Replay with Stochastic Neural Transformation for Online Continual EEG Classification

Tiehang Duan^{1,3}, Zhenyi Wang², Gianfranco Doretto³, Fang Li⁴, Cui Tao⁴, Donald A. Adjeroh³

¹ Meta AI, Bellevue, WA 98005, United States

² University of Maryland, College Park, MD 20742, United States

tiehang.duan@gmail.com, wangzhenyineu@gmail.com, gianfranco.doretto@mail.wvu.edu fang.li@uth.tmc.edu, cui.tao@uth.tmc.edu, donald.adjeroh@mail.wvu.edu

Abstract-Brain computer interface (BCI) systems used for clinical assistance purposes such as wheelchair control require decoding of streaming brain signals i.e. electroencephalography (EEG) signals over a long period of time with subject shift in the middle. Numerous challenges arise during this online continual brain signal decoding process: 1) the EEG decoder needs to deal with streaming EEG signals from sequentially arriving subjects, with no data available beforehand for largescale pretraining; 2) the EEG decoder should avoid catastrophic forgetting on previous subjects after learning on a new subject; 3) the EEG decoder should perform well on noisy signals with high variance across subjects. We proposed a principled replay-based approach for this general decoding scenario, forming a bi-level optimization framework with stochastic neural transformation for dynamic memory evolution, making them representative in feature space and encouraging the model to generalize well. The evolved signal segments are stored and replayed during later decoding stages to achieve optimal model performance on all previous subjects. The stochastic neural transformation performed in inner sup of bi-level optimization significantly enhances the diversity of stored signal segments and improves model robustness during online continual decoding. We perform detailed theoretical analysis on model's generalization ability in addition to the empirical evaluations. We construct multiple new benchmarks to mimic real-world online sequential EEG decoding scenarios with underlying subject shifts. The extensive evaluation of the proposed approach shows it outperforms related strong baselines by a large margin.

I. INTRODUCTION

Development of brain computer interface (BCI) systems for clinical assistance purposes such as robotic wheelchair control or digital interface interaction requires proper decoding of electroencephalography (EEG) signals for a long period of time [II], [2]. Given the non-invasive nature of EEG recording, the signal usually endures significant noise and its patterns differ significantly across different subjects. This poses challenge to the BCI system when it needs to sequentially decode streaming EEG signals from different subjects.

Here, We formulate the problem as that of online continual EEG decoding. The problem setting is characterized by the following challenges: 1) the model decodes streaming EEG signals in an online manner and no data is available for pretraining; 2) the streaming signals are from sequentially arriving

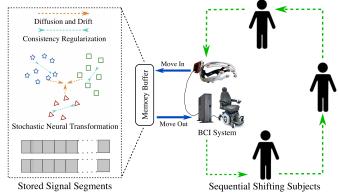


Fig. 1: An overview of the replay-based ReSNT model for online continual decoding of streaming EEG signals over a long period of time with subjects shift. We propose to perform stochastic neural transformation on the stored segments following the bi-level optimization target (eq. 67) to help the model generalize to all previous subjects.

subjects with subject shifts; 3) the EEG decoder should not forget knowledge of previous subjects after learning new ones, aka. catastrophic forgetting, and achieve optimal performance on all subjects after sequential decoding ends^[1].

Deep learning approaches have demonstrated effectiveness for EEG decoding in general and achieved state-of-the-art performance [3][4]. Prior efforts have focused mostly on developing new model architectures [5][3][6] or exploring domain adaptation/transfer learning techniques [7][8] for classic EEG decoding settings e.g. cross-subject EEG decoding, etc. To the best of our knowledge, little exploration has been made for the challenging online continual EEG decoding problem. Recently, continual learning techniques emerged to be promising for mitigation of the problem of catastrophic forgetting [9], [10], with approaches including replay based methods [11], regularization based methods [12] and dynamic expandable network architectures [13]. In this work, we focus on replay-based approaches for continual decoding

¹Classic cross subject EEG decoding models are not readily applicable to this problem setting, as the signal is streaming in an online manner and the data is not jointly available beforehand for large scale pre-training.

