referred to herein as RanTask. This task is an extension of the mental workload loading task used in the work of Rantanen et al. [25]. The participant was presented with three different tones in a random order. The participant was to count the tones that were presented and press the key corresponding to the tone when the desired count was reached. The intended count differed with the difficulty condition. The tones corresponded to the notes B at 493 Hz, F at 698 Hz, and A at 880 Hz. Each tone was presented for one second with a one second interval between tones, giving the participant two seconds to respond to a given tone. Each time a tone was presented, a log entry was created, recording the tone presented, its number in the sequence, and whether, if a participant pressed a key, it was correct or incorrect. If a key was pressed at an incorrect time, a false-positive was logged and the sequence count for the corresponding tone was reset to prevent cascading failure, e.g. counting every five correctly but being off by one from the ground truth sequence would result in all false-positive without the reset.

In the baseline trial, the participant was not required to count any of the tones, only listen. Presented in Table 2 are the desired counts for each tone in each load condition. This differs from the originally intended counts due to feedback from

the pilot study members. Feedback indicated that the original low- and moderate-load conditions, both requiring the counting of a single tone, were too similar in difficulty. As such, the moderate-load condition was replaced with the existing high-load condition and an additional condition was added that required keeping count of all three tones. The participants were given a copy of the count information from Table 2 at the time of the trial to ensure they were aware of which tones were to be counted. The participants were also informed that the use of fingers for counting or any other external methods of keeping track of tones was prohibited.

5.3 NASA Task Load Index (TLX)

The NASA Task Load Index (TLX) assessment is a means of determining the workload for a given task based on participant responses [26]. The procedure allows for exploring the different aspects of a task and where the underlying cognitive workload originates. The workload rating obtained from participants through the TLX process was used to compare the subjective difficulties of the various tasks, similar to the procedure used in existing studies [11, 24]. The current TLX procedure consists of six subscales: Mental Demands, Physical Demands, Temporal Demands, Own Performance, Effort, and Frustration. The Mental Demands subscale measures how much mental activity was required and its complexity. Physical Demands measures how much physical activity was required and how laborious it was. Temporal Demands measures how much time pressure was present in the task. In the Own Performance subscale, the participant remarks on how successful they felt they were and how satisfied they are with their performance. The Effort subscale measures how hard the participant had to work to accomplish the task. The final subscale, Frustration, measures how stressed or relaxed the participant felt during the task.

Individual calibration was required from each participant for both of the tasks that were performed. During this calibration, the participant performed pair-wise comparisons between the TLX subscales, indicating which subscales had more effect on the overall workload of the task. The calibration was then used as a method of weighting the individual TLX survey responses for a trial. The idea behind the weighting is that a component that is calibrated as having a large impact on the overall workload should have its rating affect the overall workload metric more heavily.

5.4 Structure

The study took place over four days for each participant. Each day required roughly an hour and a half from the participant. The time is approximate as it was dependent on how quickly the sensors were placed, which was, at times, delayed due to conditions such as hair. The first day consisted of an introduction to the system: what sensors were to be utilized, the intent of the study, and what was to be expected of the participant. Training was also conducted on each task to familiarize the participant with what was to be expected. Ideally, stable performance on both tasks was achieved for all participants as a result of training. Other studies [5] have included lengthier training periods, however, due to time constraints, the time required for complete stability was not feasible. As such, the task performance results of each participant were recorded and examined as a factor in the overall effectiveness of this study.

The second day marked the beginning of data collection for the participants. The participants were attached to the sensors and a five minute baseline recording using MATB was taken. Following the baseline, the participants each completed a ten minute MATB trial on the low-load condition. A NASA-TLX survey was then administered to record the participants' rating of the task workload. Next, the participants completed a second, ten minute MATB trial on the low-load condition. This again was followed by a NASA-TLX survey. Two TLX surveys were taken for each condition to evaluate the stability of the rating as participants should be recording similar ratings for both trials of a condition. After a ten minute break, the same procedure was repeated, including the baseline, substituting RanTask for MATB as the task being performed.

The third and fourth days followed the same format as the second day. The third day consisted of the tasks being performed on the moderate-load condition while the fourth day consisted of the tasks being performed on the high-load condition. Each participant performed all their trials at similar times of day to reduce variations caused by circadian rhythm. The order in which the load conditions were presented follows that of a previous study [5]. While the sensors were removed and replaced on multiple occasions as the data collection occurred over many days, which differs from some studies [4, 5, 11], this should not have an effect on the overall validity of the data [24].

6 System

The system that was created encompassed a full classification workflow. It covered the collection of data from the sensors, preprocessing to restructure the data into a more usable format, processing to generate features, training of classifiers, and evaluating the classifiers with the appropriate datasets. The system is modular, such that each task can be performed individually, allowing for flexibility with future extensions of functionality. In Figure 2 an overview of the entire system is shown. In the following subsections, the various components of the system will be discussed in detail.

6.1 Sensors

Data was collected through the use of two sensors, an EEG monitor and a heart rate monitor. The collection of psychophysiological data through EEG and heart rate monitors has been performed in numerous existing studies [5, 6, 12] and has been assessed [27]. Following with the theme of practicality in a production system, the sensors used are more readily available than those used in other research settings. Data collection was structured such that there was one heart rate data file and one EEG data file per trial period, as opposed to recording multiple trials in a single file. As such, the data recording process was started shortly before and stopped shortly after each trial, resulting in a small amount of extraneous, noisy data preceding and trailing the trial data.

6.1.1 Electroencephalogram (EEG)

The electroencephalogram (EEG) monitor used is a wet-contact, wireless headset called the Emotiv EPOC+¹. The EEG monitor is a research grade headset with 14 individual data channels and two reference channels. Electrodes on the headset are located at the following placement sites as defined by the International 10-20 System for electrode placement [28]: AF3, F7, F3, FC5, T7, P7, O1, O2, P8, T8, FC6, F4, F8, and AF4. Electrodes placed at P3 and P4 serve as references for the headset. The predecessor to this headset, which is very similar to the current headset, has been proven effective in many studies including work by Knoll et al. [29]. The design of the headset allows for accurate placement of electrodes over many trials

