

Towards Designing a Socially Assistive Robot for Adaptive and Personalized Cognitive Training

Konstantinos Tsiakas
CSE Department
University of Texas, Arlington
konstantinos.tsiakas
@mavs.uta.edu

Michalis Papakostas
CSE Department
University of Texas, Arlington
michalis.papakostas
@mavs.uta.edu

Cheryl Abellanoza
Department of Psychology
University of Texas, Arlington
cheryl.abellanoza
@mavs.uta.edu

Tasnim Makada
CSE Department
University of Texas, Arlington
tasniminayat.makada
@mavs.uta.edu

Maher Abujelala
CSE Department
University of Texas, Arlington
maher.abujelala
@mavs.uta.edu

Fillia Makedon
CSE Department
University of Texas, Arlington
makedon@uta.edu

ABSTRACT

There is a growing body of multidisciplinary research on how robotic systems can be deployed in education and training by providing personalized tutoring session to the user. Socially Assistive Robotics (SAR) is an efficient tool for educational and health-care purposes. In this work, we present our SAR system for personalized and adaptive cognitive training. More specifically, we present the sequence learning task that provides measures for executive function assessment, which may indicate learning or even behavior disabilities in children. This work outlines the designing and evaluation process of such a system, including data collection and analysis. The long-term goal of this research is to develop interactive machine learning methods towards the design of an adaptive SAR system that provides a personalized training session by adjusting the session parameters and the robot's behavior to maximize user engagement and performance.

1. INTRODUCTION

Socially Assistive Robotics (SAR) is an area that studies how robots can be deployed to assist users through social interaction, as they perform a cognitive or physical task [7]. The goal of such robotic agents is to build an effective interaction with the user, so as to enhance their task engagement and thus the effectiveness of the training session. Such agents can be deployed to various tasks as cognitive or physical training [5, 12], language learning [9], rehabilitation exercises [13] and others.

As technology advances and becomes more affordable, SAR systems can be considered to be an efficient tool for educational and tutoring purposes. A key feature of SAR systems is their ability to provide personalized interaction to the user. Personalization is essential for an effective train-

ing or tutoring session, since it can enhance the effectiveness of the session, maximizing user's training and learning potentials [3]. Personalization can be achieved among different dimensions. In this work, we focus on two dimensions; task parameters and robot behavior. More specifically, the system can adjust the difficulty of the task in order to provide a training session that fits user's abilities and skills, resulting in an "optimally challenging activity" [4]. On the other hand, the instructor (a social robot) can impact the user's intrinsic motivation and task engagement through verbal feedback during the training/tutoring session [6].

It is important that such systems monitor the user's attention (engagement, concentration) and adjust both the task parameters and their behavior in order to increase user engagement and compliance. However, quantifying engagement is not trivial, since it depends and overlaps with several user states as interest, sustained attention, immersion and (attentional and emotional) involvement [16]. Recently, Brain Computer Interfaces (BCI) have been used towards this purpose [8]. More specifically, we propose the use of *passive* BCI to measure and utilize user engagement and concentration during a cognitive task, to improve user experience and system efficiency. Such sensors can be used in order to maximize the training session effectiveness, by personalizing the session. We propose to use the Muse EEG headset¹. Muse has been used to measure concentration and relaxation [11], task enjoyment [1], adaptive storytelling [18], pain detection through self-calibrating protocols [10] and others.

This paper outlines the experimental design and process for developing a SAR system for personalized and adaptive cognitive training. We present our use case, the sequence learning task (Section 3) and the experimental procedure for the data collection (Section 4). The main contribution of this work is a publicly available dataset along with a set of data analysis tools using machine learning and statistical analysis (Section 5). Our goal through this extensive analysis is to identify links between user engagement, user performance, task difficulty, and robot behavior. We discuss our findings towards the implementation of an autonomous adaptive SAR system, focusing on the computational as-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](http://permissions.acm.org).

HRI '17 March 6–9, 2017, Vienna, Austria

© 2017 ACM. ISBN XXX-XXXX-XX-XXX/XX/XX.