³ West Virginia University, Morgantown, WV 26506, United States

⁴ University of Texas Health Science Center at Houston, Houston, TX 77030, United States

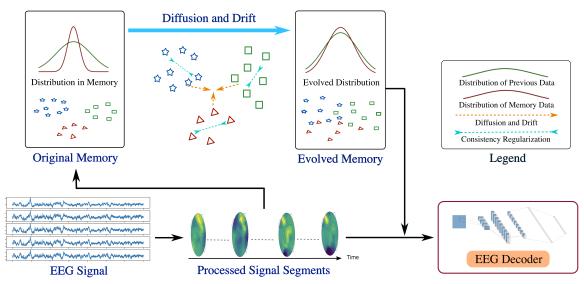


Fig. 2: Illustration on the evolution of signal segments following the stochastic neural transformation. The divergence and drift operations produce robust and diversified features for generalization, and the consistency regularization preserves similarity between transformed signal and original signal.

across sequentially arriving different subjects. The proposed approach is significantly different from existing methods, with its emphasis on generalization improvement and overfitting mitigation during online continual EEG decoding.

Specifically, the proposed replay-based approach stores a small number of representative signal segments from previous subjects for replay purpose. Given that the amount of stored segments are very limited, overfitting could be a major problem and cause the model not generalize well on previous data. As explored in previous work, the diversity of training data is a major factor for model's generalization ability [14]. We form a bi-level optimization framework with stochastic neural transformation on signal segments to keep them diverse and representative for replay purposes. The transformation process consists of two aspects: 1) diffusion and drift operations that evolve the memory data to improve generalization ability. 2) consistency regularization that enforces similarity with original data and produces smooth model response, which allows the model to achieve optimal performance on testing data. The neural transformation process can be depicted as stochastic neural ordinary differential equations, which enable the data to go through an infinite number of stochastic transformations and improve data diversity in the process.

The contribution of this work is summarized as follows:

- We propose a replay-based approach with dynamic memory evolution for online continual decoding of streaming EEG signals over a long period of time with subject shifts.
- The proposed approach performs bi-level optimization with memory evolution depicted as a stochastic process.
 The evolution process is regularized with consistency constraints to preserve similarity with original signal.
- We developed online sequential EEG decoding benchmarks, and performed extensive evaluation on the model performance. The results show the proposed approach outperformed strong baselines by a significant margin.

II. METHOD

We first offer a high-level description of the problem setting and form the overall target function, then provide details of the proposed online continual decoding approach named **Re**play with **Stochastic Neural Transformation** (ReSNT) under the bilevel optimization framework.

A. Preliminaries

For decoding over a long period of streaming EEG signals with subject shift in the middle, decoding performance can quickly deteriorate on previous subjects after learning on later subjects. Here we seek to achieve optimal decoding performance on all previously learnt subjects and mitigate forgetting during online sequential EEG decoding. Note we assume there is no jointly available data from the subjects for pretraining beforehand, which suits to real world scenarios of BCI applications.

We adopt the replay based approach, which keeps a small memory buffer \mathcal{M} to store representative data samples of previous subjects and replay in later decoding stages. With streaming EEG signal from different subjects $\mathcal{S}_1, \mathcal{S}_2, \ldots, \mathcal{S}_J$ sequentially input into the EEG decoder, the memory buffer adopts the reservoir sampling mechanism for sample selection, which has an equal probability of selecting incoming signal segments. It stores the first M samples until it is full, after which for data x_k arriving, it generates a random number i in the range of [1, k]. If i < M, then the i th data in memory will be replaced by x_k .

B. Method Overview

The training target for the replay based approach can be formulated as

$$\min_{\forall \boldsymbol{\theta} \in \boldsymbol{\Theta}} \left[\mathcal{L}(\boldsymbol{x}_i, y_i, \boldsymbol{\theta}) + \mathbb{E}_{(\boldsymbol{x}_m, y_m) \sim \mathcal{M}} \mathcal{L}(\boldsymbol{x}_m, y_m, \boldsymbol{\theta}) \right], \quad (1)$$

where x_i, y_i is the current streaming data and x_m, y_m is the stored data in memory. When we only keep a small amount of data in memory, the conventional replay will be less effective when sequentially learning over a longer period of time, as the model would easily overfit onto the limited memory data and may not successfully generalize to previous subjects. Here we utilize the stochastic neural transformation with consistency constraints to achieve improved diversity in stored signal segments for effective replay.