¹Emotiv EPOC+ product details available at <https://emotiv.com/epoc.php>

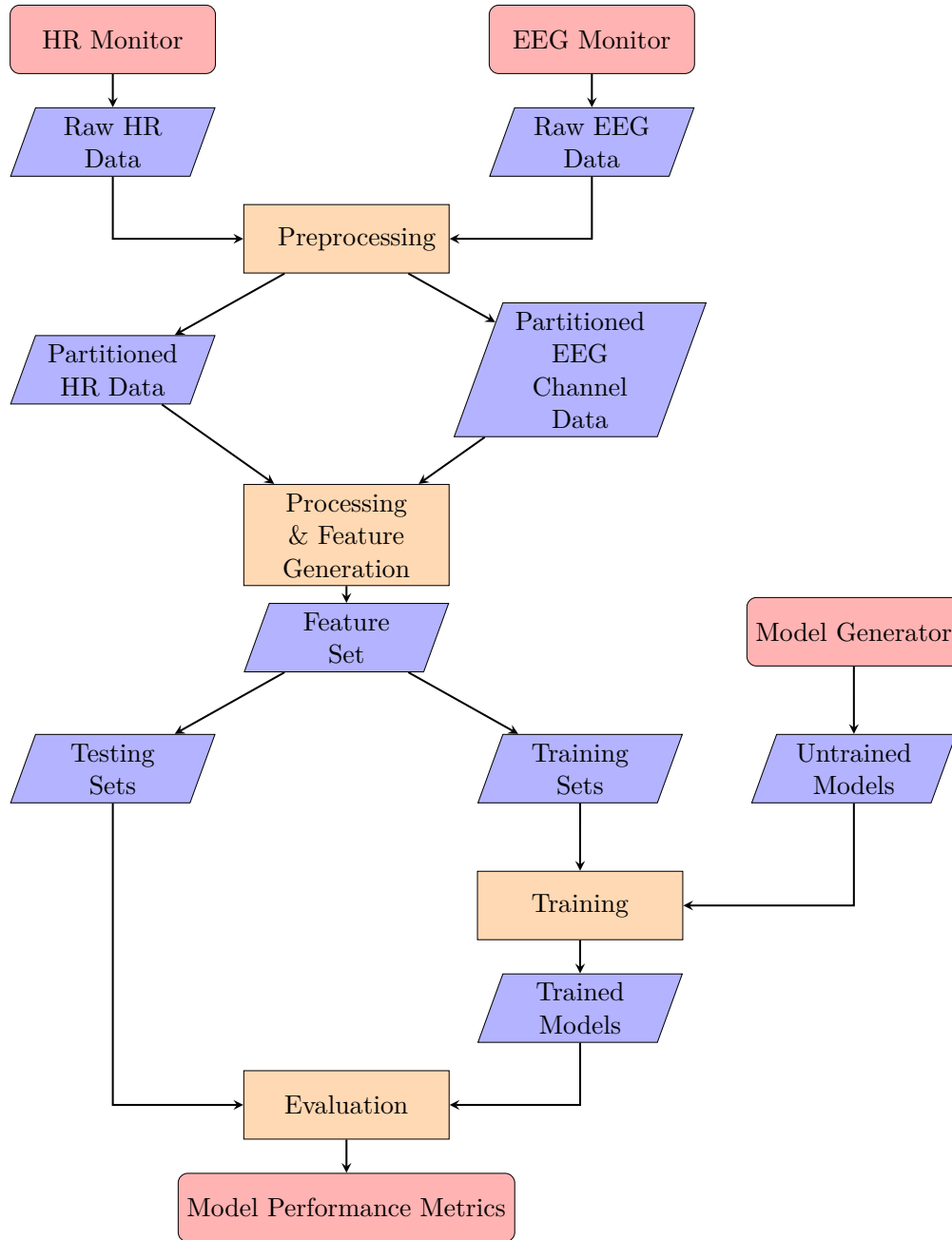


Figure 2: A general overview of the entire system workflow.

and does not need an expert for proper setup due to the presence of rubber guides that are placed on the mastoids of the participant. The use of a headset as opposed to stand-alone electrode contacts does lead to some issues regarding proper scalp contact with individuals having smaller-than-average skulls or large amounts of hair. Data collection software for the EEG headset allowed for monitoring of contact quality of the electrodes, ensuring proper readings were taken. Thanks to the work done by Estepp et al. [24], it can be assured that the removal and replacement of electrodes between data collection trials had little to no effect on the quality of the measurements. The headset operates at a resolution of 128 samples per second. This high resolution resulted in more data and reduced the amount of variation caused by instantaneous noise.

6.1.2 Heart Rate (HR)

The heart rate monitor used was an optical, arm-based monitor called RHYTHM+ by Scosche ². The monitor was designed for collection of vitals during fitness training, not research, but it was sufficient for the monitoring of participants. At a resolution of four samples per second, much less data was collected for heart rate than EEG. This was not an issue, however, as heart rate change is gradual and unlikely to vary rapidly. An arm-based monitor was chosen as it was more comfortable and convenient to participants than a chest-strap monitor. The monitor communicated with the data recording computer through a protocol called ANT ³, a protocol commonly used within fitness equipment. The monitor can easily be adjusted and was unobtrusive enough to be worn for extended periods of time.

6.2 Data Preprocessing

The preprocessing of data files was necessary to transform the data logged from the sensors into a form that was more easily used to produce features. Rather than modifying the format in which the sensors logged the data, this step was used to extract the information useful for feature generation and write it out to multiple files. This step also encompassed cleaning of the data, addressing issues caused by communication between the sensors and receivers. Each participant trial was preprocessed independently and required no information that was not already

²RHYTHM+ product details available at <http://www.scosche.com/rhythm+>

³Additional information on the ANT protocol can be found at <http://www.thisisant.com>

contained within the trial data. The majority of the preprocessing for heart rate and EEG data was separate, as is described in the following two subsections, however once the data had been cleaned, joint processing was required to properly align the times of the two datasets. As mentioned previously, the data files included leading and trailing noise that was not collected during the trial period itself. This data needed to be removed as it was not obtained under the load conditions and had no appropriate label. Part of the information that was logged during task trials was the start time of the trial. This time was used as a starting point when looking for a timestamp on the heart rate data that matched a timestamp in the EEG data. Once a time was found in which EEG and heart rate data existed that was after the start time of the trial, five-second intervals of data were taken until ten-minutes, the length of a trial, of data was processed. This not only removed the leading and trailing noise, it also synchronized the timing between the heart rate and EEG data so that the first five-second interval in the heart rate data corresponded to the first interval in the EEG data.

The data was split into intervals in an attempt to reduce the effect of noise. The use of segmentation is common when dealing with psychophysiological data with a variety of segmentation conditions having been used in previous studies, ranging from two-second intervals with no overlap [22] to 40-second intervals with 35-seconds of overlap [12]. An interval length of five-seconds, a length that has been used in a previous study [21], was chosen as it is long enough to benefit from averaging and noise reduction but short enough that updates to the participant state are useful to other systems in an on-line manner.

6.2.1 EEG

The preprocessing of the EEG data involved multiple steps to transform the data from a single EDF file into multiple, separate channel text files. The process that was followed is outlined in Figure 3. The first step using EEGLAB [30], an open-source MATLAB toolbox for processing EEG signals, to read the information out of the EDF file created by the sensor logging and to separate the individual channel streams. A baseline removal process was performed using EEGLAB on each channel, where the baseline of a channel is defined as the mean value of the channel for the time period. Once the baseline was removed, each channel was written out to its own intermediate, tab-separated text file. The channel data was then partitioned into five-second intervals in the method previously described such that each interval

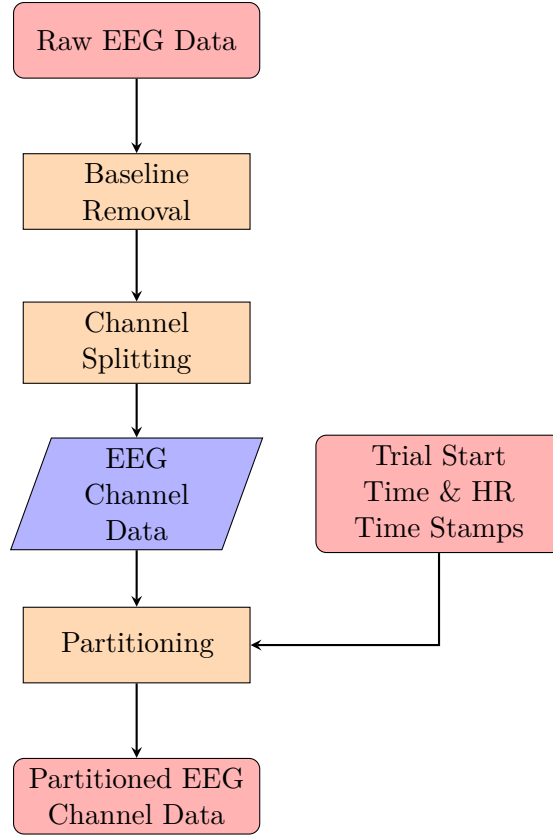


Figure 3: An overview of the steps involved in the EEG preprocessing workflow.

fell within the time frame of the trial and aligned with the heart rate intervals. Each channel is then, once again, written out to its own tab-separated text file. The advantage of multiple, intermediate text files is that modifications can be made to later portions of the workflow and applied without a complete reevaluation of earlier steps.