DOI: XX.XXX/XXX_X

¹<http://www.choosemuse.com/>

pects and methods for real-time adaptation and personalization, using Interactive Reinforcement Learning (Section 6).

2. RELATED WORK

SAR systems have been designed for providing assistance to users during a physical or cognitive task. The target population varies from children to elderly for several applications as language tutoring, post-stroke rehabilitation and personalized education for children with ASD.

In [9], the authors proposed a social robotic platform that provided personalized language learning sessions to preschool children, adapting its affective behavior to maximize child’s engagement and valence. In [6], the authors presented a social robot designed to engage elderly users in physical exercise, through motivational behavior. In an educational setting, a SAR system has been proposed for personalized number concept learning. In this work [3], the authors conducted an initial data collection from users in order to extract essential information towards the definition of personalized user models for an adaptive SAR tutoring system. Moreover, SAR systems have been deployed to enhance reading engagement through the physiological reading method, as an educational tool [2]. Our work describes the steps towards the design, implementation and evaluation of a personalized SAR system for the sequence learning task to enhance working memory skills, maximizing task engagement and performance.

3. SEQUENCE LEARNING TASK

For our experimental setup, we deploy the NAO² robot as a socially assistive robot that instructs, monitors and evaluates user’s performance during a cognitive task. We present the Sequence Learning task; a working memory task that evaluates the ability of a human to remember and repeat a sequence of items (e.g., letters, number, actions). Sequencing is the ability to arrange language, thoughts, information and actions in an effective order. It has been shown that many children with learning and attention issues have trouble with sequencing [17]. Towards this direction, we present a cognitive training task for sequence learning.

During the training task, the user has three buttons in front of them (“A”, “B”, “C”) and the robot asks the user to repeat a given sequence of these letters. The difficulty of the task is proportional to the sequence length $L = [3, 5, 7, 9]$. Moreover, the robot can provide feedback after the user completes a sequence. We examine the influence of different feedback styles and their relationship with user’s engagement and performance [15]. More specifically, the robot, after each turn, can either provide positive, negative or no feedback.

We have defined three different task modes, based on the way the user needs to reply: (B)uttons, (S)peech and (F)lanker Test:

1. **Buttons Mode:** User has to press the correct sequence of buttons
2. **Speech Mode:** User has to repeat verbally the correct sequence of buttons
3. **Flanker Test Mode:** User has to identify the middle letter of the sequence announced (e.g., in ABBAC)

²<https://www.ald.softbankrobotics.com/en/cool-robots/nao>

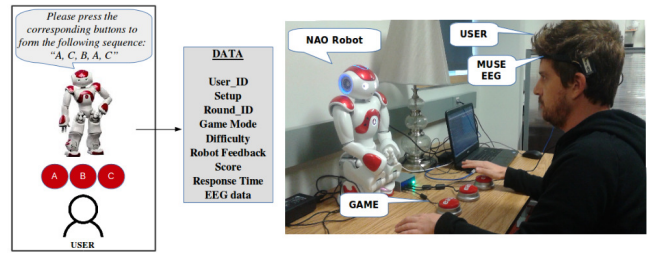


Figure 1: The Sequence Learning training task. The user is asked to repeat a given sequence, by pressing the buttons in the correct order. The NAO robot instructs and monitors the user during the task. Muse EEG signals are used to estimate concentration and engagement levels.

Each task mode requires different user skills and abilities. We plan to examine user reaction and performance under task switching conditions. Task switching is an executive function, and a kind of cognitive flexibility, that involves the ability to shift attention between one task and another [14]. However, the scope of this paper is preliminary data collection and analysis to identify different user skills and preferences for different task modes and difficulties, as well as for different feedback types, towards the definition of personalized models for the sequence learning task.

4. EXPERIMENTAL PROCEDURE

As a first step, in order to collect and analyze interaction data during the sequence learning task, we conducted a user study. During the experiment, the NAO robot instructs and monitors the user during the training session, collecting interaction data as described in Section 5.