C. Stochastic Neural Transformation during Replay

The proposed stochastic neural transformation creates diversified data with the diffusion and drift operations that depict the evolution of a stochastic dynamic system. With $\sigma_{\phi}: [0,G] \times \mathbb{R}^d \to \mathbb{R}^{d \times w}$ being the diffusion term and $\mu_{\phi}: [0,G] \times \mathbb{R}^d \to \mathbb{R}^d$ being the drift term, both generated by a network parameterized with ϕ and G is the transformation intensity, the stochastic neural transformation process can be formulated as

$$dX_t = \mu_{\phi}(X_t)dt + \sigma_{\phi}(X_t) \circ dW_t, \tag{2}$$

where X_t is the distribution of transformed data at t, with $t \in [0,G]$. and $X_t \in \mathbb{R}^d$. $W_t \to \mathcal{W}^w$ is a w-dimensional variable following Brownian motion, with the property that $W_{t+r} - W_r$ follows a Gaussian distribution with variance t. The initial state of the neural transformation is the original stored segments. The transformed data can be expressed as

$$X_G = X_0 + \int_0^G \mu_{\phi}(X_t)dt + \int_0^G \sigma_{\phi}(X_t) \circ dW_t$$
 (3)

The \circ operation between $\sigma_{\phi}(X_t)$ and dW_t denotes the Stratonovich integral.

In order for the transformed data to be consistent with the original data, we utilize the Jenson-Shannon divergence regularization to produce smooth model response on the transformed data, which functions between the posterior distribution of transformed data and original data.

$$p_{mean} = (p_{x_0} + p_{x_C})/2$$
 (4)

$$JS(\boldsymbol{x}_G, \boldsymbol{x}_0) = (KL(p_{\boldsymbol{x}_0}||p_{mean}) + KL(p_{\boldsymbol{x}_G}||p_{mean}))/2 \quad (5)$$

where \mathbb{KL} is the KL divergence, p_{x_0} and p_{x_G} are the network output probability of classifying the original and transformed data into different classes.

Parameters of the stochastic neural transformer ϕ is learned in an end-to-end manner together with the model parameters θ with bi-level optimization. Here, the target loss function can be expressed as

$$\min_{\boldsymbol{a}} [\mathcal{L}(\boldsymbol{x}_i, y_i, \boldsymbol{\theta}) + \mathcal{L}(\boldsymbol{x}_m(G), \boldsymbol{y}_m, \boldsymbol{\theta}, \boldsymbol{\phi}_*)]$$
 (6)

s.t.
$$\phi_* = \arg \max_{\phi} [\mathcal{L}(\boldsymbol{x}_m(G), \boldsymbol{y}_m, \boldsymbol{\theta}, \phi) - \mathcal{L}(\boldsymbol{x}_m, \boldsymbol{y}_m, \boldsymbol{\theta}, \phi) - \lambda \mathbb{JS}(\boldsymbol{x}_m(G), \boldsymbol{x}_m)]$$
 (7)

where $x_m(G)$ and x_m are the transformed and original stored segments respectively. The minimization target of eq. [6]

performs optimization on model parameter θ , and the stochastic neural transformation parameter ϕ is optimized with eq. $\boxed{2}$ which diverges the transformed data not to be memorized by model and at the same time being consistent and semantically similar to original data.

To perform proper update on $\sigma_{\phi}(X_t)$ and $\mu_{\phi}(X_t)$ based on the target function specified in eq. $\boxed{1}$ we perform the following: With $A_t = \frac{d\mathcal{L}(X_t)}{dX_t}$, the relationship between A_t and $\sigma_{\phi}(X_t)$, $\mu_{\phi}(X_t)$ can be expressed as

$$dA_t^i = -A_t^j \frac{\partial \mu_\phi^j(X_t)}{\partial X^i} dt - A_t^j \frac{\partial \sigma_\phi^{j,k}(X_t)}{\partial X^i} \circ dW_t^k \tag{8}$$

where i, j, k are the index of matrix A_t and $\sigma_{\phi}(X_t)$. The gradient on ϕ can be obtained accordingly based on eq. \blacksquare Different from conventional data augmentation and evolution approaches which rely on empirical evaluation for the effectiveness to tackle overfitting, the proposed approach offers a principled formulation that enables theoretical analysis on its generalization improvement. Details omitted for brevity.

D. Algorithm

Our overall approach is summarized in Algorithm $\boxed{1}$ The learning algorithm alternates between update of EEG decoder parameters θ and neural transformation parameters ϕ . We use reversible Heun method $\boxed{15}$ to perform update on ϕ .

