

# Offline Learning of Closed-Loop Deep Brain Stimulation Controllers for Parkinson Disease Treatment

Qitong Gao  
Electrical and Computer Engineering  
Duke University  
qitong.gao@duke.edu

Stephen L. Schmidt  
Biomedical Engineering  
Duke University  
stephen.schmidt@duke.edu

Afsana Chowdhury  
Electrical and Computer Engineering  
Duke University  
afsana.chowdhury@duke.edu

Guangyu Feng  
Electrical and Computer Engineering  
Duke University  
guangyu.feng@duke.edu

Jennifer J. Peters  
Biomedical Engineering  
Duke University  
jennifer.peters@duke.edu

Katherine Genty  
Neurosurgery  
Duke University  
katherine.genty@duke.edu

Warren M. Grill  
Biomedical Engineering  
Duke University  
warren.grill@duke.edu

Dennis A. Turner  
Neurosurgery  
Duke University  
dennis.turner@duke.edu

Miroslav Pajic  
Electrical and Computer Engineering  
Duke University  
miroslav.pajic@duke.edu

## ABSTRACT

Deep brain stimulation (DBS) has shown great promise toward treating motor symptoms caused by Parkinson's disease (PD), by delivering electrical pulses to the Basal Ganglia (BG) region of the brain. However, DBS devices approved by the U.S. Food and Drug Administration (FDA) can only deliver continuous DBS (cDBS) at a fixed amplitude; this energy inefficient operation reduces battery lifetime of the device, cannot adapt treatment dynamically for activity, and may cause significant side-effects (e.g., gait impairment). In this work, we introduce an offline reinforcement learning (RL) framework, allowing the use of past clinical data to train an RL policy to adjust the stimulation amplitude in real time, with the goal of reducing energy use while maintaining the same level of treatment (i.e., control) efficacy as cDBS. Moreover, clinical protocols require the safety and performance of such RL controllers to be demonstrated ahead of deployments in patients. Thus, we also introduce an offline policy evaluation (OPE) method to estimate the performance of RL policies using historical data, before deploying them on patients. We evaluated our framework on four PD patients equipped with the RC+S DBS system, employing the RL controllers during monthly clinical visits, with the overall *control efficacy* evaluated by severity of symptoms (i.e., bradykinesia and tremor), changes in PD biomarkers (i.e., local field potentials), and patient ratings. The results from clinical experiments show that our

RL-based controller maintains the same level of control efficacy as cDBS, but with significantly reduced stimulation energy. Further, the OPE method is shown effective in accurately estimating and ranking the expected returns of RL controllers.

## CCS CONCEPTS

• **Applied computing** → **Life and medical sciences**; • **Computing methodologies** → **Control methods**.

## KEYWORDS

Deep Brain Stimulation, Offline Reinforcement Learning, Offline Policy Evaluation

### ACM Reference Format:

Qitong Gao, Stephen L. Schmidt, Afsana Chowdhury, Guangyu Feng, Jennifer J. Peters, Katherine Genty, Warren M. Grill, Dennis A. Turner, and Miroslav Pajic. 2023. Offline Learning of Closed-Loop Deep Brain Stimulation Controllers for Parkinson Disease Treatment. In *ACM/IEEE 14th International Conference on Cyber-Physical Systems (with CPS-IoT Week 2023) (ICCPs '23)*, May 9–12, 2023, San Antonio, TX, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3576841.3585925>

## 1 INTRODUCTION

Currently, around 1.05 million individuals in the United States are affected by Parkinson's disease (PD) [44]. Deep brain stimulation (DBS) is an effective treatment to reduce PD symptoms such as tremor and bradykinesia [3, 12, 13, 49]. A DBS system consists of electrodes that are placed into the Basal Ganglia (BG) region of the brain, and a pulse generator implanted in the chest to generate trains of short electrical pulses (see Fig. 1). Existing FDA-approved DBS solutions are limited to continuous DBS (cDBS). These devices are programmed to stimulate at a fixed amplitude, with the specific parameters determined by clinicians through trial-and-error [52]. However, such stimuli usually lead to extensive energy consumption, significantly reducing the battery lifetime of the device. Moreover, over-stimulated patients, even intermittently, may suffer from side-effects such as dyskinesia and speech impairment [5]. As a result, developments of closed-loop DBS controllers that are more

This work is sponsored in part by the NSF CNS-1837499 award and the National AI Institute for Edge Computing Leveraging Next Generation Wireless Networks, Grant CNS-2112562, as well as by NIH UH3 NS103468. Investigational Summit RC+S systems and technical support provided by Medtronic PLC. Apple Watches were provided by Rune Labs.

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 the author(s) 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](mailto:permissions@acm.org). ICCPS '23, May 9–12, 2023, San Antonio, TX, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 979-8-4007-0036-1/23/05...\$15.00 <https://doi.org/10.1145/3576841.3585925>

responsive to activity and patient state (i.e., context) are of considerable interest to clinicians, patients, and the community.

Existing DBS control methods focus on simply switching on/off the stimulation or scaling up/down its intensity in a proportional control approach, conditioned on the change of specific biomarkers, i.e., when they cross over some pre-determined thresholds [1, 2, 5, 40, 41]. Biomarkers include local field potentials (LFPs) and electroencephalography (EEG) from the BG, as well as accelerometry data and electromyography obtained from wearable devices [50]. Though such methods have improved energy efficiency [25, 41], they still require substantial efforts to experiment and fine-tune the thresholds for each specific patient. Moreover, the patient may suffer from sub-optimal DBS settings in between clinical visits with poor symptom control due to varying patient state. For example, exercise or fluctuations in medication dosage or timing could affect their PD symptoms and DBS control, so the tuning results may be biased. Consequently, the **challenge (I)** of developing closed-loop DBS controllers is to ensure that the control policy can perform consistently over diverse and dynamic patient contexts and states.

Reinforcement learning (RL) has shown considerable potential in control over complicated systems [15, 21, 22, 46], and various RL-based approaches have been proposed to facilitate closed-loop DBS [19, 23, 48, 52]. Specifically, several approaches [23, 48, 52] model EEG and LFP as the state space of the RL environment and use temporal difference learning or fitted Q-iteration to design control policies adapting stimulation amplitudes/frequencies to conserve energy usage. The deep actor-critic based approach proposed in [19] further allows the temporal pattern of the stimuli to be adapted over time, benefiting from the use of deep RL techniques capable of searching in larger state and action space. Although such methods achieve satisfactory control of efficacy and energy savings jointly, *they have only been evaluated in simulations*, i.e., on computational BG models [30, 58]. One may assume that unlimited training data can be obtained from such models, which is contrary to the real-world case where the device programming is done in clinics and the patient only participates sparsely over time.

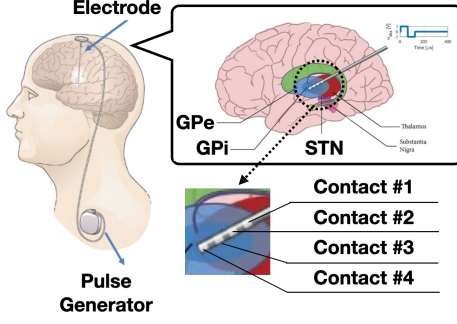
Another limitation of directly using deep RL methods for real-time DBS control is the computational complexity of evaluating the RL policies *in vivo*, as they are usually represented by deep neural networks (DNNs) that may require millions of multiplications in a single forward pass. The resource-constrained implantable devices (e.g., Fig. 1) may not support or facilitate such computations. Thus, the **challenge (II)** of closed-loop DBS is to ensure that the controller can be designed with limited training samples and executed without the need of extensive computing resources. Further, in contrast to simulated or robotic environments where most RL policies can be deployed directly for performance evaluation, the safety and control efficacy of the controllers directly used on patients need to be thoroughly evaluated before each test condition starts [51]. Hence, the **challenge (III)** of enabling closed-loop DBS therapies in patients is being able to proactively provide accurate estimations of the expected performance of the controllers.