At the beginning of the experiment, each user is asked to take place in front of the robot and wear the Muse EEG sensor. After the task administrator ensures the correct placement of the Muse sensor, the NAO robot greets the user and describes to them the sequence learning task and the different modes³. After the introduction, the robot asks the user if the process was clear to them. During the task, the robot performs an action (A0 – A5) that defines the difficulty level or the feedback type of the next turn, as shown in Figure 2.

Difficulty (Length L)	A0	Easy	L = 3	
	A1	Medium	L = 5	
	A2	Normal	L = 7	
	A3	Hard	L = 9	
Feedback	A4	Positive	“Very good!! Keep going” or “Oh, you missed but don’t worry!!”	Same difficulty
	A5	Negative	“Maybe that was too easy” or “Aren’t you paying any attention?”	

Figure 2: Robot Actions. For each action, the robot announces a sequence of the corresponding difficulty.

We defined two different experimental designs, in terms of how the task difficulty changes. Each user performs the task for both designs: *blocked* and *mixed*. We followed these two different designs in order to capture user data under different variations of difficulty. The robot action sequence

³<https://www.youtube.com/watch?v=giTwZGaBUtE>

(including feedback actions) is predefined and same for all users.

1. **Blocked Design:** In this design, the difficulty levels are gradually increasing from the lowest ($L = 3$) to the highest ($L = 9$) difficulty, for each task mode. Each user has to perform the task for 9 rounds for each task mode, resulting in $9 \cdot 3 = 27$ rounds.
2. **Mixed Design:** In this design, the difficulty levels are mixed and change during the task. Each user has to perform the task for 12 rounds for each task mode, resulting in $16 \cdot 3 = 48$ rounds.

5. DATA COLLECTION AND ANALYSIS

In our user study, we recruited 15 participants (10 males, 5 females) between the age of 24 and 37. Each experimental session lasted for about 30 minutes, including the completion of the consent form (IRB 2017-0375), the task introduction, the training task and a post-experiment user survey. The data collected during each experiment are depicted in the database schema in Figure 3.

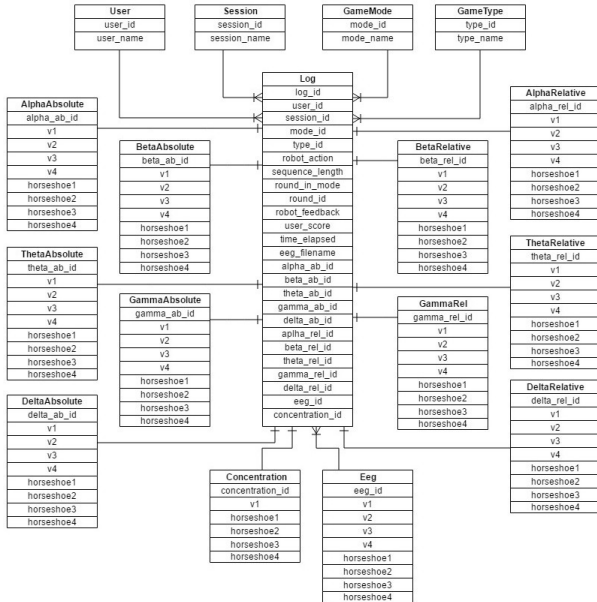


Figure 3: The database schema. The database stores the collected data for each experiment. The DB structure allows for data mining through querying.

As shown in the schema, the system keeps the following information for each task session: user ID, task design, task mode, round ID, robot action, sequence length, completion time, user performance and EEG raw data. More specifically, the system stores all absolute and relative EEG band and concentration values, as well as the headband connection status indicators (as derived by MUSE). The dataset and the analysis can be found online⁴.

In order to get an insight of how we can leverage the collected data to develop personalized models for the sequence learning task, we provide an extensive and multi-aspect analysis, including machine learning and statistical analysis. Moreover, a user survey was conducted to evaluate the task itself and the user interaction with the robot.