Algorithm 1 Replay with Stochastic Neural Transformation

- 1: **REQUIRE:** EEG decoder parameters θ , stochastic neural transformation parameters ϕ , EEG decoder learning rate η , neural transformation update rate β ; transformation intensity G at each iteration, stored signal segments for replay \mathcal{M} ;
- 2: **for** i = 1 to K **do**
- 3: current signal segment (x_i, y_i) input to model.
- 4: retrieve stored signal segment for replay, i.e., $(x_m, y_m) \sim \mathcal{M}$
- 5: perform stochastic neural transformation $x_m(G) = \operatorname{Transform}(x_m, 0, G)$ (Eq. (3))
- 6: update ϕ with gradient ascent on $\max_{\boldsymbol{\phi}} [\mathcal{L}(\boldsymbol{x}_m(G), \boldsymbol{y}_m, \boldsymbol{\theta}, \boldsymbol{\phi}) \mathcal{L}(\boldsymbol{x}_m, \boldsymbol{y}_m, \boldsymbol{\theta}, \boldsymbol{\phi}) \lambda \mathbb{JS}(\boldsymbol{x}_m, \boldsymbol{x}_m(G))]$
- 7: replay of $(\boldsymbol{x}_m(G), \boldsymbol{y}_m)$ and jointly train with $(\boldsymbol{x}_i, \boldsymbol{y}_i)$ for update of model parameters $\boldsymbol{\theta}$: $\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i \eta \nabla_{\boldsymbol{\theta}} [\mathcal{L}(\boldsymbol{\theta}_i, \boldsymbol{x}_m(G), y) + \mathcal{L}(\boldsymbol{\theta}_i, \boldsymbol{x}_i, y_i)]$
- 8: update memory \mathcal{M} by reservoir sampling (RS), $\mathcal{M} = \text{RS}(\mathcal{M}, (\boldsymbol{x}_i, y_i))$
- 9: end for

III. EMPIRICAL EVALUATION

We constructed numerous sequential EEG decoding benchmarks for the problem setting on top of three large public datasets (BCI-IV 2a[16]] DEAP[17] and SEED[18] , which mimics the real world scenario of online sequential EEG decoding. Detailed ablation studies are performed in terms of transformation depth, memory size and different subject

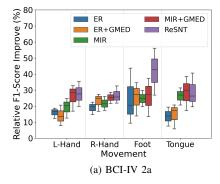
²http://bnci-horizon-2020.eu/database/data-sets

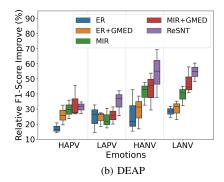
https://www.eecs.qmul.ac.uk/mmv/datasets/deap/download.html

⁴http://bcmi.sjtu.edu.cn/~seed/downloads.html

TABLE I: Performance on accuracy and BWT evaluated after decoder finished sequential learning on all subjects.

Dataset	BCI-IV 2a		DEAP		SEED	
Method	Accuracy	BWT	Accuracy	BWT	Accuracy	BWT
EWC	41.18±0.75	-12.64 ± 1.03	40.63±0.98	-11.41 ± 0.57	46.53±1.74	-13.68 ± 0.73
UCB	40.72 ± 0.66	-13.14 ± 0.58	38.24 ± 1.31	-12.16 ± 2.42	47.55 ± 2.67	-12.94 ± 2.35
ER	43.21 ± 0.49	-12.56 ± 0.72	41.45 ± 1.80	-8.28 ± 1.73	54.21±0.90	-10.05 ± 1.16
ER+GMED	46.57 ± 0.73	-12.10 ± 0.51	43.19 ± 2.03	-9.41 ± 0.96	56.25 ± 2.41	-11.98 ± 0.85
MIR	48.35 ± 1.14	-9.43 ± 0.87	43.68 ± 1.27	-7.55 ± 2.53	58.64 ± 1.56	-9.75 ± 1.47
MIR+GMED	49.42 ± 1.53	-7.49 ± 1.96	44.52 ± 0.74	-8.08 ± 2.35	59.97 ± 1.45	-8.36 ± 2.33
ReSNT	52.84±0.85	-6.13 ± 1.24	46.20±1.96	-7.33 ± 2.94	61.29±1.18	$-6.95{\pm}1.34$
sequential	36.46±0.39	-14.78 ± 0.41	33.41±0.92	-16.46 ± 0.90	38.53±0.27	-18.84 ± 0.91
joint train	81.70±0.52		71.57 ± 0.69		82.50±0.86	—





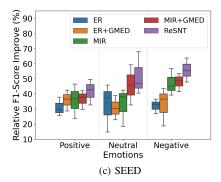


Fig. 3: F1-score improvement for different classes, relative to the baseline of sequential learning directly on base decoder. The x axis corresponds to different classes of decoding tasks. For DEAP dataset, the high/low arousal are abbreviated as HA/LA, positive/negative valence are abbreviated as PV/NV, e.g. high arousal positive valence corresponds to "HAPV".

ordering, etc., which offers in-depth understanding of model performance.