6.2.2 Heart Rate

The recorded heart rate data required less preprocessing than the EEG data due to it being only a single channel. In Figure 4, the overall process followed is shown. Once the data was read in, it was smoothed to ensure there was at least one data point per second. While ideally there would be multiple heart rate records per second, issues in communication between the sensor and the receiver resulted in some periods of time in which no readings were recorded. Since heart rate undergoes gradual

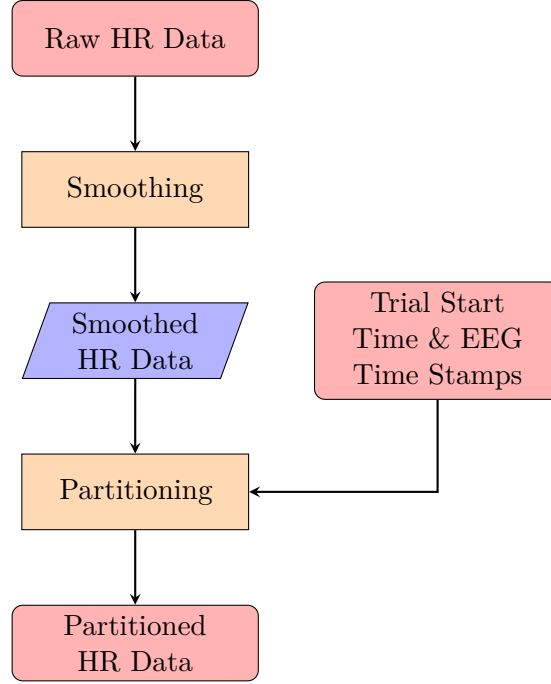


Figure 4: Overview of the heart rate preprocessing workflow.

change, missing data points were interpolated from the previous and following heart rate readings. After the data was smoothed, the heart rate records were partitioned into five-second intervals, as previously described, and written to a tab-separated text file.

6.3 Data Processing and Features

Data processing to produce the desired features was done after the heart rate and EEG data had been cleaned and preprocessed. A total of 72 features were generated for each five-second interval. Of the 72 features, 70 were EEG features and two were heart rate features. Figure 5 presents an overview of the processing steps that changed the preprocessed data into a collection of feature vectors. In addition to the feature values, each feature vector record is labeled with a class value corresponding to the load-condition of the trial from which the data originated, either low, moderate, or high. The data for each task of a participant, at this point labeled with the appropriate class, was combined together such that it could be handled more easily by the classifiers. Data records of each task for each participant were kept separate

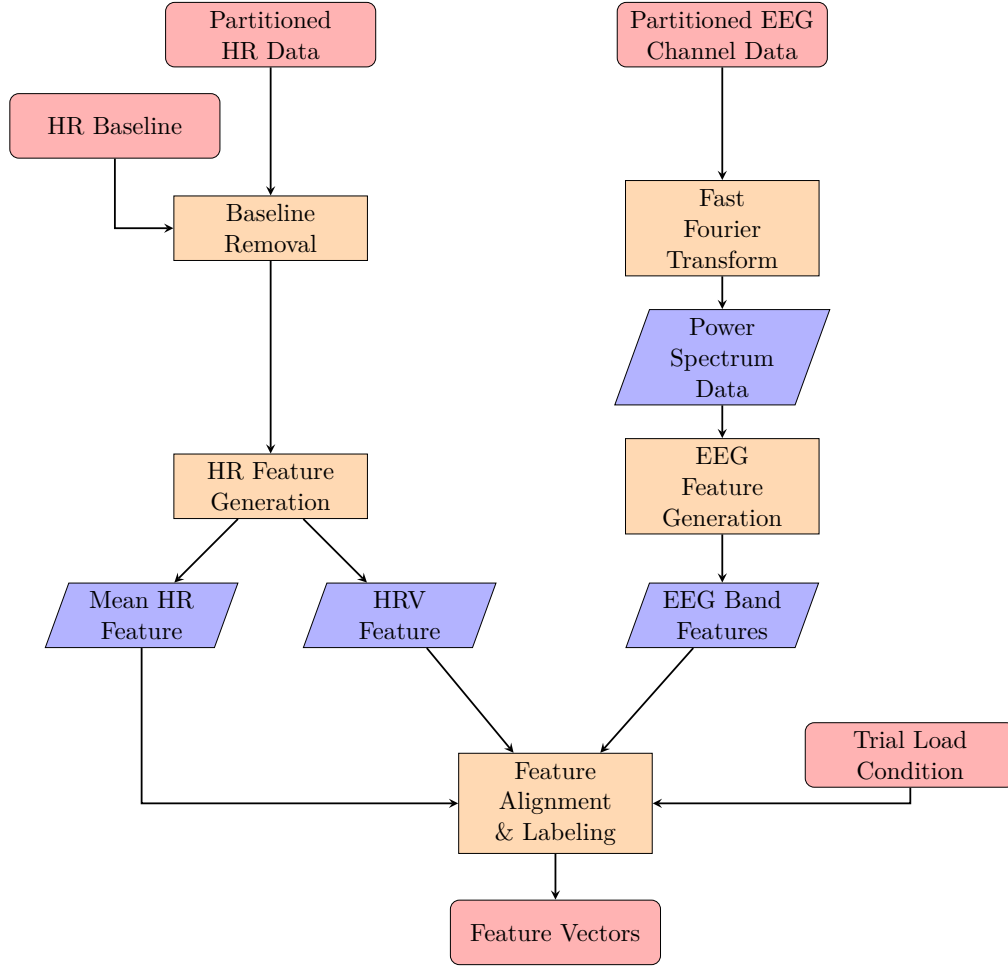


Figure 5: Overview of the processing and feature generation workflow.

to ensure the appropriate training and testing datasets could be created during the classifier phase.

6.3.1 EEG

The end goal of processing the EEG data was the production of band information resulting from the spectral power analysis of each channel. The processing and feature generation was performed on each channel independently before the features were joined together into one record. The features for a channel were created through the application of a Fast Fourier Transform on each five-second interval of data. This transform produced power spectrum information of the EEG channel for the time

interval. The bands used as features consisted of delta (1-3 Hz), theta (4-7 Hz), alpha (8-13 Hz), beta (14-30 Hz), and gamma (31-42 Hz). These are the same bands that were used in the work done by Wilson et al. [5], as well as various other studies [?, 24]. The band values were computed by summing the logarithm of the power spectrum values that corresponded to the frequencies in each band. The logarithm of the values was taken to normalize the numbers. The five band values for each channel were then combined to produce the complete vector of 70 EEG features. It is important to note that no empirical noise reduction or artifact correction was done between the collection of EEG data and the creation of the features. While this modification would likely increase the effectiveness of the models and is indeed done in many other studies [20, 22, 24], the driving factor behind this study was to have a system feasible for production, and an expert will not always be available to hand filter the data in a live system.