⁴<http://heracleia.uta.edu/%7Etsiakas/SLdataset.html>

5.1 User Survey Analysis

The post-experiment survey discusses the participant’s experience during the experiment, including their overall interaction with the robot, as well as self-reported enjoyment, engagement, difficulty and concentration needed during the experiment.

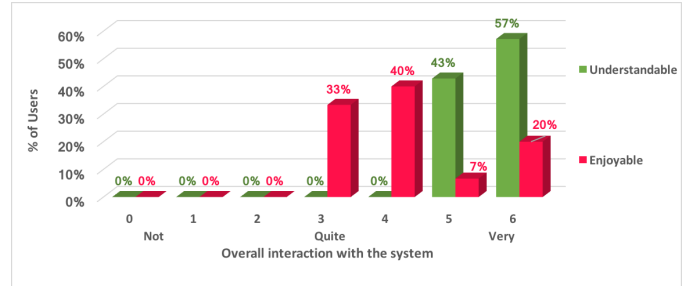


Figure 4: Users Evaluation for their interaction with the system.

Figure 4 illustrates that all participants find the interaction with the system at least quiet enjoyable and the majority finds the system very easy to understand. Also, the participants reported the difficulty, engagement, and their concentration at each level of the task as can be seen in Figure 5. It can be observed that the users’ ratings shift from (rating of 1) not difficult, not engaging and no concentration needed to (rating of 6) very difficult, very engaging and high concentration needed as the difficulty level increases.

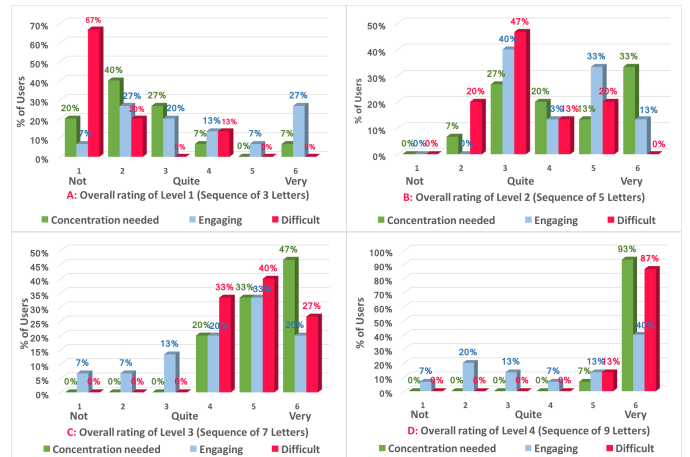


Figure 5: User rating on Task Difficulty, Task Engagement and Concentration Needed for the different difficulty levels

In particular, in Figure 5A most of ratings are below “Quite”, while in Figure 5B most of the ratings are between “Not” and “Very”. In Figure 5C most of the ratings are between “Quite” and “Very”, and in Figure 5D most of the ratings are at “Very”. The results of this figure indicate that Level 1 might be less effective in keeping most of the participants engaged and concentrated, and Level 4 might be too challenging for most participants. However, Level 2 and Level 3 might be more effective in keeping the participants engaged and concentrated, and in keeping their performance high as illustrated in Figure 9.

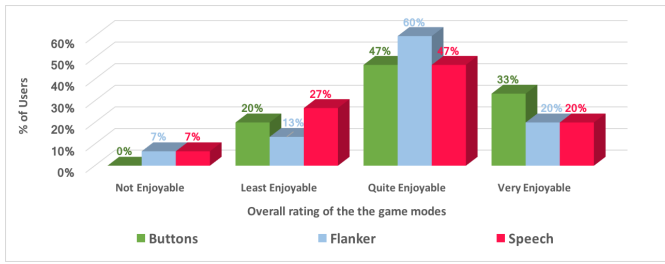


Figure 6: Overall enjoyment rating for the different task modes (Buttons, Flanker, and Speech)