Baselines We incorporated a wide range of baselines for comparison in the experiment, including: 1) upper and lower bound: lower bound is sequential learning across different subjects directly with base decoder, and the upper bound is joint training with data from all subjects simultaneously available. 2) currently widely used continual learning approaches, including both regularization approaches such as EWC [19], UCB [20], and replay-based approaches including ER [11], MIR [21] and GMED [22].

Benchmark and Metrics We created new benchmarks to mimic the scenario of continual decoding over a long period of time across sequentially arriving subjects. Three different scenarios of subject ordering is explored in the ablation study, including 1) sequential order based on subject id, which is the default setting in experiments, 2) ascending order based on decoding difficulty, and 3) descending order based on decoding difficulty. We fed the streaming EEG signal of different subjects sequentially into the decoder based on the aforementioned subject ordering.

We evaluated the average testing accuracy on all subjects after completion of sequential learning, and performed measurement on the improvement of F1 score per individual class. We also performed evaluation on backward transfer (BWT). BWT is defined as BWT = $\frac{1}{N-1}\sum_{i=1}^{N-1}a_{N,i}-a_{i,i}$, with $a_{j,i}$ being the accuracy evaluated on subject j after sequential learning of subject i. N is the total number of

subjects. Negative value of BWT indicates the occurrence of catastrophic forgetting, and positive value shows learning of later subjects helps improve model performance on previous subjects.

Model Settings We adopt the shallow ConvNet architecture introduced in [3] to serve as base model for the decoding task. We used a small memory to store 200 signal segments by default for both datasets. The diffusion and drift network for the neural transformation is formed with two layers of ConvNet with filter size 4×4 for first layer and 2×2 for second layer. We set the transformation intensity G to be 0.1 by default, with its influence on model performance studied in ablation study. The memory size is 200 for BCI-IV 2a and DEAP datasets, and we used a memory size of 500 for SEED dataset due to it is significantly larger. We set the regularization weight λ to be 2.0 by default, and its sensitivity is also examined in ablation study. Other hyperparameter settings follow [3]. The reported results are repeated for 10 runs in each setting.

A. Performance Evaluation

Table I shows the performance in terms of test accuracy and BWT on ReSNT and comparison baselines. The accuracy is evaluated on all subjects after sequential learning finishes. ReSNT outperforms the baselines by a significant margin on all three datasets. The gain on accuracy is 3.42%, 1.68% and 1.32% for BCI-IV 2a, DEAP and SEED respectively. It has a margin of 1.36%, 0.75% and 1.41% on the three datasets in terms of BWT. The model is better at maintaining

TABLE II: Ablation study on the influence of different subject ordering towards replay-based approaches. We explored three scenarios, 1) sequential order of subject id, 2) ascending order of decoding difficulty, 3) descending order of decoding difficulty.

Scenario	ER	ER+GMED	MIR	MIR+GMED	ReSNT
		BCI-IV	2a Dataset		
Sequential Ascending Descending	$\begin{array}{c} 43.21{\pm}0.49 \\ 43.56{\pm}0.72 \\ 42.75{\pm}1.08 \end{array}$	$46.57{\pm0.73}\atop46.44{\pm1.25}\atop46.10{\pm0.51}$	$48.35{\pm}1.14 \\ 48.69{\pm}0.92 \\ 48.13{\pm}0.76$	$^{49.42\pm1.53}_{49.51\pm0.87}_{49.24\pm1.08}$	52.84 ± 0.85 53.10 ± 1.39 52.43 ± 0.61
		DEAP	Dataset		
Sequential Ascending Descending	$\begin{array}{c c} 41.45{\pm}1.80 \\ 42.28{\pm}0.74 \\ 40.57{\pm}1.13 \end{array}$	$\begin{array}{c} 43.19{\pm}2.03 \\ 42.63{\pm}1.35 \\ 42.26{\pm}1.27 \end{array}$	$43.68{\pm}1.27 \\ 44.20{\pm}0.85 \\ 43.32{\pm}2.28$	$44.52{\pm}0.74 \\ 45.95{\pm}1.19 \\ 44.03{\pm}1.62$	$\begin{array}{c} 46.20{\pm}1.96 \\ 46.71{\pm}1.55 \\ 46.14{\pm}0.80 \end{array}$
		SEED	Dataset		
Sequential Ascending Descending	$\begin{array}{c} 54.21{\pm}0.90 \\ 55.53{\pm}1.42 \\ 53.64{\pm}1.75 \end{array}$	$\begin{array}{c} 56.25{\pm}2.41 \\ 57.71{\pm}1.25 \\ 55.91{\pm}1.03 \end{array}$	58.64 ± 1.56 58.89 ± 0.73 58.28 ± 1.39	$59.97{\pm}1.45$ $59.85{\pm}1.28$ $59.36{\pm}0.93$	$61.29{\pm}1.18\\61.84{\pm}1.63\\60.97{\pm}1.10$
70%	Sequential FR Sequential Sequential FR Sequential Sequential FR Sequential Se	MIR ReSNT		ER+GMED MIR+GMED 2 3 4 5 6 subject	7 8 9
70% — — — — — — — — — — — — — — — — — — —	Sequential ER	(a) BCI-IV	7 2a Dataset 70% 50% 50% 888 808 808 808 808 808 808 808 808 8	■ ER+GMED ■ MIR+GMED	ReSNT
0 5	10 15 20 subject	25 30	20% 0 :	5 10 15 20 subject	25 30
		(b) DEA	AP Dataset		
90% — — — — — — — — — — — — — — — — — — —	Sequential ER	MIR ReSNT 12 14	900%	ER+GMED MIR+GMED 4 6 8 10 subject	12 14
		(c) SEE	D Dataset		