Consequently, in this paper, we first introduce an offline RL framework to address the challenges (I) and (II) above, resulting in a closed-loop DBS system that is both *effective (in terms of therapy)* and *energy-efficient*. Specifically, we model the BG regions of the brain

as a Markov decision process (MDP), capturing the underlying neuronal activities in response to the stimuli. Then, the deep actor-critic algorithm [39] is adapted to adjust the amplitude of the stimuli according to the changes in LFPs. A total of four patients, equipped with the Medtronic Summit RC+S DBS devices [59], participated in the data collection and testing trials in clinics. Given that the deep actor-critic framework is considered offline RL and can leverage all historically collected trajectories, i.e., experience replay to facilitate optimizing the control policy, we address challenge (I) by varying the level of activities, medications etc. of the patients before and during the trials. Similarly, experience collected from non-RL controllers can also be used to update the policy; for example, in the early stage of learning, a controller that generates uniformly random amplitudes (within some range) can facilitate exploring the state and action space. We also introduce model distillation/compression [26] techniques specifically for the DBS systems, such that the RL policies can be captured by deep neural networks (DNNs) with significantly fewer nodes, whose forward passes can be executed within the required control rates, addressing challenge (II).

To address challenge (III), we introduce a model-based offline policy evaluation (OPE) method that captures the underlying dynamics of the considered MDP, where the expected returns of the control policy can be estimated by the mean return of the trajectories rolled out from the learned model, without directly deploying the policy to the patient. In each DBS trial, the control efficacy is evaluated from various sources, including LFP biomarkers recorded from the implantable DBS device, patient responses to bradykinesia tests, satisfaction level reported by the patient, and the overall tremor severity quantified from accelerometry data collected by external wearable devices (e.g., smart watch). Note that each of the latter three criteria is only evaluated once at the end of each trial; yet they are imperative for evaluating the control efficacy from the patient's side. These efficacy metrics are thus considered sparsely available compared to the LFPs that can be sensed in each time step, which limits the use of existing OPE methods, including importance sampling (IS) [16, 54], distributional correction estimations (DICE) [47], and the model-based OPE [20], as these do not allow for explicitly capturing/modeling such end-of-session rewards. Our OPE method can capture such behaviors through a specially designed architecture and training objective, outperforming existing methods as we show in clinical experiments.

The contributions of this work are three-fold: (i) to the best of our knowledge, this is the first '*full-stack*' offline RL methodology that facilitates both *optimizing* and *evaluating* RL-based DBS control policies using historical data; (ii) we developed an RL-based DBS controller whose performance is validated through clinical trials with PD patients, demonstrating *reduced energy consumption with non-inferior control efficacy compared to cDBS* – **this is the first effective closed-loop DBS control that is not an ON/OFF switching, or scaling up/down proportionally, and has been extensively tested in clinic (i.e., on patients)**; (iii) our OPE method effectively captures the end-of-session rewards, leading to accurate estimations of control efficacy using the data collected in clinic; thus, helps demonstrate the effectiveness of the policies to be tested proactively, and can be used to prioritize the policies that could lead to better performance within the limited amount of testing time.



**Figure 1: An implantable deep brain stimulation (DBS) device.** The stimuli, generated by the pulse generator at a given amplitude and frequency, are delivered to the basal ganglia (BG) through multi-contact electrodes. Each electrode has four contacts; two stimulate the BG and two sense local field potentials (LFPs) that may be used for control feedback.

This paper is organized as follows. Sec. 2 provides the basics of DBS, RL, and OPE, before our clinical closed-loop DBS setup is introduced in Sec. 3. In Sec. 4, the offline RL framework is introduced, enabling training and updating RL controllers with historical data. Sec. 5 introduces the model-based OPE approach to estimate performance of RL policies. Sec. 6 presents the results of the experimental evaluations on patients, before concluding remarks in Sec. 7.

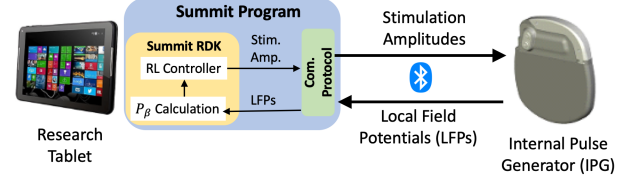
## 2 PRELIMINARIES AND MOTIVATION

In this section, we first introduce DBS, before presenting in the next section the DBS experimental setup we developed for clinical trials, including sensing, communication and control. Also, preliminaries for offline RL and OPE are briefly introduced; more comprehensive reviews of RL and OPE can be found in [19, 20, 39, 57].

### 2.1 The Need for Closed-Loop DBS

PD is caused by progressive death of dopaminergic neurons in the substantia nigra region of the brain. This change in dopaminergic signaling results in pathological activity in the BG regions targeted by DBS, *globus pallidus pars interna* (GPi), *globus pallidus pars externa* (GPe) and subthalamic nucleus (STN); see Fig. 1. Given the reduced number of neurons, the level of dopamine generally decreases in BG, leading to various motor symptoms such as bradykinesia and tremor [7, 11, 34]. Physiologically, the effect of PD can be captured by the changes in LFPs in GPi, GPe and STN. Specifically, PD can cause abnormal neuron firings in these regions, and lead to increased beta-band (13-35 Hz) amplitude ( $P_\beta$ ), referred to as the beta amplitude, of the LFPs [20].

Existing research-only DBS devices are capable of capturing the changes in LFPs through the multi-contact electrodes implanted in the BG. As illustrated in Fig. 1, we used 4-contact electrodes placed in the STN and GP regions. Monopolar stimulation was delivered on a single contact on each lead (with the case serving as counter-electrode). The two contacts surrounding the stimulation contact were used for sensing LFPs (i.e., sandwich sensing). Existing devices providing open-loop cDBS stimulate pulses at a fixed amplitude, which in most cases can correct the abnormal neuronal activity [37].



**Figure 2: The overall architecture of the RC+S DBS system.** The Summit research and development kit (RDK) can be used to configure the Summit program, allowing us to compute the beta amplitude ( $P_\beta$ ) and execute the RL controller.

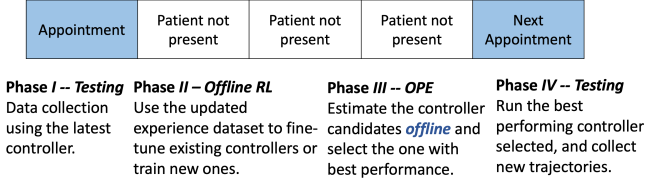
However, constantly stimulating with high amplitudes significantly reduces the battery lifetime of the DBS device and may cause serious side-effects such as speech impairment [4, 40, 60]. Consequently, it is important to design DBS controllers that are *effective* (from the control, i.e., therapy, perspective) and *energy-efficient*.

As discussed in Introduction, current aDBS approaches require considerable time and effort for the patients and their healthcare providers to determine the thresholds through trial-and-error [63]. Several deep-RL-based controllers have been proposed for closed-loop DBS, which can adapt the amplitude of the stimulation pulses in real time [19, 20] in response to changes in the feedback signals (e.g.,  $P_\beta$ ). However, such frameworks are only validated through numerical simulations, i.e., on *simplified* computational BG models, instead of clinical trials with human participants. In real world, *substantial* historical experience, or trajectories collected from past interactions between the controller and the environment (patient), may be necessary to learn an RL policy with suitable control efficacy and patient satisfaction [38]. Offline RL holds promise to resolve this challenge, as it can use the data collected from any type of controllers, including cDBS or simply a policy switching between arbitrary stimulation amplitudes/frequencies, to optimize an RL control policy. Moreover, each time before a new control policy is deployed to the patient, the clinicians need to assess its effectiveness and may require justifications toward its estimated control efficacy and performance [51]. OPE can facilitate such use cases, as it is capable of estimating the expected return of RL policies using historical trajectories, bridging the gap between the offline RL training and evaluations. Preliminaries for offline RL and OPE are presented in two subsections below.

### 2.2 Offline Reinforcement Learning

Offline RL has proven useful in many domains, including robotics [18, 22], healthcare [21], etc., since it can optimize the control policies without requiring the environment to be presented, which guarantees the safety of the learning process. Further, it does not require the training data to be exclusively collected by the control policy being updated, leading to improved sample efficiency. To facilitate offline RL, the underlying dynamical environments are firstly modeled as Markov decision processes (MDPs).