6.3.2 Heart Rate

The heart rate data contributed only two features to the overall feature vector. Before any features were generated, however, the mean heart rate of the baseline trial was removed from every data point. The idea is that each participant likely has a different resting heart rate; removing the mean of the base heart rate will create offset readings that are caused by the load of the trial. Since removing the mean heart rate could potentially create non-positive numbers and complicate further computations, a constant offset was added across all participants. This shift factor did not affect any patterns that may exist in the data but prevented calculation errors such as divide-by-zero. The first feature to be created from this modified data was the average heart rate for each five-second interval, calculated using mean. The second feature used was heart rate variation or HRV, the ratio between the standard deviation and mean of heart rate for the five-second time interval, similar to that used by Zhang et al. [4]:

$$HRV = \frac{\theta_{HRV}}{\mu_{HRV}}$$

where θ_{HRV} is the standard deviation and μ_{HRV} is the mean of the HR data for a particular time segment.

6.4 Configurations

The feature vectors generated from the psychophysiological data was combined in a variety of ways to create different training and test dataset configurations. While the end goal was a model that was both participant- and task-independent, a number of benchmarks along the way provided additional information on the models’ generalizability. Each component of a configuration has three possibilities: same, all, and cross. ‘Same’ refers to being trained and tested on the same individual or task. This setting serves as a control, preventing the component from affecting the generalization. ‘All’ refers to combining all the data from the participants or tasks to create both the training and test datasets. This setting explores the classifiers’ ability to differentiate between the individual patterns or create a single pattern that is representative of all the participants or tasks. ‘Cross’ refers to training on a subset of the participants or tasks and testing on the remaining participant or task. This setting is the end goal of generalizability as it allows the model to be trained with a specified dataset and then used with a variety of novel participants and tasks. All training and testing dataset splits were stratified with respect to load condition to ensure models could be adequately trained. Additionally, in any situation where both training and testing datasets were created from the same data, as is the case with ‘same’ and ‘all’ configurations, three-fold cross validation was used to ensure results were not dependent on anomalies in the specific data split.

The following subsections described each of the configurations that were explored in more detail. For convenience, a naming system in the style of acronyms has been created for referencing particular configurations; e.g. “SP-ST” refers to “Same Participant - Same Task”.

6.4.1 Same Participant - Same Task (SP-ST)

Each classifier was first trained and evaluated in a manner similar to previous studies [3–5, 21]. This served as a baseline to compare the overall performance of the models to those created in other studies. The main use of this configuration was in determining the effectiveness of the created features. As part of this configuration, the data for each participant on each task was used for both training and testing. This means the classifier was dependent on both the participant and the task. This configuration involves 14 classifiers per classification method, one for each task for each person.

6.4.2 All Participants - Same Task (AP-ST)

One step up in generalizability is the creation of a classifier that is trained on multiple subjects, similar to that of a previous study [12]. The data for each participant for each task was split into training and testing subsets. The training and testing subsets for all participants on a particular task were then combined to create joint datasets. This explored the possibility of finding patterns in the data that match multiple participants, as opposed to a unique pattern for each participant. This configuration involved two classifiers per classification method, one for each task.

6.4.3 Same Participant - All Tasks (SP-AT)

Another configuration that was explored is similar to the previous configuration, All Participant - Same Task, except the task data was merged, rather than the participant data. The configuration started the same way as the previous two configurations, with the data of each task for each participant being split into training and testing subsets. Then, the training and testing datasets for each task for a participant were combined. A total of seven classifiers for each classification method were trained for this configuration. This explored finding patterns that exist between the tasks, given a particular participant.

6.4.4 All Participants - All Tasks (AP-AT)

The configuration of combined participants and combined tasks is another step up in generalizability. This configuration is a combination of the previous two configurations. In this configuration, all the data was split into training and testing subsets which were then combined across all participants and all tests. The models trained from this configuration explored increased generalizability over existing studies but did not explore the model’s capabilities with completely novel data. Only a single classifier per classification method was necessary for this condition.

6.4.5 Cross Participant - Same Task (CP-ST)

This configuration is the first that deals with true generalizability. For this configuration, the data for six of the participants for each of the tasks was combined to form the training set. The last participant served as the test set. Rather than differentiating between patterns previously seen, as the All Participants - Same Task

configuration does, this configuration explored how well the patterns learned for a set of participants can generalize to novel data that may contain a distinct, and previously unseen pattern. This configuration resulted in 14 classifiers for each classification method, accounting for the two different task datasets and each participant being the unseen test data for a model.

6.4.6 Same Participant - Cross Task (SP-CT)

Another configuration that was explored is a variation of the previous configuration, substituting participant crossing with task crossing. In this configuration, each classifier was trained on the data for one task for a participant and tested of the data for the second task. This classifier explored how well the patterns discovered in one task generalize to a separate, unseen task. A total of 14 classifiers per classification method were created, allowing for each task to be the test data.

6.4.7 All Participant - Cross Task (AP-CT)

One step down from complete generalization involves combining the previously described options, ‘all’ and ‘cross’. This configuration involved combining the data for all participants, resulting in one dataset for each task. A single classifier was then created that was trained on the data for one task and tested on the second task. This configuration is similar to the previous configuration, but rather than a unique classifier for each participant, only one was created. This configuration is similar to a previous study [11]. This explored the possibility of a general pattern being found for cognitive workload that can be used independent of task. The configuration resulted in two models per classification method.

6.4.8 Cross Participant - All Tasks (CP-AT)

This configuration is similar to the previous configuration, however it explored the independence of participants instead of tasks. The configuration involved combining data for all tasks, resulting in seven datasets, one for each participant that included both tasks. A classifier was then trained on the data from six participant and tested on the last participant. After repeating the configuration, to test on each participant, a total of 14 models were created.

Table 3: Model Count Per Configuration

	SP-ST	AP-ST	SP-AT	AP-AT	CP-ST	SP-CT	AP-CT	CP-AT	CP-CP	Total
Number of models	14	2	7	1	14	14	2	7	14	75

6.4.9 Cross Participant - Cross Task (CP-CT)

The final configuration represents complete generalizability. It explored the independence of both the participant and the task at the same time. The setup for this configuration involved training a classifier on the data of one task for six participants. The classifier was then tested on the data for the last participant performing the other task. To achieve full test coverage, this condition was repeated 14 times; once for each participant / task test combination.

6.4.10 Overview

A large number of models were created to test different configurations of training and testing data. Table 3 displays how many models were created for each configuration as well as the total number of models created. The table only shows the number of final models, during the training and tuning process many more models were created, evaluated, and discarded due to suboptimal accuracy rates. The table also only shows the number of models per classification method, the actual number of final models is twice what is shown in the table due to two classifier methods being evaluated.