Fig. 4: Model performance on individual subject after sequential learning ends. The proposed model maintains relatively good performance on earlier subjects with forget mitigation mechanisms. (a) BCI-IV 2a dataset, (b) DEAP dataset, (c) SEED dataset.

performance of previous subjects during sequential learning as it produces more robust and diversified data in feature space for effective replay. We performed further decomposition of model performance at class-level, in terms of the different motor imagery and emotion recognition tasks. The result is summarized in fig. 3 The relative improvement in F1 score for the different classes shows the proposed model has varying levels of effectiveness for the different recognition tasks. Specifically, the model shows it is more effective on capturing foot movement compared to other baselines

TABLE III: Effect of transformation intensity. We observed the performance tend to converge with transformation intensity G > 0.1.

Trans. Intensity	G = 0.05	G = 0.1	G = 0.15	G = 0.2
BCI-IV 2a DEAP SEED	$\begin{array}{c} 52.25{\pm}0.92 \\ 45.57{\pm}0.63 \\ 61.04{\pm}1.35 \end{array}$	$\begin{array}{c} 52.84{\pm}0.85 \\ 46.20{\pm}1.96 \\ 61.29{\pm}1.18 \end{array}$	$\begin{array}{c} 53.29{\pm}1.31 \\ 46.41{\pm}1.22 \\ 61.67{\pm}0.80 \end{array}$	53.36 ± 0.62 46.44 ± 1.59 61.75 ± 1.37

when performing motor imagery decoding on BCI-IV 2a. For emotion recognition tasks, the model performs better than baselines in detecting high arousal negative valence emotions. We performed T-SNE visualization of transformed data, shown in fig. [5]. We observed the transformed data is in general more scattered and shows a loosely mixed pattern in feature space. Fig. [4] shows the performance of individual subjects after the sequential decoding finished. The proposed ReSNT model is better at maintaining performance of earlier subjects, with more effective knowledge retaining during the sequential decoding process.

B. Ablation Study

Influence of Different Subject Ordering

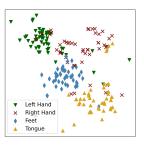
We explored the influence of different subject ordering on model performance. We run the model with three different subject ordering scenarios, ordered by 1) subject ID, 2) descending of decoding difficulty, 3) ascending of decoding difficulty respectively. The decoding difficulty is reflected as decoding accuracy during test. We summarized the result in Table III. Overall, the ascending order based on decoding difficulty turns out to yield the best overall performance, and descending order is the most challenging setting. This shows the model benefits from learning harder subjects at the beginning during the sequential decoding of different subjects.

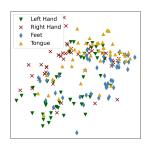
Effect of Transformation Intensity

We performed evaluation of model performance with respect to different transformation intensities depicted by G. The result is provided in Table $\boxed{\text{III}}$ We observed the model performance improves monotonically with the increase of G, this shows the effectiveness of consistency regularization which helps preserve the similarity between transformed data and original data. The model performance is trending towards convergence with G>0.1. We used G=0.1 as default setting in our experiments.