**DEFINITION 2.1 (MDP).** An MDP is a tuple  $\mathcal{M} = (\mathcal{S}, s_0, \mathcal{A}, \mathcal{P}, R, \gamma)$ , where  $\mathcal{S}$  is a finite set of states;  $s_0$  is the initial state;  $\mathcal{A}$  is a finite set of actions;  $\mathcal{P}$  is the transition function defined as  $\mathcal{P} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ ;  $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is the reward function, and  $\gamma \in [0, 1)$  is a discount factor.



**Figure 3: Timeline for training RL-based DBS controllers in clinical studies. Since only limited data can be collected during each clinical visit, offline RL can be used to fine-tune existing or train new controllers using all the historical data. Then, offline policy evaluation (OPE) facilitates choosing the possible top-performing ones to be tested in the next visit.**

Then, the RL policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  determines the action  $a = \pi(s)$  to be taken at a given state  $s$ . The accumulated return under a policy  $\pi$  can be defined as follows.

**DEFINITION 2.2 (ACCUMULATED RETURN).** Given an MDP  $\mathcal{M}$  and a policy  $\pi$ , the accumulated return over a finite horizon starting from the stage  $t$  and ending at stage  $T$ , for  $T > t$ , is defined as

$$G_t^\pi = \sum_{k=0}^{T-t} \gamma^{t+k} r_{t+k}, \quad (1)$$

where  $r_{t+k}$  is the return at the stage  $t+k$ .

The goal of offline RL can now be defined as follows.

**PROBLEM 1 (OFFLINE REINFORCEMENT LEARNING).** Given an MDP  $\mathcal{M}$  with unknown transition dynamics  $\mathcal{P}$ , a pre-defined reward function  $R$ , and a experience replay buffer  $\mathcal{E}^\mu = \{(s_0, a_0, r_0, s_1), \dots, (s_{T-1}, a_{T-1}, r_{T-1}, s_T)\}^{(0)}, [(s_0, a_0, r_0, s_1), \dots]^{(1)}, \dots | a_t \sim \mu(a_t|s_t)\}$  containing trajectories collected over an unknown behavioral policy  $\mu$ , find the target policy  $\pi^*$  such that the expected accumulative return starting from the initial stage over the entire horizon is maximized, i.e.,

$$\pi^* = \underset{\pi}{\operatorname{argmax}} \mathbb{E}_{s,a \sim \rho^\pi, r \sim R} [G_0^\pi]; \quad (2)$$

here,  $\rho^\pi$  is the state-action visitation distribution under policy  $\pi$ .

The deep actor-critic RL framework [39] can be leveraged to solve (2). Other value-based RL methods such as conservative Q-learning [36] and implicit Q-learning [33] could also be considered; however, actor-critic methods can in general reduce the variance of gradient estimations and result in faster convergence [19, 45, 64]. Here, we specifically consider the deterministic version of actor-critic [39], instead the one producing stochastic policies [24], as it would be easier to demonstrate the effectiveness of deterministic policies in clinics, as well as via OPE methods introduced below. Details on the deep actor-critic algorithm [39] are provided in Appendix C.1.

### 2.3 Offline Policy Evaluation for DBS

OPE allows the use of experience replay buffer to estimate the expected return of RL policies, without the need of deploying them to the environment directly. Fig. 3 illustrates the use case of OPE in the context of DBS clinical testing. Specifically, during phase I and II, offline RL uses all trajectories collected historically to train RL policies following different hyper-parameters etc. Then, in phase

III, OPE can be used to estimate and rank the expected return of these policies, where the top-performing ones can be deployed during the next clinic visit (phase IV). Consequently, OPE can effectively reduce the number of testing sessions needed, so the policies that show promise attaining better performance can be thoroughly tested within the short time frame. Also, it can demonstrate the effectiveness of the policies to be deployed in clinics.

The goal of OPE can be defined as follows.

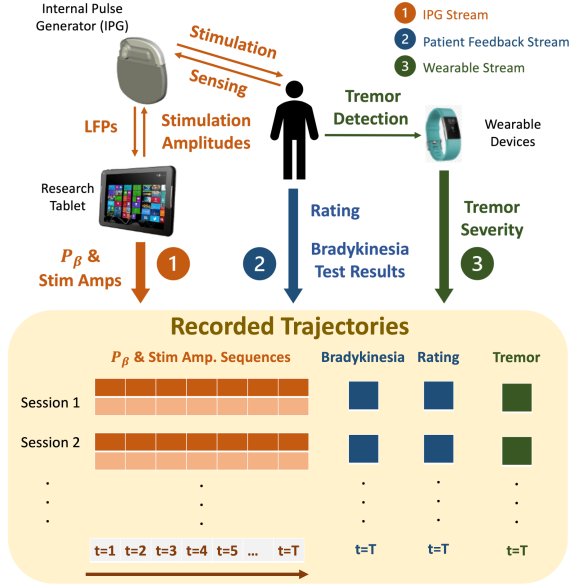
**PROBLEM 2 (OFFLINE POLICY EVALUATION).** Consider a target policy  $\pi$ , and *off-policy* trajectories  $\mathcal{E}^\mu = \{(s_0, a_0), (s_1, a_1), \dots | a_t = \mu(s_t)\}$ , collected following a behavioral policy  $\mu \neq \pi$ , over an MDP  $\mathcal{M}$ . The OPE goal is to estimate the expected return of the target policy  $\pi$ , i.e.,  $\mathbb{E}_{s,a \sim \rho^\pi, r \sim R} [G_0^\pi]$ .

Most existing OPE methods, such as [10, 14, 16, 29, 42, 54, 61, 62, 65], are heavily based on importance sampling (IS) and could result in inconsistent estimations due to the high variance of the IS weights [10, 42]. On the other hand, model-based OPE methods have shown strengths in estimating the expected returns more accurately [14, 20], by directly capturing the MDP transitions and rewards. The variational encoding-decoding based deep latent MDP model (DLMM) introduced in [20] is shown to be effective evaluating control policies for a computational BG model. Specifically, DLMM is derived following the variational inference framework from [32]. The basics of DLMM are provided in Appendix C.2, and we refer the readers to [32] for basics of variational inference. In Sec. 5, we extend it toward the clinical use case considered in this work, to allow for including the QoC metrics that can be only evaluated once in each session, such as the bradykinesia results, patient ratings, and tremor severity, which will be available as illustrated in Fig. 4.

### 3 DBS SETUP USED IN CLINICAL TRIALS

We build on the research-only Medtronic's Summit RC+S system [59] to enable testing of RL-based controllers in clinical trials. The overall architecture of the RC+S-based system we developed is illustrated in Fig. 2. Specifically, Medtronic provides the code and communication APIs (Summit program), which enable the stimulation amplitude of the pulses delivered by the internal pulse generator (IPG) to be adapted over time. The Summit program is developed using the C# language under the .NET framework, which we extended to execute RL policies leveraging the provided Summit research development kit (RDK), requiring the use of a Windows OS.

Thus, a research tablet is used for the execution of the developed DBS controllers; the desired stimulation amplitude is computed for each control cycle (every 2 seconds) and sent to the IPG over Bluetooth<sup>TM</sup>, using proprietary communication and security protocols. On the other hand, the IPG transmits to the controller the LFPs captured from the BG, from which the beta amplitude of the LFPs, denoted by  $P_\beta$ , is calculated and used as a quality of control (QoC) metric as well as potential control feedback signals (i.e., inputs to the RL controller). Each clinical trial session lasts 5-20 minutes depending on the schedule of the visit, and multiple controllers can be tested across different sessions. All the computed  $P_\beta$  and stimulation amplitudes applied over time are logged for future training and evaluation purposes, as summarized in Fig. 4. For the developed system design, we obtained the FDA's Investigative



**Figure 4: Setup of the developed DBS clinical testing procedure.** A total of three data streams are collected: (1) the LFPs and stimulation amplitudes are recorded over time; the logged trajectories are used to evaluate the performance of deployed RL controllers, as well as training data for further fine-tuning; (2) patient feedback including results from bradykinesia tests and a rating on the scale between 1-10; (3) patient tremor severity captured by wearable devices.