6.5 Classifiers

Two classification methods were used in the training and evaluation of models, artificial neural networks and random forests. During the training process, a grid search was used to evaluate various combinations of input parameters and tune the model. This tuning occurred for each final model that was created, rather than at a configuration level; e.g. for the 14 models required for the SP-ST configuration, the tuning process was repeated 14 times. Therefore optimal parameters for one model in a configuration are not necessarily consistent with the optimal parameters of another model in the same configuration. This process was automated, using the classification accuracy to determine which model was most effective, and no manual analysis of the tuning results was required. As such, the tuning is still

viable for a real-world classification system . The classification accuracy was used in determining which of the created models was most effective for a dataset. The following subsections detail what libraries were used in the creation of models as well as what parameters were used.

6.5.1 Random Forest

Ensemble classification methods, such as random forest, are not common in previous studies of this kind, even though they perform at least as well as other types of classifiers in most situations. The Scikit-Learn [31] Python library was used for the training of random forest models as it is easy-to-use and can be parallelized effectively. Only two input parameters were tuned in the training of the models: tree depth and number of trees. There are many other parameters that could have been tuned but they have less of an effect on the overall accuracy and each additional parameter being tuned greatly increases the number of models to be trained and evaluated. The number of trees was varied across the following values: 50, 100, 150, 300, 500, 750, 1000, 1500, and 2000. The number of trees controls how many individual decision trees are present in the random forest model and allowed to vote for a classification. The more trees in a model, the more accurate the approximation of the sample space and any patterns that are contained within, however too many trees leads to a large training time and possible overfitting to the training data. The second parameter that was tuned was max tree depth, which was varied across the following values: 50, 100, 200, 300, 400, 500, 600, 700, and 800. This parameter controls how deep an individual tree in the model can go, i.e. how many branch points a tree can have. Too few branches and decisions are made without looking at enough information but too many branches and the model becomes either overfit to the training data or reliant on features that do not affect the classification. All other parameters that were not actively tuned were kept at the default values as defined by the Scikit-Learn library.

6.5.2 Artificial Neural Network (ANN)

Artificial neural networks (ANN) have been proved in previous studies [5, 24] to be successful in classifying cognitive workload levels of single participants. This classification technique is now expanded to evaluate how effective it is at a generalized model. The ANN library used was the Fast Artificial Neural Network (FANN) [32]

library with Python bindings. The library handled all components of the ANN model, from training and parameter setting to evaluation. As the training of ANN models is able to be parallelized easily, multiple models were trained simultaneously, each model training on a single core. This, as well as the nature of ANN training, resulted in models that took far longer to train than the corresponding random forest models.

Artificial neural networks have a large number of parameters that can be changed to achieve different results. Due to the already lengthy amount of time required to train an ANN, the number of parameters that were tuned was kept to a minimum. A three layer neural network was used, having only a single hidden layer. The input layer contained 72 nodes, corresponding to the 72 features in the feature vectors. The output layer contained three nodes, each node representing one of the possible load condition classes. The number of nodes in the hidden layer was tuned from the following values: 72, 60, and 40. Another parameter that was tuned was the connection rate between the nodes. Connection rates of 0.7, 0.9, and 1.0 were evaluated. The last parameter that was tuned was the desired error rate. In the initial configurations, values of 0.01, 0.001, and 0.0005 were evaluated, however this was removed and a value of 0.01 was used for the later configurations as the training time was high and none of the models were able to reach the desired error rate, meaning further training was not constructive. While the initial models were trained with a maximum of 100,000 iterations, the majority of the ANN models were trained with a maximum of 50,000 iterations as little improvement was seen between 50,000 and 100,000 iterations and this decrease greatly reduced the amount of time spent training. The learning rate of the ANN models was not set as the FANN default training algorithm of RPROP [33] was used, which does not rely on a learning rate.

7 Results

While the main focus of this study is the accuracy of the trained models, results from the participant study are also analyzed to provide insight into the effectiveness of the participant study during data generation. The results of participant performance and of reported subjective workload is first examined to determine to what degree the tasks were successful in inducing the desired workload. Then the accuracies of the trained models are examined to determine under which configurations the

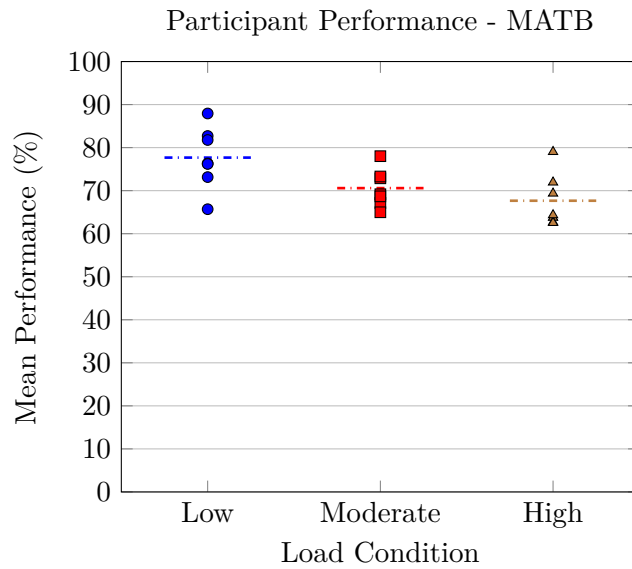


Figure 6

models were successful and if generalization occurred.

The graphs are currently all over the place. Once I have the results text in I will worry about straightening it out so it is not so messy. All of the graphs are currently within the results section, even if they do not appear to be. I am in the process of created a table of mean accuracies of the models that I will place in the discussion section, similar to the Participant Study information.