Influence of Memory Size

We performed ablation study on different memory sizes for all replay based approaches. The result is shown in Table IV. The performance improves for the replay based approaches with larger number of signal segments stored in memory. We observed the ReSNT model has a larger margin compared to comparison baselines in small memory settings, e.g. for BCI-IV 2a dataset, ReSNT outperforms other baselines by at least 3.93% when only 100 segments are stored for replay, and this margin gradually shrinks to 0.57% when 500 segments are allowed in memory. This shows the ReSNT model suits better to scenarios with smaller memory settings. We used a memory size of 200 for BCI-IV 2a and DEAP datasets, and memory size of 500 for SEED dataset as default settings.





- (a) Original Segments
- (b) Transformed Segments

Fig. 5: TSNE visualization at feature level of (a) original stored segments, (b) transformed segments in memory for the different motor imagery classes of BCI-IV 2a dataset. The stochastic neural transformation on stored segments produces more diverse and robust features for model generalization improvement.

TABLE IV: Model performance with respect to different number of signal segments in memory. The model shows improved avg. accuracy on all subjects with larger number of segments stored in memory.

Mem. Size	ER	ER+GMED	MIR	MIR+GMED	ReSNT	
BCI-IV 2a Dataset						
100 200 500	$ \begin{vmatrix} 40.63 \pm 1.32 \\ 43.21 \pm 0.49 \\ 51.47 \pm 0.80 \end{vmatrix} $	$\begin{array}{c} 42.49{\pm}0.85 \\ 46.57{\pm}0.73 \\ 53.28{\pm}1.79 \end{array}$	$\begin{array}{c} 45.96{\pm}2.01 \\ 48.35{\pm}1.14 \\ 56.82{\pm}1.53 \end{array}$	$\begin{array}{c} 46.22{\pm}0.79 \\ 49.42{\pm}1.53 \\ 57.56{\pm}0.85 \end{array}$	$\begin{array}{c} 50.15{\pm}1.46 \\ 52.84{\pm}0.65 \\ 58.13{\pm}0.92 \end{array}$	
DEAP Dataset						
100 200 500	35.80±0.76 41.45±1.80 48.64±0.59	$37.64\pm0.95\ 43.19\pm2.03\ 50.42\pm0.39$	39.56 ± 2.89 43.68 ± 1.27 52.26 ± 2.16	$40.10{\pm}3.24 \\ 44.52{\pm}0.74 \\ 53.97{\pm}1.91$	$42.63{\scriptstyle\pm1.70}\atop46.20{\scriptstyle\pm1.96}\\55.15{\scriptstyle\pm1.34}$	
SEED Dataset						
200 500 800	54.21±0.90 63.51±1.25 68.82±1.62	$56.25{\pm}2.41$ $65.17{\pm}0.57$ $69.39{\pm}0.94$	58.64 ± 1.56 66.69 ± 2.04 71.82 ± 1.59	$59.97{\pm}1.45$ $67.36{\pm}2.73$ $72.28{\pm}2.52$	$61.29{\pm}1.18 \\ 68.53{\pm}1.45 \\ 72.71{\pm}1.74$	

IV. CONCLUSION

For clinical applications of BCI systems, it is necessary for the system to perform decoding over a long period of time on streaming brain signals for patient assistance tasks. In this work, we propose a replay based approach for effective brain signal decoding in this challenging scenario. The model performs stochastic neural transformation on stored segments to generalize well on previous subjects and avoids overfitting. We formed a bi-level optimization framework for end to end training of the stochastic neural transformer together with the EEG decoder. Its effectiveness is evaluated on numerous newly formed benchmarks that mimic real world prolonged decoding scenarios. Further exploration is needed for the following items: 1) decoding of subject sequences with heterogeneous classes, 2) online continual decoding without usage of stored signal segments in memory.

V. ACKNOWLEDGEMENTS

This material is based upon work supported in part by the US National Science Foundation, Award #s: 1920920, 2125872 and 2223793.