Device Exception (IDE) G180280, which has allowed us to perform human experiments according to an Institutional Review Board (IRB) protocol approved by Duke University Medical Center.

In addition to  $P_\beta$ , three other QoC metrics are collected from every patient at the end of each session. Specifically, near the end of each session, the patient is asked to perform 10 seconds of hand grasps (rapid and full extension and close of all fingers) maneuver [55] to evaluate the severity of the possible bradykinesia caused by PD. Such hand motions are captured by a leap motion sensor by Ultraleap [8]. Then, the elapsed time between any two consecutive open fist is captured and recorded by the sensor, after which the grasp frequency can be calculated as

$$QoC_{grasp} = \frac{1}{\frac{1}{N-1} \sum_{i=1}^{N-1} t_{(i,i+1)}}; \quad (3)$$

here,  $N$  is the total number of open fists throughout the 10 s test, and  $t_{(i,i+1)}$  is the time spent between the  $i$ -th and  $i+1$ -th grasp. Further, at the end of each session, the patient provides a score between 1-10, with 10 indicating the highest level of satisfaction with the treatment received in the past session, and 1 being the lowest, i.e.,

$$QoC_{rate} \in [1, 10] \subset \mathbb{Z}^+. \quad (4)$$

The grasp frequency and rating for each session are also recorded, which corresponds to the patient feedback stream in Fig. 4.

Throughout all sessions, an Apple watch is worn by the patient at their wrist, where the Apple's movement disorders kit [53] is used to analyze the accelerometry movements, classifying the patient's tremor severity as no-tremor, slight, mild, moderate and strong

every 1 minute, following StrivePD's implementation [9]. At the end of each session, an overall tremor severity is recorded as the fraction of time the patient experiencing mild ( $T_{mild}$ ), moderate ( $T_{moderate}$ ) or strong ( $T_{strong}$ ) tremor over the entire session with length  $T_{session}$ , i.e.,

$$QoC_{tremor} = \frac{T_{mild} + T_{moderate} + T_{strong}}{T_{session}} \times 100\%. \quad (5)$$

The three data streams are collected from all trial sessions after each clinical visit. Moreover, each time a patient may come into the clinic with slightly different PD conditions (e.g., pathology progression over time), medication prescriptions, activity levels etc.; thus, our goal is to capture impact of such changes by the data collection process, in order to facilitate the training and testing the offline RL and OPE frameworks for DBS.

## 4 OFFLINE RL DESIGN OF DBS CONTROLLERS

In this section, we employ offline RL for learning control policies for DBS clinical trials, starting from the formulation of an MDP  $\mathcal{M}$  capturing the underlying neurological dynamics in the BG, and the policy distillation technique that allows for reducing the computational time and resource needed to evaluate the RL policies (represented by DNNs).

### 4.1 Modeling the BG as an MDP

We now define the elements of an MDP  $\mathcal{M} = (\mathcal{S}, s_0, \mathcal{A}, \mathcal{P}, R, \gamma)$ .

**State Space  $\mathcal{S}$  and the Initial State  $s_0$ .** As discussed in Sec. 2.1 and 3, our DBS controller supports calculation of  $P_\beta$  from LFPs, and the changes in  $P_\beta$  can be used as a biomarker for PD-levels for some patients. Thus, we consider the MDP state, at a *discrete* time step  $t$ , as a historical sequence of  $P_\beta$  sampled at a fixed intervals, captured by  $m \in \mathbb{Z}^+$ , over a sliding queue of size  $W \in \mathbb{Z}^+$ , i.e.,

$$s_t = [\beta_{(\tilde{t}-(W-1)m)}, \beta_{(\tilde{t}-(W-2)m)}, \dots, \beta_{(\tilde{t}-2m)}, \beta_{(\tilde{t}-m)}, \beta_{(\tilde{t})}]. \quad (6)$$

Here,  $\beta_{(\tilde{t})}$ 's are the  $P_\beta$  evaluated at the elapsed time  $\tilde{t}$  since the clinical trial starts,  $m$  is configurable in our system design (Fig. 2), and we used  $m = 2$  corresponding to calculating  $P_\beta$  every 2 s, resulting in 20 s time-windows for  $W = 10$  elements in the queue; finally,  $s_t \in \mathbb{R}^W$  is the state at  $t$ -th (discrete) step of the MDP. The initial state  $s_0$  is considered to be the  $\beta$  sequence collected right before the clinical trial starts, i.e., from  $\tilde{t} = -(W-1)m$  to  $\tilde{t} = 0$ .

**Action Space  $\mathcal{A}$ .** The amplitude of DBS stimulation pulses can be changed in pre-defined (discrete) time steps, i.e., every 2 seconds for the developed controllers. We consider the actions  $a_t$  as the percentage of the cDBS amplitude determined by clinicians; i.e.,  $a_t \in [0, 1] \subset \mathbb{R}$ , where  $a_t = 0$  and  $a_t = 1$  correspond to no-DBS and stimulation with the same amplitude as in cDBS, respectively.

**Transition Dynamics  $\mathcal{P} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ .** Every time after the stimulation amplitude is adjusted following  $a_t$ , the system computes the latest  $\beta_{(\tilde{t}+m)}$  using the LFPs sent back from the IPG; this leads to the MDP state at the  $(t+1)$ -th (discrete) step as

$$s_{t+1} = [\beta_{(\tilde{t}-(W-2)m)}, \beta_{(\tilde{t}-(W-3)m)}, \dots, \beta_{(\tilde{t})}, \beta_{(\tilde{t}+m)}], \quad (7)$$

i.e., the left-most element in (6) is pushed out, with  $\beta_{(\tilde{t}+m)}$  appended to the right-end. Note that we define the MDP states  $s_t$  and actions



$a_t$  over discrete time steps,  $t$ 's, instead the elapsed time  $\tilde{t}$ , for the conciseness of equations and presentations below. Now, the MDP transitions are captured to directly follow  $s_{t+1} \sim \mathcal{P}(s_t, a_t)$ .

**Reward Function  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ .** Following from the setup of the DBS system (Sec. 3), we define the rewards as

$$R(s_t, a_t, s_{t+1}) = \begin{cases} r_a - C_1 \cdot a_t, & \text{if } \tilde{\beta}_{(\tilde{t}+m)} < \xi_\beta; \\ r_b - C_1 \cdot a_t, & \text{if } \tilde{\beta}_{(\tilde{t}+m)} \geq \xi_\beta; \end{cases} \quad (8)$$

specifically, if the beta amplitude received at the  $(t+1)$ -th step,  $\tilde{\beta}_{(\tilde{t}+m)}$ , is less than some threshold  $\xi_\beta$ , then a non-negative reward  $r_a$  is issued along with the term  $-C_1 \cdot a_t$  ( $C_1 > 0, C_1 \in \mathbb{R}$ ) penalizing over-usage of large stimulation amplitudes (for better energy efficiency). On the other hand, if  $\tilde{\beta}_{(\tilde{t}+m)}$  is greater than the threshold  $\xi_\beta$ , a negative reward  $r_b$  will be used to replace  $r_a$  above.