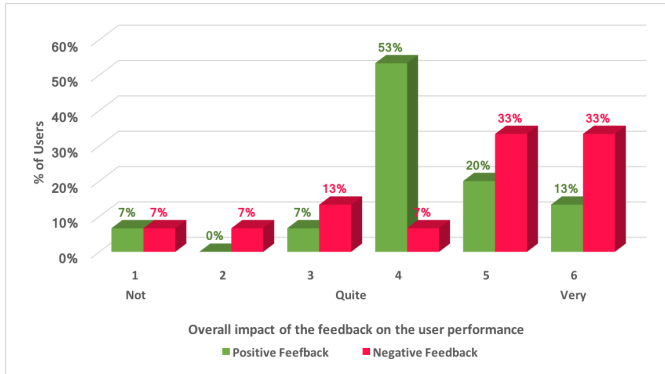


Figure 7: User Evaluation on the impact of robot's positive/negative feedback on their interaction with the task and the robot

With regard to different task modes (Buttons, Flanker, Speech), the majority of participants (80%) reported that they at least quite enjoyed the Buttons (47%+33%) and the Flanker (60%+ 20%) modes compared to 67% [47%+20%] who at least quite enjoyed the Speed mode (see Figure 6). This finding does not match with the users performance, since the majority of the participants performed much better in the Flanker mode compared to the other two modes (see Figure 8). Similarly, the participants reported that both of the positive and negative feedback have at least quite affected their performance (see Figure 7); however, the feedback did not affect their performance as discussed in section 5.2.

5.2 Statistical Analysis

In this section, we describe the statistical analysis on the interaction data. We are interested in investigating how different features correlate to each other (e.g., task difficulty, robot feedback, task engagement, user performance, completion time, etc.). We provide an initial analysis using histograms, to get an insight of the data distribution, supported by an additional statistical analysis to identify possible patterns towards defining personalized user models.

Significant effects were found in accuracy, completion time, and concentration with respect to Mode ($ps < .001$), Difficulty ($ps < .001$), and Mode \times Difficulty ($ps < .001$). Generally, users were most accurate, fastest, and concentrated least on the Flanker test. Also, users were most accurate, fastest, and concentrated least on the easiest level of difficulty (3 units in a sequence). These patterns are further supported by the histograms in Figures 8, 9 and

10. For instance, the somewhat platykurtic distribution for the Flanker task in Figure 8 also shows that more users had higher accuracy; because the Flanker task only requires users to remember one letter, this may have been the easiest task. Also, there is a difference of distributions of correct answers based on difficulty, as seen in Figure 9, with the most amount of accurate responses in the 3-unit sequence.

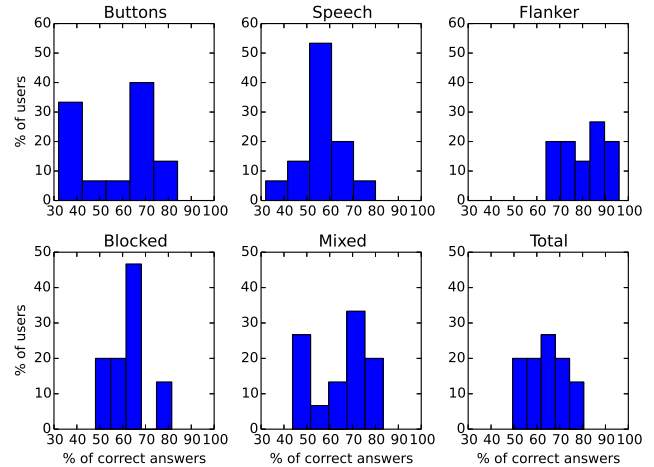


Figure 8: Percentage of correct answers per design and mode.