REFERENCES

- I. Iturrate, J. Antelis, and J. Minguez, "Synchronous eeg brain-actuated wheelchair with automated navigation," in 2009 IEEE International Conference on Robotics and Automation, 2009, pp. 2318–2325.
- [2] A. Campbell et al., "Neurophone: Brain-mobile phone interface using a wireless eeg headset," in *Proceedings of the Second ACM SIGCOMM* Workshop on Networking, Systems, 2010.
- [3] R. T. Schirrmeister et al., "Deep learning with convolutional neural networks for eeg decoding and visualization," *Human Brain Mapping*, vol. 38, no. 11, pp. 5391–5420, 2017.
- [4] T. Duan, Z. Wang, S. Liu, Y. Yin, and S. N. Srihari, "Uncer: A framework for uncertainty estimation and reduction in neural decoding of eeg signals," *Neurocomputing*, vol. 538, 2023.
- [5] V. J. Lawhern, A. J. Solon, N. R. Waytowich, S. M. Gordon, C. P. Hung, and B. J. Lance, "EEGNet: a compact convolutional neural network for EEG-based brain-computer interfaces," *Journal of Neural Engineering*, vol. 15, no. 5, p. 056013, jul 2018.
- [6] T. Duan and S. N. Srihari, "Layerwise interweaving convolutional LSTM," in Advances in Artificial Intelligence - 30th Canadian Conference on Artificial Intelligence, Canadian AI 2017, Edmonton, AB, Canada, vol. 10233, 2017, pp. 272–277.
- [7] W.-L. Zheng and B.-L. Lu, "Personalizing eeg-based affective models with transfer learning," in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, ser. IJCAI'16. AAAI Press, 2016, pp. 2732–2738.
- [8] T. Duan, M. A. Shaikh, M. Chauhan, J. Chu, R. K. Srihari, A. Pathak, and S. N. Srihari, "Meta learn on constrained transfer learning for low resource cross subject eeg classification," *IEEE Access*, vol. 8, pp. 224 791–224 802, 2020.
- [9] D. Lopez-Paz and M. Ranzato, "Gradient episodic memory for continual learning," Advances in Neural Information Processing Systems, 2017.
- [10] Z. Wang, T. Duan, L. Fang, Q. Suo, and M. Gao, "Meta learning on a sequence of imbalanced domains with difficulty awareness," 2021 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 8927–8937, 2021.
- [11] A. Chaudhry, M. Rohrbach, M. Elhoseiny, T. Ajanthan, P. K. Dokania, P. H. S. Torr, and M. Ranzato, "Continual learning with tiny episodic memories," *CoRR*, vol. abs/1902.10486, 2019.
- [12] Z. Wang, L. Shen, T. Duan, D. Zhan, L. Fang, and M. Gao, "Learning to learn and remember super long multi-domain task sequence," 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 7972–7982, 2022.
- [13] J. Yoon, E. Yang, J. Lee, and S. J. Hwang, "Lifelong learning with dynamically expandable networks," in *International Conference on Learning Representations*, 2018.
- [14] R. Gontijo-Lopes, S. Smullin, E. D. Cubuk, and E. Dyer, "Tradeoffs in data augmentation: An empirical study," in *International Conference on Learning Representations*, 2021.
- [15] P. Kidger, J. Foster, X. C. Li, and T. Lyons, "Efficient and accurate gradients for neural sdes," in *Advances in Neural Information Processing Systems*, vol. 34. Curran Associates, Inc., 2021, pp. 18747–18761.
- [16] M. Tangermann et al., "Review of the bci competition iv," Frontiers in Neuroscience, vol. 6, p. 55, 2012.
- [17] S. Koelstra, C. Muhl, M. Soleymani, J. Lee, A. Yazdani, T. Ebrahimi, T. Pun, A. Nijholt, and I. Patras, "Deap: A database for emotion analysis ;using physiological signals," *IEEE Transactions on Affective Computing*, vol. 3, no. 1, pp. 18–31, Jan 2012.
- [18] R.-N. Duan, J.-Y. Zhu, and B.-L. Lu, "Differential entropy feature for EEG-based emotion classification," in 6th International IEEE/EMBS Conference on Neural Engineering (NER). IEEE, 2013, pp. 81–84.
- [19] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell, "Overcoming catastrophic forgetting in neural networks," *Proceedings of the national* academy of sciences, 2017.
- [20] S. Ebrahimi, M. Elhoseiny, T. Darrell, and M. Rohrbach, "Uncertainty-guided continual learning with bayesian neural networks," arXiv preprint arXiv:1906.02425, 2019.
- [21] R. Aljundi, L. Caccia, E. Belilovsky, M. Caccia, M. Lin, L. Charlin, and T. Tuytelaars, *Online Continual Learning with Maximally Interfered Retrieval*. Red Hook, NY, USA: Curran Associates Inc., 2019.
- [22] X. Jin, A. Sadhu, J. Du, and X. Ren, "Gradient based memory editing for task-free continual learning," ArXiv, vol. abs/2006.15294, 2020.