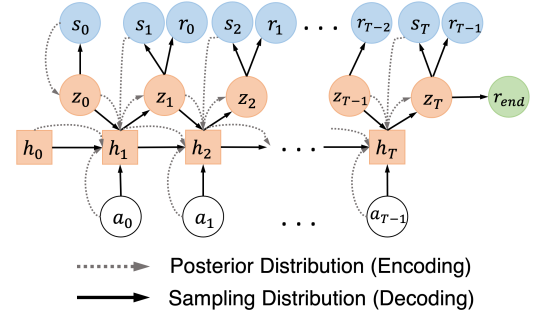
**REMARK 4.1.** *The reward functions used for RL training do not consider the QoC metrics that are available not at every step of the control execution (i.e., every 2 s) but only at the end of each clinical session, i.e.,  $QoC_{grasp}$ ,  $QoC_{rate}$ ,  $QoC_{tremor}$  from (3), (4), (5). The reason is that the horizon  $T$  is usually large and the their coverage can be very sparse. Instead, these QoC metrics serve as great measurements quantifying how well the policies perform, which are thus leveraged by the OPE techniques introduced in Sec. 5.*

For the introduced MDP  $\mathcal{M}$ , we leverage the offline RL framework introduced in Sec. 2.2 to search for the target policy  $\pi^*$ . Following from Problem 1, it requires an experience replay buffer  $\mathcal{E}^\mu$  that consists of historical trajectories collected over some behavioral policy  $\mu$ . At the beginning of offline RL training, exploration of the environment is deemed more important than exploitation [28]. Hence, a controller that generates random actions uniformly from  $[B, 1]$  is used to constitute  $\mathcal{E}^\mu$  at earlier stage of clinical trials, where  $B$  is the lower bound from which the random  $a_t$  can be generated, for the sake of patient's safety and acceptance.

Once the RL policies can attain satisfactory overall performance, i.e., quantified as achieving significantly improved QoCs (introduced in Sec. 3) compared to the random controller above, we consider including into  $\mathcal{E}^\mu$  the trajectories obtained from such RL policies. From this point onward, the replay buffer  $\mathcal{E}^\mu$  will be iteratively updated and enriched with the RL-induced trajectories after each trial. Consequently, the behavioral policy  $\mu$  can be considered as a mixture of random control policy and several RL policies deployed in past trials in general. With  $\mathcal{E}^\mu$  being defined, the objective for training RL policies, (20), can be optimized using gradient descent [19, 20, 39].

## 4.2 Policy Distillation

Our system design (Fig. 2) is set to process various tasks in each 2 s stimulation (i.e., control) period, facilitating communication between the research tablet and IPG, computing  $P_\beta$  from LFPs, evaluating the RL controller, data logging, and other basic functionalities that ensure the safety and functionality of DBS. Hence, it was critical to reduce the overall computation requirements, such that each task meets the required timings, as well as prolong the battery lifetime. As introduced in Sec. 2.2, the RL policies are parameterized as DNNs; although a forward pass of a DNN would not require as much computational resources as for training (through back-propagation), it may still involve hundreds of thousands of



**Figure 5: Architecture of the new deep latent sequential model (DLSM). The conditional dependencies between the variables from the posterior and sampling distributions are shown in dashed and solid lines, respectively.**

multiplication operations. For example, consider the recommended DNN size as in [39], it takes at least 120,000 multiplications to evaluate a two-layer NN with 400 and 300 nodes each. Hence, we integrate into our system the model/policy distillation techniques [26], allowing smaller sized NNs to be used to parameterize RL policies.

We build on a similar approach as in [56], originally proposed to reduce the size of DNNs used in deep Q-learning [46], which only works for a discrete action space. In particular, our extension allows for the use in the deterministic actor-critic cases considered in this work. Consider the original policy (*teacher*)  $\pi_{\theta_a}$  parameterized by a DNN with weights  $\theta_a$ . We train a smaller-sized DNN (*student*) with weights  $\tilde{\theta}_a$  to learn  $\theta_a$ 's behavior, by minimizing the mean squared error

$$\min_{\tilde{\theta}_a} \|\pi_{\theta_a}(s_t) - \pi_{\tilde{\theta}_a}(s_t)\|^2, \quad (9)$$

for all state samples contained in the experience replay  $s_t \in \mathcal{E}^\mu$ . We also consider augmenting the data used to optimize (9) to smooth out the learning process. We introduce synthetic states,  $\tilde{s}_t$ 's, where each  $\tilde{s}_t$  is generated by adding noise to each dimension of a state sample  $s_t$  that is originally in  $\mathcal{E}^\mu$ ; the noise is sampled from a zero-mean Gaussian distribution,  $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$  with  $\sigma$  being a hyper-parameter.

## 5 OPE OF DBS CONTROLLERS INCLUDING PATIENT FEEDBACK AND TREMOR DATA

As discussed in Remark 4.1, besides the reward function introduced in Sec. 4.1, for OPE we employ QoC metrics  $QoC_{grasp}$ ,  $QoC_{rate}$ , and  $QoC_{tremor}$  defined in (3), (4), (5), respectively, which are only available at the end of each session. As these well-capture performance (i.e., therapy effectiveness) of the considered policy, for OPE we additionally consider the end-of-session rewards defined as

$$\begin{aligned} r_{end} &= R_{end}(s_0, a_0, s_1, a_1, \dots, s_{T-1}, a_{T-1}, s_T) \\ &= C_2 \cdot QoC_{grasp} + C_3 \cdot QoC_{rate} - C_4 \cdot QoC_{tremor}, \end{aligned} \quad (10)$$

with  $C_2, C_3, C_4 > 0$  real constants. Without loss of generality, we slightly modify the total return under policy  $\pi$  (from Problem 2) as

$$G_0^\pi = r_{end} + \sum_{t=0}^T \gamma^t r_t, \quad (11)$$

where  $r_t$  and  $r_{end}$  follow from (8) and (10), respectively.

As discussed in Sec. 2.3, the DLMM introduced in [20], falls short in dealing with long horizons and predicting the end-of-session rewards  $r_{end}$ . To address these limitations, in this section we introduce the *deep latent sequential model* (DLSM) that directly enforces the transitions over the LVS. The overall model architecture is shown in Fig. 5. First, the latent prior  $p_\psi(z_0)$  is defined only over the initial latent variable at step  $t = 0$ ,  $z_0$ , which follows a multivariate Gaussian distribution with zero mean and identity covariance matrix.

Then, the encoder (approximated posterior) is defined over each trajectory (from  $t = 0$  to  $T$ ) as

$$q_\phi(z_{0:T}|s_{0:T}, a_{0:T-1}) = q_\phi(z_0|s_0) \prod_{t=1}^T q_\phi(z_t|z_{t-1}, a_{t-1}, s_t). \quad (12)$$

Further, the second term  $q_\phi(z_t|z_{t-1}, a_{t-1}, s_t)$ , which enforces the transitions between  $z_{t-1}$  and  $z_t$  conditioned on  $(a_{t-1}, s_t)$  and enables the encoder to capture the dynamical transitions in the LVS, can be obtained iteratively following

$$z_0^\phi \sim q_\phi(z_0|s_0), \quad h_t^\phi = f_\phi(h_{t-1}^\phi, z_{t-1}^\phi, a_{t-1}, s_t), \quad z_t^\phi \sim q_\phi(z_t|h_t^\phi); \quad (13)$$

here,  $q_\phi(z_0|s_0)$  and  $q_\phi(z_t|h_t^\phi)$  are parameterized by multivariate diagonal Gaussian distributions, each with mean and covariance determined by a feedforward DNN [6]; moreover,  $h_t^\phi$  is the hidden state of a recurrent DNN, such as long short-term memory (LSTM) [27], capturing the historical transitions among  $s_t$ ,  $a_t$  and  $z_t^\phi$  for all past steps up until  $t - 1$  within each trajectory.