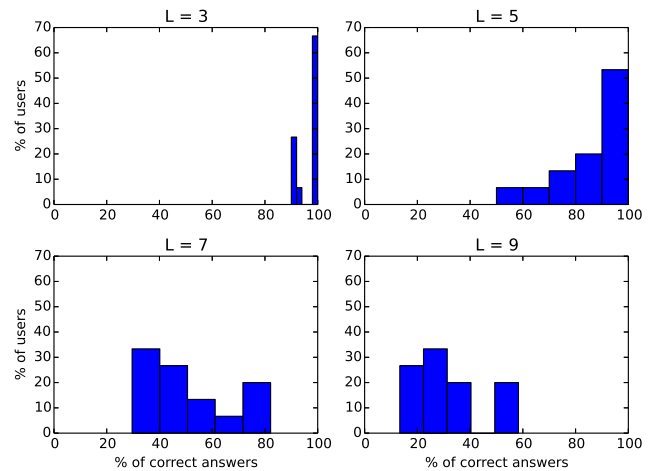


Figure 9: Percentage of correct answers per difficulty.

To determine whether feedback (positive, negative, none) influenced overall performance, a one-way between-subjects Analyses of Variance (ANOVA) were run. To address violations of assumptions of normality and homogeneity of variance (assessed using Shapiro-Wilks and Levene's tests), Welch's F tests of equality of means were calculated, corrected values were reported, and tasks-Howell post-hoc tests were used. Users were more accurate after receiving no feedback, $F(2,387.78) = 8.48, p < .001$. Users also had faster reaction times after receiving no feedback, $F(2,362.48) = 10.39, p < .001$. No significant differences were found be-

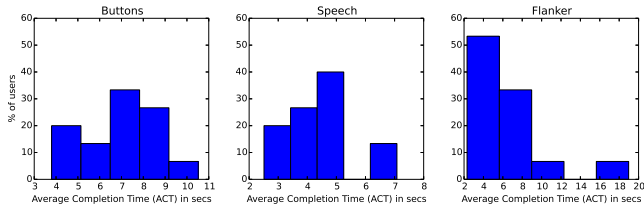


Figure 10: Average Completion Time (ACT) for the different task modes.

tween positive and negative feedback. Means and standard errors from this analysis are reported below in Table 1 (significant difference is noted by asterisk). Figure 11 shows the distribution (histogram) of correct answers per feedback type.

Table 1: User Performance

Feedback Type	Accuracy (%)	Completion Time (ms)
*None	70 \pm 2	5.00 \pm 0.17
Negative	58 \pm 4	6.40 \pm 0.46
Positive	57 \pm 3	6.36 \pm 0.29

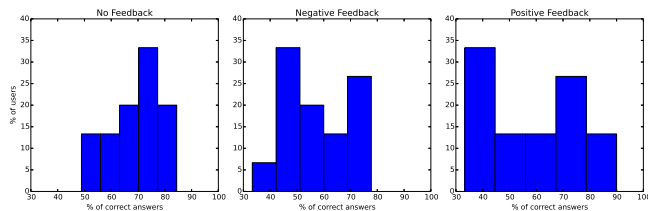


Figure 11: Percentage of correct answers per feedback type.

Though feedback did not have a significant effect on overall performance, the analysis suggests some trends. For instance, the mean accuracy scores for negative feedback trials was higher than for no feedback or positive feedback trials during mid-levels of task difficulty. In 5-unit sequences, the average accuracies were: no feedback = 86.4%, positive feedback = 81%, and negative feedback = 87.9%. In 7-unit sequences, the average accuracies were: no feedback = 52.2%, positive feedback = 49.4%, and negative feedback = 52.3%. This could suggest that users did not benefit from feedback at easier levels because the task may have been too easy, or because they perceived the feedback to be unnecessary. Similarly, users may not have benefited from feedback at harder levels (>7 units) because the task was too difficult, or because they were focusing more on the required task. The statistical analysis on our collected data denotes the need of personalized models for personalized training sessions.

5.3 Machine Learning Analysis

The purpose of the following ML analysis is to verify the validity of the proposed dataset and the capability of the captured data to sufficiently model patterns of user behavior for personalized training sessions. In particular, we perform two different types of experiments. In the first experiment, we train a Random Forest (RF) classifier to *predict user performance* for a single task round. In the second experiment,

we train a Linear-Regression (LR) model to *estimate task completion time* for a given sequence. For both experiments, we use the collected data for each round (Fig.3), including the average band and concentration values (EEG).