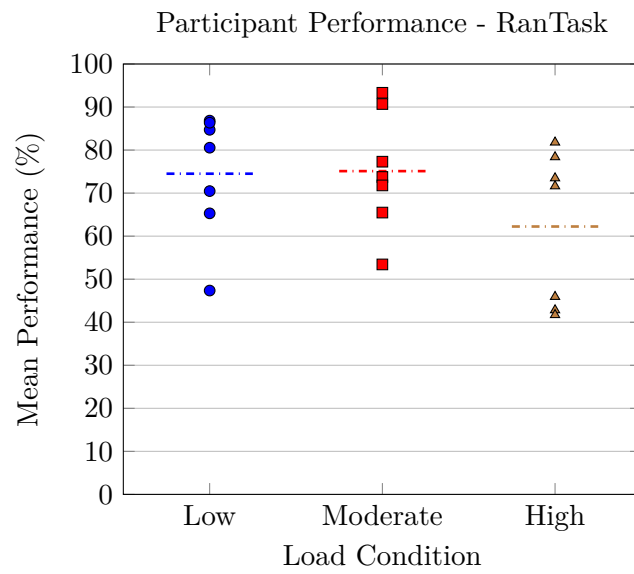


Figure 7

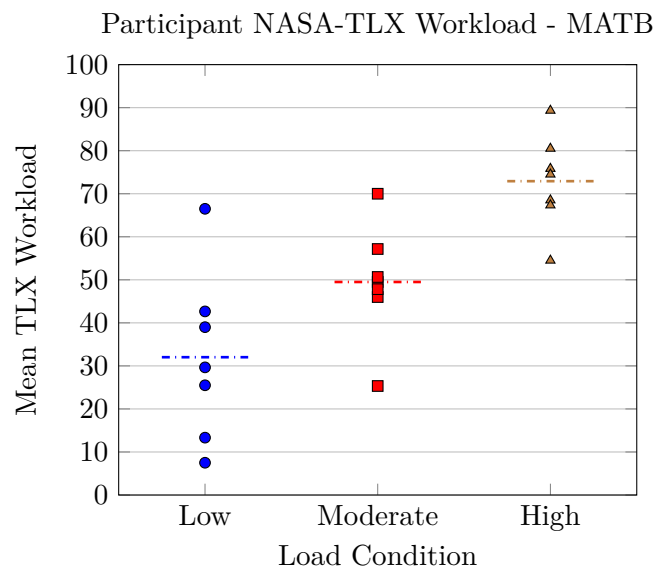


Figure 8

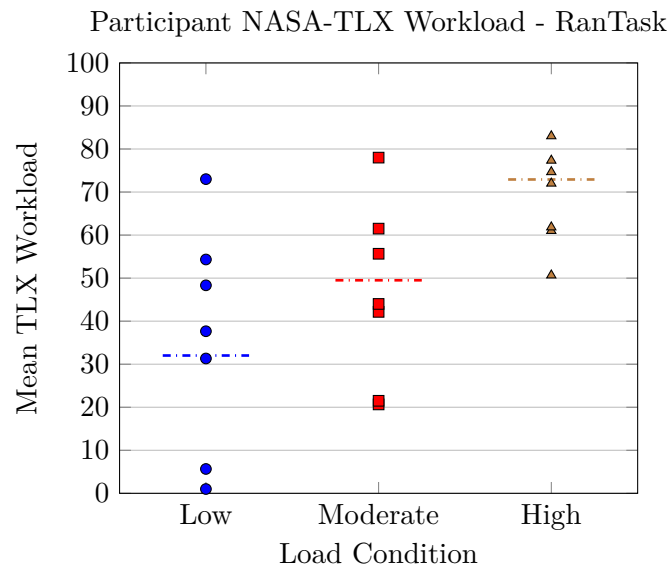


Figure 9

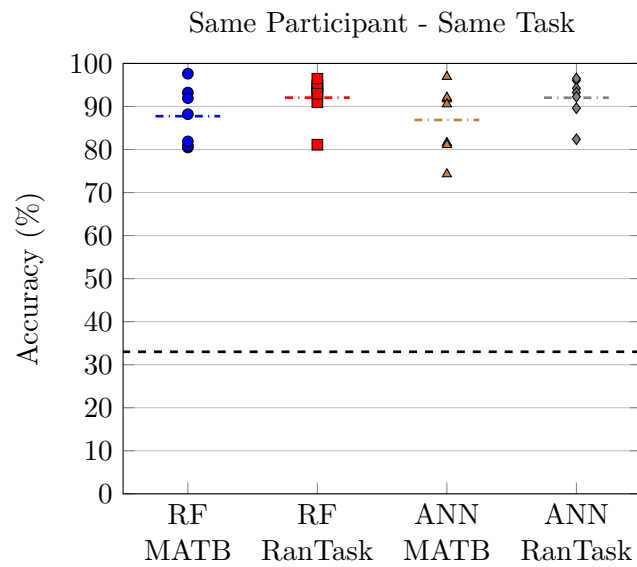


Figure 10

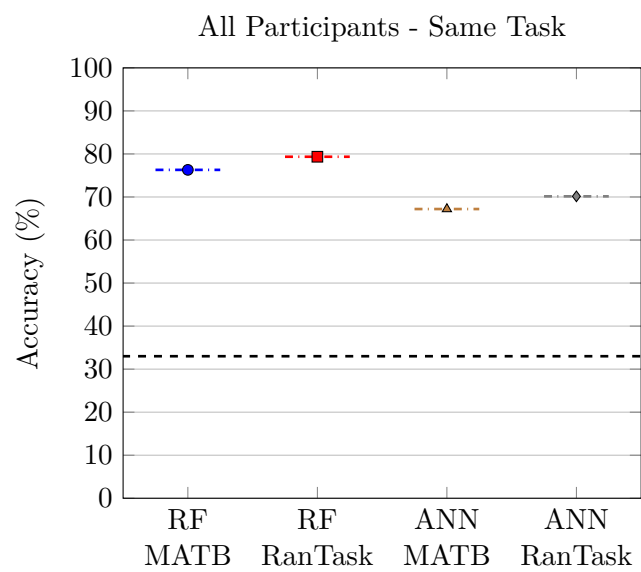


Figure 11

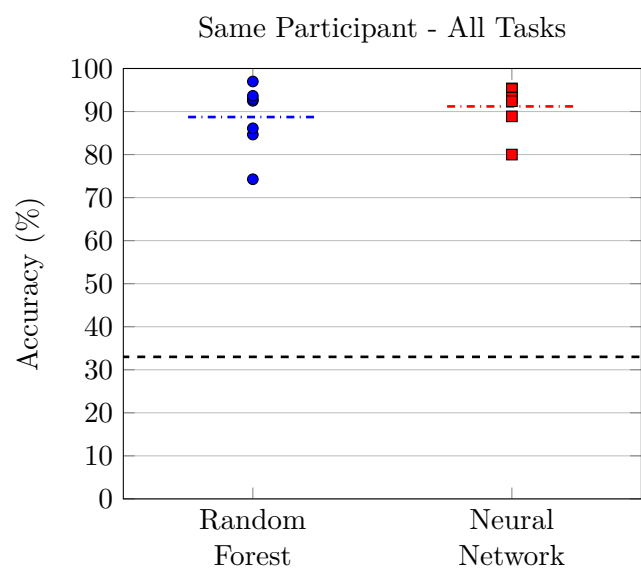


Figure 12

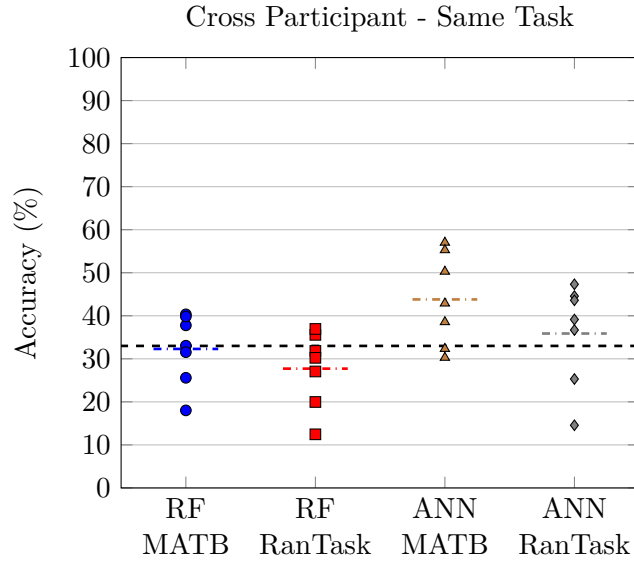


Figure 13

7.1 Participant Study

7.1.1 Performance

7.1.2 NASA Task Load Index (TLX)

7.2 Classifier Models

7.2.1 Same Participant - Same Task (SP-ST)

7.2.2 All Participants - Same Task (AP-ST)

7.2.3 Same Participant - All Tasks (SP-AT)

7.2.4 All Participants - All Tasks (AP-AT)

7.2.5 Cross Participant - Same Task (CP-ST)

7.2.6 Same Participant - Cross Task (SP-CT)

7.2.7 All Participant - Cross Task (AP-CT)

7.2.8 Cross Participant - All Tasks (CP-AT)

7.2.9 Cross Participant - Cross Task (CP-CT)

8 Discussion

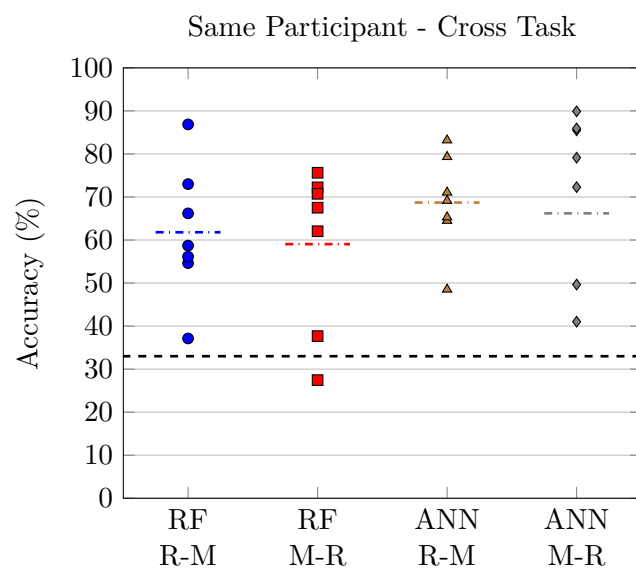


Figure 14