The decoder (sampling distribution) is responsible for interacting with the target policies to be evaluated, from which the expected returns can be estimated as the mean return obtained by the simulated trajectories. Specifically, the decoder is defined as follows, *i.e.*,

$$p_\psi(z_{1:T}, s_{0:T}, r_{0:T-1}, r_{end}|z_0) = p_\psi(r_{end}|z_T) \cdot \prod_{t=0}^T p_\psi(s_t|z_t) \prod_{t=1}^T p_\psi(z_t|z_{t-1}, a_{t-1}) p_\psi(r_{t-1}|z_t); \quad (14)$$

here,  $p_\psi(r_{end}|z_T)$  estimates the end-of-session rewards given the latent variable at  $t = T$ ,  $z_T$ ;  $p_\psi(s_t|z_t)$ ,  $p_\psi(r_{t-1}|z_t)$  reconstruct the states and rewards;  $p_\psi(z_t|z_{t-1}, a_{t-1})$  enforces the transitions over the latent variables,  $z_t$ 's, conditioned on the actions; and  $z_0 \sim p_\psi(z_0)$  is sampled from the prior. As a result, each simulated trajectory can be generated by the decoder following

$$h_t^\psi = f_\psi(h_{t-1}^\psi, z_{t-1}^\psi, a_{t-1}), \quad z_t^\psi \sim p_\psi(z_t|h_t^\psi), \quad s_t^\psi \sim p_\psi(s_t|z_t^\psi), \\ r_{t-1}^\psi \sim p_\psi(r_{t-1}|z_t^\psi), \quad a_{t-1} \sim \pi(a_{t-1}|s_{t-1}^\psi), \quad r_{end}^\psi \sim p_\psi(r_{end}|z_T); \quad (15)$$

here,  $h_t^\psi$  is the hidden state of a recurrent DNN;  $p_\psi(z_t|h_t^\psi)$ ,  $p_\psi(s_t|z_t^\psi)$ ,  $p_\psi(r_{t-1}|z_t^\psi)$  and  $p_\psi(r_{end}|z_T)$  are multivariate diagonal Gaussians with means and covariances determined by four feedforward DNNs separately. Hence,  $s_t^\psi$ 's and  $r_{t-1}^\psi$ 's can be sampled iteratively following the process above, using the actions obtained from the target policy  $a_{t-1} \sim \pi(a_{t-1}|s_{t-1}^\psi)$  accordingly, which constitute the simulated trajectories; and  $r_{end}^\psi$  is sampled at the end of each simulated trajectory.

The theorem below derives an ELBO for the joint log-likelihood  $\log p_\psi(s_{0:T}, r_{0:T-1}, r_{end})$ , following the above DLSM architecture.

**THEOREM 5.1 (ELBO FOR DLSM).** *An ELBO of the joint log-likelihood  $\log p_\psi(s_{0:T}, r_{0:T-1}, r_{end})$  can be obtained as*

$$\mathcal{L}_{ELBO}(\psi, \phi) = \mathbb{E}_{z_t \sim q_\phi} \left[ \sum_{t=0}^T \log p_\psi(s_t|z_t) + \sum_{t=1}^T \log p_\psi(r_{t-1}|z_t) + \log p_\psi(r_{end}|z_T) - KL(q_\phi(z_0|s_0)||p(z_0)) \right. \\ \left. - \sum_{t=1}^T KL(q_\phi(z_t|z_{t-1}, a_{t-1}, s_t)||p_\psi(z_t|z_{t-1}, a_{t-1})) \right] \quad (16)$$

$$\leq \log p_\psi(s_{0:T}, r_{0:T-1}, r_{end}); \quad (17)$$

here, the first three terms are the log-likelihood of the decoder to reconstruct  $s_t$ ,  $r_{t-1}$  and  $r_{end}$  correctly, and the two terms that follow regularize the transitions captured by the encoder over the LVS, with  $KL(\cdot||\cdot)$  being the Kullback–Leibler (KL) divergence [35].

The proof of Theorem 5.1 can be found in Appendix F. Empirically, similar to the DLMM [20], the ELBO can be evaluated using the trajectories from the experience replay  $\mathcal{E}^\mu$ , by replacing the expectation as the mean over all trajectories, after which the objective  $\max_{\psi, \phi} \mathcal{L}(\psi, \phi)$  can be achieved using gradient descent [31] following the algorithm in Appendix D. Moreover, the reparameterization trick [32] is used, which allows for the gradients to be back-propagated when sampling from Gaussian distributions with means and covariances determined by DNNs. Details on reparameterization can be found in [20, 32].

## 6 CLINICAL EVALUATIONS

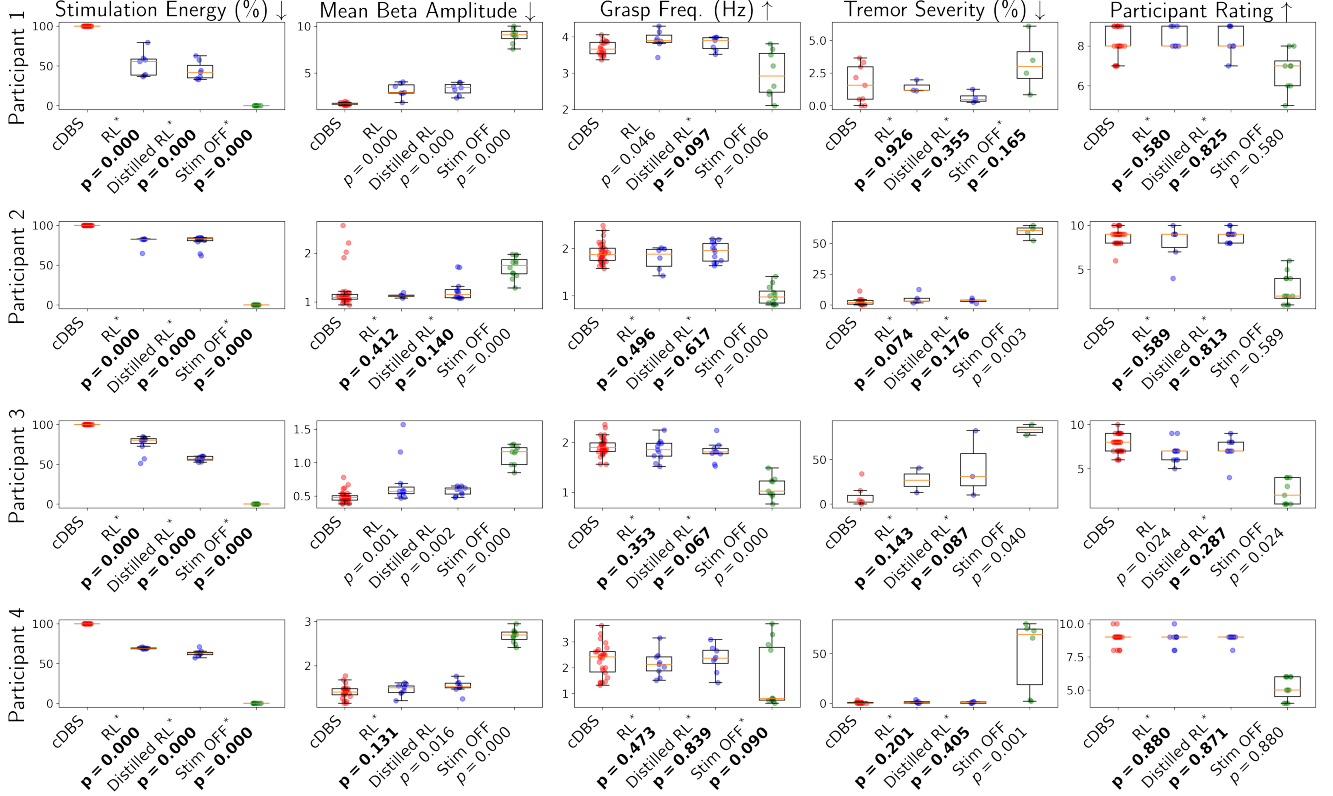
Using our closed-loop DBS system presented in Sec. 3, we evaluated the developed RL-based control framework in clinical trials on four PD patients, at Duke University Medical Center. In particular, we evaluated and compared four different types of controllers: cDBS, RL, RL with policy distillation (*i.e.*, distilled RL), and no-DBS (*i.e.*, without stimulation). The electrodes of the DBS device were placed in STN and GPi brain regions for all four participants; LFPs were sensed from STN and stimuli were delivered to both STN and GP.

Each participant also has had *different PD symptoms and severity*; their characteristics are summarized in Appendix E. All trials were conducted under close supervision of clinical experts, strictly following the process approved by the Duke University Medical Center IRB protocol complying with the obtained FDA IDE (G180280). Further, all participants provided informed written consent.

### 6.1 Therapy Efficacy and Energy-Efficiency of the RL Control Policies

We follow the offline RL and policy distillation methodology introduced in Sec. 4 to train and update (distilled) RL policies iteratively over time. Specifically, each participant had monthly clinical visit, where during each day of trials a total of 2-4 RL policies would be tested. A cDBS session was placed in between any two RL sessions as a control group. A small number of no-DBS sessions, with DBS stimulation fully off, were also tested, to validate our choice of the employed QoCs metrics – *i.e.*, whether they significantly change when the participants are not stimulated.