For the first experiment, we chose RFs against other traditional methods due to their consistency, when tested on the proposed dataset (Table-2). However, we provide online complementary ML approaches that can be used by other researchers for modeling and comparison, as SVMs, ANNs and other classifiers. The RF classifier, consists of 100 randomly designed estimators. We evaluate each model using 5-Fold cross validation, each time using 80% of the available data for training and the rest 20% for testing. Each classifier is trained two times, one without including completion-time as a feature and one where completion-time is included. In total, we trained and tested twelve classifiers based on six different sub-datasets. The first dataset consists of all the available data, two sub-datasets were extracted including only data associated with a specific task design and the last three sub-datasets were created based only on data related to a specific task mode. Classification results are shown in Table-2

Table 2: Classification Accuracy

	RF Classifier	
	With Time	Without Time
All Data	0.81 \pm 0.09	0.77 \pm 0.07
Blocked Design	0.86 \pm 0.08	0.81 \pm 0.11
Mixed Design	0.75 \pm 0.11	0.74 \pm 0.11
Flanker Mode	0.79 \pm 0.15	0.78 \pm 0.13
Buttons Mode	0.86 \pm 0.07	0.77 \pm 0.05
Speech Mode	0.78 \pm 0.12	0.75 \pm 0.08

For the second experiment, we deployed a Linear Regression model (Table-3) using the LASSO method (with $\alpha = 0.1$). As before, we trained twelve LR models, two for each of the aforementioned sub-datasets. The first model takes the user performance into account, whereas the second does not. We evaluate each LR model using the Mean Squared Error and the Variance Score, evaluation metrics. As it can be derived by equations 1 and 2, a good regression model is specified by a MSE close to 0 and a VS close to 1.

$$MSE(y, \hat{y}) = \frac{1}{n_{samples}} \sum_{i=0}^{n_{samples}-1} (y_i - \hat{y}_i)^2 \quad (1)$$

$$VS(y, \hat{y}) = 1 - \frac{Var\{y - \hat{y}\}}{Var\{y\}} \quad (2)$$

These preliminary results indicate that the selected features can be sufficiently used to predict user performance and completion time. During both experiments, we trained our models using features from all users, resulting in a generalized model. We argue that training a model using data from a specific individual or cluster of users can result to a better accuracy. For example, we repeated the classification experiments using the data from the cluster of users with total accuracy above 70%, in both designs. Using only user performance as a dimension for user clustering, we observe that the user score classification accuracy increases to 90% from 81% (Table 2). We argue that personalized models

Table 3: Linear Regression Results

	LASSO LR			
	With Score		Without Score	
	MSE	VS	MSE	VS
All Data	13.01	0.45	13.84	0.41
Blocked Design	8.86	0.61	10.69	0.53
Mixed Design	14.49	0.33	16	0.26
Flanker Mode	5.91	0.07	5.96	0.06
Buttons Mode	7.42	0.61	9.88	0.49
Speech Mode	10.04	0.45	10.72	0.41

across multiple dimensions would increase the classification and regression accuracy.

6. DISCUSSION AND FUTURE WORK

In this work, we outlined the implementation procedure of a SAR system employed as a tutor for the sequence learning task. We provided a detailed description of the experimental and data collection process, as well as a set of multi-aspect data analysis. The preliminary results showed that task parameters and robot feedback may have an effect on the training session and the interaction with the robot, that varies from user to user. Our ongoing work includes the training of a Reinforcement Learning agent using the collected data, learning different user-specific policies, as described in [19], applying our proposed Interactive Learning and Adaptation framework for real-time adaptation. Further studies will be conducted to evaluate and refine the proposed framework.

7. ACKNOWLEDGMENTS

This material is based upon work supported by the National Science Foundation under Grant CHS 1565328. The authors would also like to thank Chris Collander for helping with the hardware of the experimental setup.