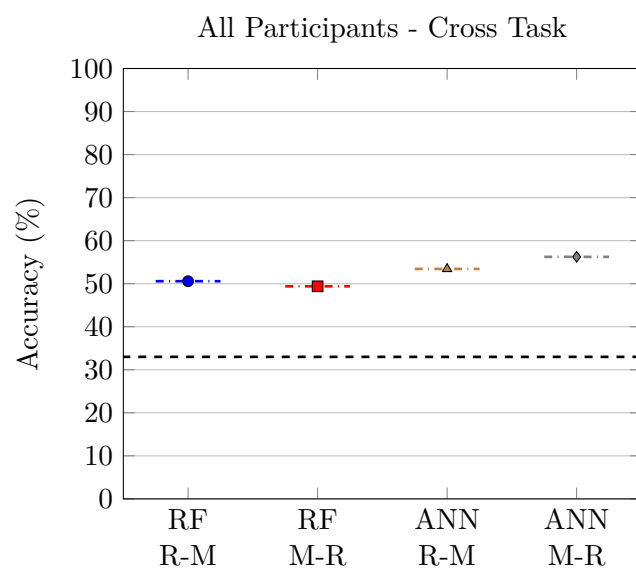


Figure 15

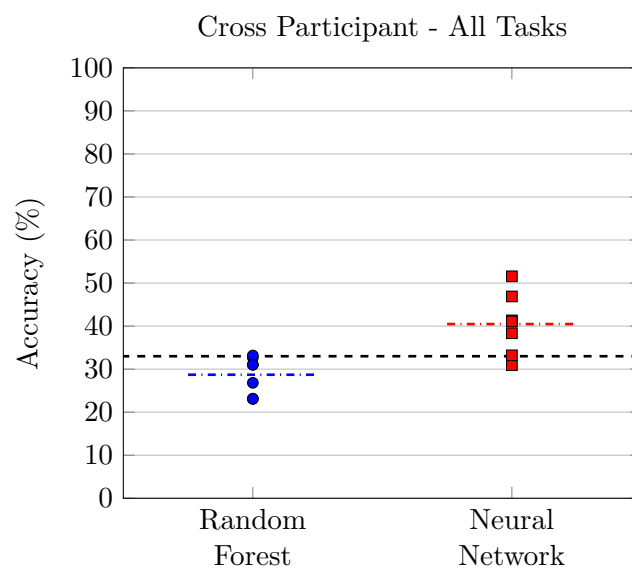


Figure 16

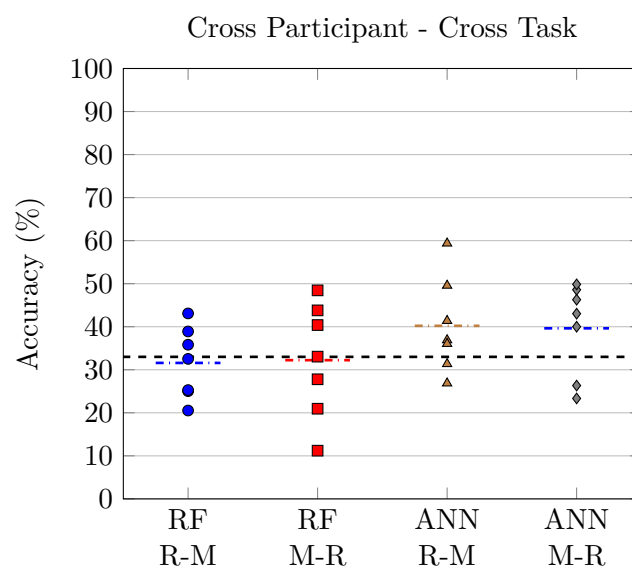


Figure 17

Table 4: Mean Performance Difference Between Conditions

Condition Change	MATB	RanTask
Low to Moderate	7.08 (5.14)	7.65 (2.82)
Moderate to High	3.22 (1.78)	22.27 (15.52)

Table 5: Mean NASA-TLX Difference Between Conditions

Condition Change	MATB	RanTask
Low to Moderate	17.47 (8.98)	13.78 (8.38)
Moderate to High	23.48 (10.73)	24.14 (16.80)

8.1 Participant Study

8.2 Classifier Models

8.3 Potential Sources of Error

9 Future Work & Conclusion

b

Table 6: Mean Model Accuracy of Random Forest(RF) and Neural Network(ANN)

	SP-ST	AP-ST	SP-AT	AP-AT	CP-ST	SP-CT	AP-CT	CP-AT	CP-CT
RF	89.89 (6.20)	77.83 (2.16)	88.72 (7.70)	77.03 ()	30.01 (8.49)	60.42 (16.69)	49.98 (0.84)	28.70 (4.33)	31.91 (10.61)
ANN	89.45 (6.90)	68.68 (2.09)	91.19 (5.38)	66.57 ()	39.84 (11.61)	67.46 (13.75)	54.83 (1.98)	40.47 (7.23)	39.92 (10.46)