After each trial day was completed, the trajectories collected from all the sessions were added to the experience replay buffer  $\mathcal{E}^\mu$  unique to each participant. Between two consecutive visits of each participant, her  $\mathcal{E}^\mu$  was used to fine-tune the top-performing policies determined from the last trial (using smaller learning rates



**Figure 6: Quality of control (QoC) results from all clinical trials across participants. Wilcoxon rank-sum tests [43] between cDBS and each of the other controllers are used to test the null hypothesis that two sets of measurements are drawn from the same distribution, resulting in the  $p$ -values reported above. The null hypothesis is rejected when consider the stimulation energy consumed by both RL controllers, illustrating that they lead to significant energy reduction compared to cDBS. For all other QoCs, the null hypothesis is accepted in majority cases, showing that both RL controllers can in general attain similar control efficacy to cDBS. The controllers that lead to the acceptance/rejection of the null hypothesis in the desired direction are highlighted with asterisks and bold  $p$ -values.**

between  $[10^{-7}, 10^{-5}]$ ) or to train new policies from scratch (with learning rates between  $[10^{-5}, 10^{-3}]$ ); such policies were then tested in the next visit. We followed [39] and used two-layer NNs with 400 and 300 nodes each to parameterize the RL policies; moreover, a distilled version (student) of each corresponding full-sized RL policy (teacher) were trained as introduced in Sec. 4.2, with each represented as a two-layer NN with 20 and 10 nodes. The constants in (8) were set to  $r_a = 0, r_b = -1, C_1 = 0.3$  for all participants.

In each testing session, to evaluate the overall performance of the employed control policy, a total of 5 metrics were considered: the energy used by the IPG for stimulation, the mean beta amplitude over the session, and the 3 QoCs introduced in Sec. 3; for  $QoC_{grasp}$ , we captured the grasp frequencies of the hand that best correlates with the PD symptom for the participant (see Appendix E for details).

Fig. 6 summarizes the obtained results, and Table 1 documents the total amount of time each controller was tested in clinic. Wilcoxon rank-sum tests [43] between cDBS and each of the other controllers were used to test the null hypothesis – if two sets of measurements were drawn from the same distribution (i.e., that the controllers perform similarly over the considered metrics); from this,  $p$ -values can

	cDBS	RL	Distilled RL	No-DBS
Participant 1	84	97	97	36
Participant 2	145	80	182	52
Participant 3	135	115	115	39
Participant 4	124	119	98	48

**Table 1: Overall time, in minutes, spent toward testing each type of controller in clinical trials. Each testing session lasted 5-20 minutes, and no-DBS sessions were usually 5-min long to minimize the discomfort participants may experience.**

be calculated. The  $p$ -values accepting/rejecting the null hypothesis in the desired direction are highlighted in Fig. 6. Specifically, it can be observed that, compared to cDBS, the RL policies and their distilled version can save significant (20%-55%) stimulation energy across participants; as  $p < .05$  achieved for all participants, which rejected the null hypothesis.

When considering the other 4 metrics, there exist a great majority of results with  $p \geq .05$ , accepting the null hypothesis and indicating that both RL controllers attain control (i.e., therapy) efficacy similar to cDBS. In contrast, for the no-DBS sessions, the null hypothesis



	RL Policy (400×300 NN)	Distilled RL Policy (20×10 NN)
Mean of Computation Time	4.78 ms	2.98 ms
Std of Computation Time	32.26 ms	1.72 ms

**Table 2: Computation time of the original RL versus the distilled RL policy.**

	RL	Distilled RL	Random Controller
Battery Runtime (m)	227 ± 5	220 ± 6	247 ± 4

**Table 3: Overall battery runtime of the DBS system when the RL, distilled RL or random controllers were used.**

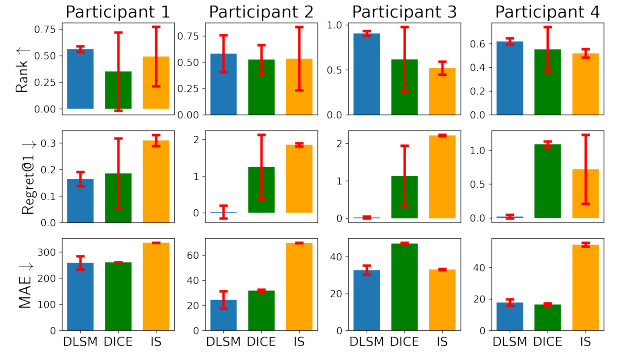
is rejected in most cases. Specifically,  $p < .05$  attained by no-DBS over the mean beta amplitude, for all participants, show that beta amplitudes can change significantly when sufficient DBS is received or not, which justify our choice of using the beta amplitudes to constitute MDP states. This also shows that the RL policies can follow the reward function (from Sec. 4.1) to effectively optimize the control strategies, with beta amplitudes also playing an important role. Consequently, the results show that both full and distilled RL policies can significantly reduce the stimulation energy, while achieving non-inferior control efficacy compared to cDBS.

**6.1.1 Computational Complexity and Overall Energy Consumption.** We also study the additional computation time and battery consumption of the DBS system due the use of full-sized RL policies or their distilled version. A Surface Go with an Intel Pentium Gold 4415Y CPU and 4GB RAM was used as the research tablet in Fig. 4. The computation time was quantified as the time needed to run a single forward pass of the NN that represents the RL policy. We evaluate the forward passes for both types of RL policies 200 times; Table 2 summarizes the mean and standard deviation of the obtained computation times. As can be seen, the distilled RL policy can be evaluated significantly faster than its counterpart.

Moreover, we quantify the overall battery consumption of the entire DBS system as the time for which the tablet or the IPG battery drains from 100% to 10% (whichever comes first). We compare the battery runtime among the full RL and distilled RL, as well as a random controller that sets the IPG to stimulate with an arbitrary amplitude in each control cycle. Each experiment was repeated 3 times, resulting in the statistics in Table 3 showing that the two RL-based controllers do not drastically shorten the runtime of the DBS system; *i.e.*, the energy used for RL-based control does not dominate the overall energy used by the DBS system.

## 6.2 Evaluation of the OPE Methodology

For each participant, a DSLM was trained following the methodology introduced in Sec. 5, and then used as a synthetic environment to interact with 6 policies trained using the deep actor-critic method (Sec. 4) with different hyper-parameters, over the buffer  $\mathcal{E}^\mu$  specific to the patient; these policies can in general lead to varying performance. Then, for each policy, the mean of total returns (11) over all simulated trajectories can be calculated, and was used to estimate the policy’s expected return from Problem 2. The constants in (10), balancing the scale of the QoCs (*i.e.*, grasp frequency, rating and tremor severity) were set to  $C_2 = C_3 = C_4 = 10$  for patients 2-4 who



**Figure 7: DSLM in general achieves higher ranks, lower regret@1’s and lower MAEs, compared to DICE and IS. Each method is trained and evaluated with 3 different random seeds, with the standard deviations shown by the error bars.**

can experience bradykinesia and pronounced tremor with insufficient DBS; in contrast, the symptoms of participant 1 are considered subtle, so we set  $C_2 = C_3 = C_4 = 25$  to better distinguish if sufficient DBS is provided; see Appendix E for details on patient characteristics as well as the dosage of PD medications.

DSLM’s performance was compared against the classic IS [54], as well as a state-of-the-art IS-based OPE method, dual-DICE [47]. Three metrics were considered to evaluate the performance of OPE, including mean absolute error (MAE), rank correlation, and regret@1, following from [14]. MAE evaluates the absolute error between the total return estimated by OPE, versus the *actual* returns, *i.e.*, mean total return recorded from clinical trials. Rank correlation quantifies the alignment between the rank of policies over OPE-estimated returns and the actual returns. Regret@1 quantifies the percentage loss, over the total actual returns, one would get by picking the policy with maximum OPE-estimated return, against the actual best-performing policy, showing if the OPE methods can identify the best-performing policy correctly. Their mathematical definitions can be found in Appendix G.