8. REFERENCES

- [1] M. Abujelala, C. Abellanoza, A. Sharma, and F. Makedon. Brain-ee: Brain enjoyment evaluation using commercial eeg headband. In *Proceedings of the 9th ACM International Conference on Pervasive Technologies Related to Assistive Environments*, page 33. ACM, 2016.
- [2] M. Andujar and J. E. Gilbert. Let’s learn!: enhancing user’s engagement levels through passive brain-computer interfaces. In *CHI’13 Extended Abstracts on Human Factors in Computing Systems*, pages 703–708. ACM, 2013.
- [3] C. Clabaugh, G. Ragusa, F. Sha, and M. Matarić. Designing a socially assistive robot for personalized number concepts learning in preschool children. In *IEEE International Conference on Development and Learning and Epigenetic Robotics*, pages 314–319. IEEE, 2015.
- [4] M. Csikszentmihalyi. *Beyond boredom and anxiety*. Jossey-Bass, 2000.
- [5] J. Fasola and M. J. Matarić. Robot motivator: Increasing user enjoyment and performance on a physical/cognitive task. In *2010 IEEE 9th International Conference on Development and Learning*, pages 274–279. IEEE, 2010.
- [6] J. Fasola and M. J. Mataric. Using socially assistive human–robot interaction to motivate physical exercise for older adults. *Proceedings of the IEEE*, 100(8):2512–2526, 2012.
- [7] D. Feil-Seifer and M. J. Mataric. Defining socially assistive robotics. In *9th International Conference on Rehabilitation Robotics, 2005. ICORR 2005.*, pages 465–468. IEEE, 2005.
- [8] L. George and A. Lécuyer. An overview of research on” passive” brain-computer interfaces for implicit human-computer interaction. In *International Conference on Applied Bionics and Biomechanics ICABB 2010-Workshop W1” Brain-Computer Interfacing and Virtual Reality”*, 2010.
- [9] G. Gordon, S. Spaulding, J. K. Westlund, J. J. Lee, L. Plummer, M. Martinez, M. Das, and C. Breazeal. Affective personalization of a social robot tutor for children’s second language skills. In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- [10] T. Karydis, F. Aguiar, S. L. Foster, and A. Mershin. Performance characterization of self-calibrating protocols for wearable eeg applications. In *Proceedings of the 8th ACM International Conference on Pervasive Technologies Related to Assistive Environments*, page 38. ACM, 2015.
- [11] Z. Li, J. Xu, and T. Zhu. Prediction of brain states of concentration and relaxation in real time with portable electroencephalographs. *arXiv preprint arXiv:1509.07642*, 2015.
- [12] M. J. Matarić, J. Eriksson, D. J. Feil-Seifer, and C. J. Winstein. Socially assistive robotics for post-stroke rehabilitation. *Journal of NeuroEngineering and Rehabilitation*, 4(1):1, 2007.
- [13] R. Mead, E. Wade, P. Johnson, A. S. Clair, S. Chen, and M. J. Mataric. An architecture for rehabilitation task practice in socially assistive human-robot interaction. In *19th International Symposium in Robot and Human Interactive Communication*, pages 404–409. IEEE, 2010.
- [14] S. Monsell. Task switching. *Trends in cognitive sciences*, 7(3):134–140, 2003.
- [15] E. Park, K. J. Kim, and A. P. Del Pobil. The effects of a robot instructor’s positive vs. negative feedbacks on attraction and acceptance towards the robot in classroom. In *International Conference on Social Robotics*, pages 135–141. Springer, 2011.
- [16] C. Peters, G. Castellano, and S. de Freitas. An exploration of user engagement in hci. In *Proceedings of the International Workshop on Affective-Aware Virtual Agents and Social Robots*, page 9. ACM, 2009.
- [17] P. Rosen. Trouble with sequencing: What you need to know, 2014.
- [18] D. Szafr and B. Mutlu. Pay attention!: designing adaptive agents that monitor and improve user engagement. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 11–20. ACM, 2012.
- [19] K. Tsiakas, M. Dagioglou, V. Karkaletsis, and F. Makedon. Adaptive robot assisted therapy using interactive reinforcement learning. In *International Conference on Social Robotics*, pages 11–21. Springer, 2016.