References

- [1] N. Meshkati, “Toward Development of a Cohesive Model of Workload,” in *Advances in Psychology* (P. A. H. a. N. Meshkati, ed.), vol. 52 of *Human Mental Workload*, pp. 305–314, North-Holland, 1988.
- [2] G. R. J. Hockey, “Operator functional state: the prediction of breakdown in human performance,” in *Measuring the Mind: Speed, control, and age* (J. Duncan, L. Phillips, and P. McLeod, eds.), pp. 373–394, Oxford University Press, July 2005.
- [3] R. Wang, J. Zhang, Y. Zhang, and X. Wang, “Assessment of human operator functional state using a novel differential evolution optimization based adaptive fuzzy model,” *Biomedical Signal Processing and Control*, vol. 7, pp. 490–498, Sept. 2012.
- [4] J.-H. Zhang, X.-D. Peng, H. Liu, J. Raisch, and R.-B. Wang, “Classifying human operator functional state based on electrophysiological and performance measures and fuzzy clustering method,” *Cognitive Neurodynamics*, vol. 7, pp. 477–494, Dec. 2013.
- [5] G. F. Wilson and C. A. Russell, “Real-Time Assessment of Mental Workload Using Psychophysiological Measures and Artificial Neural Networks,” *Human Factors*, vol. 45, no. 4, pp. 635–643, 2003.
- [6] S. Yang and J. Zhang, “An adaptive human-machine control system based on multiple fuzzy predictive models of operator functional state,” *Biomedical Signal Processing and Control*, vol. 8, pp. 302–310, May 2013.
- [7] G. F. Wilson, “An Analysis of Mental Workload in Pilots During Flight Using Multiple Psychophysiological Measures,” *The International Journal of Aviation Psychology*, vol. 12, pp. 3–18, Jan. 2002.
- [8] G. F. Wilson and C. A. Russell, “Operator Functional State Classification Using Multiple Psychophysiological Features in an Air Traffic Control Task,” *Human Factors*, vol. 45, no. 3, pp. 381–389, 2003.
- [9] R. J. Arnegard and J. R. Comstock, “The multi-attribute task battery for human operator workload and strategic behavior research,” Tech. Rep. NASA

- Tech. Memorandum No. 104174, National Aeronautics and Space Administration, Langley Research Center, Hampton, VA: NASA, 1992.
- [10] B. Lorenz, “Detection and prediction of an automation-induced state of impaired operator competence,” in *Proceedings of NATO ARW on Operator Functional State*, (Il Ciocco), 2002.
 - [11] Y. Ke, H. Qi, F. He, S. Liu, X. Zhao, P. Zhou, L. Zhang, and D. Ming, “An EEG-based mental workload estimator trained on working memory task can work well under simulated multi-attribute task,” *Frontiers in Human Neuroscience*, vol. 8, Sept. 2014.
 - [12] Z. Wang, R. M. Hope, Z. Wang, Q. Ji, and W. D. Gray, “Cross-subject workload classification with a hierarchical Bayes model,” *NeuroImage*, vol. 59, pp. 64–69, Jan. 2012.
 - [13] J. C. Christensen, J. R. Estepp, G. F. Wilson, and C. A. Russell, “The effects of day-to-day variability of physiological data on operator functional state classification,” *NeuroImage*, vol. 59, pp. 57–63, Jan. 2012.
 - [14] E. A. Byrne and R. Parasuraman, “Psychophysiology and adaptive automation,” *Biological Psychology*, vol. 42, pp. 249–268, Feb. 1996.
 - [15] R. Parasuraman, R. Molloy, and I. L. Singh, “Performance Consequences of Automation-Induced ‘Complacency’,” *The International Journal of Aviation Psychology*, vol. 3, pp. 1–23, Jan. 1993.
 - [16] R. Parasuraman and V. Riley, “Humans and automation: Use, misuse, disuse, abuse,” *Human Factors*, vol. 39, p. 230, June 1997.
 - [17] G. F. Wilson, “Operator functional state assessment for adaptive automation implementation,” vol. 5797, pp. 100–104, 2005.
 - [18] B. M. Huey, C. D. Wickens, and National Research Council (U.S.), eds., *Workload transition: implications for individual and team performance*. Washington, DC: National Academy Press, 1993.
 - [19] G. F. Wilson, C. Russell, J. Monnin, J. Estepp, and J. Christensen, “How Does Day-to-Day Variability in Psychophysiological Data Affect Classifier Accuracy?,” *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, vol. 54, no. 3, pp. 264–268, 2010.

- [20] C.-H. Ting, M. Mahfouf, A. Nassef, D. Linkens, G. Panoutsos, P. Nickel, A. Roberts, and G. Hockey, “Real-Time Adaptive Automation System Based on Identification of Operator Functional State in Simulated Process Control Operations,” *IEEE Transactions on Systems, Man and Cybernetics, Part A: Systems and Humans*, vol. 40, pp. 251–262, Mar. 2010.
- [21] Z. Yin and J. Zhang, “Operator functional state classification using least-square support vector machine based recursive feature elimination technique,” *Computer Methods and Programs in Biomedicine*, vol. 113, pp. 101–115, Jan. 2014.
- [22] M. E. Smith, A. Gevins, H. Brown, A. Karnik, and R. Du, “Monitoring task loading with multivariate EEG measures during complex forms of human-computer interaction,” *Human Factors*, vol. 43, no. 3, p. 366, 2001.
- [23] J. Miller, K. D. Schmidt, J. R. Estep, M. Bowers, and I. Davis, “An Updated Version of the U.S. Air Force Multi-Attribute Task Battery (AF-MATB),” Tech. Rep. AFRL-RH-WP-SR-2014-0001, Air Force Research Lab Wright-Patterson AFB OH Human Performance Wing (711th) - Human Effectiveness Directorate, Aug. 2014.
- [24] J. R. Estep and J. C. Christensen, “Electrode replacement does not affect classification accuracy in dual-session use of a passive brain-computer interface for assessing cognitive workload,” *Neuroprosthetics*, vol. 9, p. 54, 2015.
- [25] E. M. Rantanen and J. H. Goldberg, “The effect of mental workload on the visual field size and shape,” *Ergonomics*, vol. 42, pp. 816–834, June 1999.
- [26] S. G. Hart, “NASA Task Load Index (TLX). Volume 1.0; Paper and Pencil Package,” tech. rep., Jan. 1986.
- [27] J. Sweller, P. Ayres, and S. Kalyuga, “Measuring Cognitive Load,” in *Cognitive Load Theory*, no. 1 in Explorations in the Learning Sciences, Instructional Systems and Performance Technologies, pp. 71–85, Springer New York, 2011.
- [28] H. Jasper, “Report of the committee on methods of clinical examination in electroencephalography,” *Electroencephalography and Clinical Neurophysiology*, vol. 10, pp. 370–375, May 1958.

- [29] A. Knoll, Y. Wang, F. Chen, J. Xu, N. Ruiz, J. Epps, and P. Zarjam, “Measuring Cognitive Workload with Low-Cost Electroencephalograph,” in *Human-Computer Interaction – INTERACT 2011* (P. Campos, N. Graham, J. Jorge, N. Nunes, P. Palanque, and M. Winckler, eds.), no. 6949 in Lecture Notes in Computer Science, pp. 568–571, Springer Berlin Heidelberg, 2011.
- [30] A. Delorme and S. Makeig, “EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis,” *Journal of Neuroscience Methods*, vol. 134, pp. 9–21, Mar. 2004.
- [31] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, “Scikit-learn: Machine Learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [32] S. Nissen, “Implementation of a Fast Artificial Neural Network Library (fann),” tech. rep., Department of Computer Science University of Copenhagen (DIKU), 2003. <http://fann.sf.net>.
- [33] M. Riedmiller and H. Braun, “A direct adaptive method for faster backpropagation learning: The RPROP algorithm,” in *Neural Networks, 1993., IEEE International Conference on*, pp. 586–591, IEEE, 1993.

Appendices

A NASA Task Load Index Records

A.1 Multi-Attribute Task Battery (MATB)

A.2 RanTask

B Performance Data

B.1 Multi-Attribute Task Battery (MATB)

B.2 RanTask

C Model Performance

C.1 Artificial Neural Network

C.2 Random Forest