The obtained results are summarized in Fig. 7. As shown, the DSLM in general achieved significantly higher rank and lower regret, as well as non-inferior MAE, over DICE and IS.

## 7 CONCLUSION

In this paper, we introduced an offline RL and OPE framework to design and evaluate closed-loop DBS controllers using only historical data. Moreover, a policy distillation method was introduced to further reduce the computation requirements for evaluating RL policies. The control efficacy and energy efficiency of the RL controllers were validated with clinical testing over 4 patients. Results showed that RL-based controllers lead to similar control efficacy as cDBS, but with significantly reduced stimulation energy. The computation times for the RL and distilled RL controllers were compared, showing that the distilled version executed significantly faster; future work will focus on further reducing execution times of the distilled RL controllers to match capabilities of implanted devices. Finally, the DSLM is trained to estimate the expected returns of RL policies, which outperforms existing IS-based OPE methods, in terms of rank correlations, regrets and MAEs.

## REFERENCES

- [1] Mattia Arlotti, Manuela Rosa, et al. 2016. The adaptive deep brain stimulation challenge. *Parkinsonism & related disorders* 28 (2016), 12–17.
- [2] Mattia Arlotti, Lorenzo Rossi, et al. 2016. An external portable device for adaptive deep brain stimulation (aDBS) clinical research in advanced Parkinson's Disease. *Medical engineering & physics* 38, 5 (2016), 498–505.
- [3] Alim Louis Benabid. 2003. Deep brain stimulation for Parkinson's disease. *Current opinion in neurobiology* 13, 6 (2003), 696–706.
- [4] Aleksandar Beric, Patrick J Kelly, et al. 2001. Complications of deep brain stimulation surgery. *Stereotactic and functional neurosurgery* 77, 1–4 (2001), 73–78.
- [5] M Beudel and P Brown. 2016. Adaptive deep brain stimulation in Parkinson's disease. *Parkinsonism & related disorders* 22 (2016), S123–S126.
- [6] Christopher Bishop. 2006. *Pattern recognition and machine learning*. Springer.
- [7] Peter Brown, Antonio Oliviero, et al. 2001. Dopamine dependency of oscillations between subthalamic nucleus and pallidum in Parkinson's disease. *Journal of Neuroscience* 21, 3 (2001), 1033–1038.
- [8] A H Butt, E Rovini, et al. 2018. Objective and automatic classification of Parkinson disease with Leap Motion controller. *Biomedical engineering* 17, 1 (2018), 1–21.
- [9] Witney Chen, Lowry Kirkby, et al. 2021. The role of large-scale data infrastructure in developing next-generation deep brain stimulation therapies. *Frontiers in Human Neuroscience* 15 (2021), 717401.
- [10] Bo Dai, Ofir Nachum, et al. 2020. Coincidence: Off-policy confidence interval estimation. *arXiv preprint arXiv:2010.11652* (2020).
- [11] Lonneke ML De Lau and Monique MB Breteler. 2006. Epidemiology of Parkinson's disease. *The Lancet Neurology* 5, 6 (2006), 525–535.
- [12] Günther Deuschl, Carmen Schade-Brittinger, et al. 2006. A randomized trial of deep-brain stimulation for Parkinson's disease. *New England Journal of Medicine* 355, 9 (2006), 896–908.
- [13] Kenneth A Follett, Frances M Weaver, et al. 2010. Pallidal versus subthalamic deep-brain stimulation for Parkinson's disease. *New England Journal of Medicine* 362, 22 (2010), 2077–2091.
- [14] Justin Fu, Mohammad Norouzi, et al. 2020. Benchmarks for Deep Off-Policy Evaluation. In *ICLR*.
- [15] Ge Gao, Qitong Gao, et al. 2022. A Reinforcement Learning-Informed Pattern Mining Framework for Multivariate Time Series Classification. In *IJCAI*.
- [16] Ge Gao, Song Ju, Markel Sanz Ausin, and Min Chi. 2023. Hope: Human-centric off-policy evaluation for e-learning and healthcare. In *AAMAS*.
- [17] Qitong Gao, Ge Gao, Min Chi, and Miroslav Pajic. 2023. Variational Latent Branching Model for Off-Policy Evaluation. In *ICLR*.
- [18] Qitong Gao, Davood Hajinezhad, et al. 2019. Reduced Variance Deep Reinforcement Learning with Temporal Logic Specifications. In *ICCPs*. ACM.
- [19] Qitong Gao, Michael Naumann, et al. 2020. Model-Based Design of Closed Loop Deep Brain Stimulation Controller using Reinforcement Learning. In *2020 ACM/IEEE 11th Int. Conf. on Cyber-Physical Systems (ICCPs)*. IEEE, 108–118.
- [20] Qitong Gao, Stephen L Schmidt, et al. 2022. Offline Policy Evaluation for Learning-based Deep Brain Stimulation Controllers. In *2022 ACM/IEEE 13th International Conference on Cyber-Physical Systems (ICCPs)*. IEEE, 80–91.
- [21] Qitong Gao, Dong Wang, et al. 2022. Gradient Importance Learning for Incomplete Observations. In *International Conference on Learning Representations*.
- [22] Shixiang Gu, Ethan Holly, Timothy Lillicrap, and Sergey Levine. 2017. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In *Int. Conf. on robotics and automation (ICRA)*. IEEE, 3389–3396.
- [23] A. Guez, R. D. Vincent, M. Avoli, and J. Pineau. 2008. Adaptive Treatment of Epilepsy via Batch-mode Reinforcement Learning. In *AAAI*. 1671–1678.
- [24] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *ICML*. PMLR, 1861–1870.
- [25] J. Habets, M. Heijmans, et al. 2018. An update on adaptive deep brain stimulation in Parkinson's disease. *Movement Disorders* 33, 12 (2018), 1834–1843.
- [26] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, and others. [n. d.]. Distilling the knowledge in a neural network. ([n. d.]).
- [27] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [28] S. Ishii, W. Yoshida, and J. Yoshimoto. 2002. Control of exploitation–exploration meta-parameter in reinforcement learning. *Neural networks* 15 (2002), 665–687.
- [29] Nan Jiang and Lihong Li. 2016. Doubly robust off-policy value evaluation for reinforcement learning. In *ICML*. PMLR, 652–661.
- [30] Ilija Jovanov, Michael Naumann, et al. 2018. Platform for model-based design and testing for deep brain stimulation. In *ICCPs*.
- [31] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [32] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).
- [33] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. 2022. Offline Reinforcement Learning with Implicit Q-Learning. In *ICLR*.
- [34] A.A. Kühn, A. Kupsch, GH. Schneider, and P Brown. 2006. Reduction in subthalamic 8–35 Hz oscillatory activity correlates with clinical improvement in Parkinson's disease. *Euro. J. of Neuroscience* 23, 7 (2006), 1956–1960.
- [35] Solomon Kullback and Richard A Leibler. 1951. On information and sufficiency. *The annals of mathematical statistics* 22, 1 (1951), 79–86.
- [36] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. 2020. Conservative q-learning for offline reinforcement learning. In *NeurIPS*.
- [37] Alexis M Kuncel and Warren M Grill. 2004. Selection of stimulus parameters for deep brain stimulation. *Clinical neurophysiology* 115, 11 (2004), 2431–2441.
- [38] Alex X Lee, Anusha Nagabandi, Pieter Abbeel, and Sergey Levine. 2020. Stochastic latent actor-critic: Deep reinforcement learning with a latent variable model. *Advances in Neural Information Processing Systems* 33 (2020), 741–752.
- [39] Timothy P Lillicrap, Jonathan J Hunt, et al. 2016. Continuous control with deep reinforcement learning. *ICLR* (2016).
- [40] Simon Little, Alex Pogosyan, et al. 2013. Adaptive deep brain stimulation in advanced Parkinson disease. *Annals of neurology* 74, 3 (2013), 449–457.
- [41] Simon Little, Elina Tripoliti, et al. 2016. Adaptive deep brain stimulation for Parkinson's disease demonstrates reduced speech side effects compared to conventional stimulation in the acute setting. *J Neurol Neurosurg Psychiatry* 87, 12 (2016), 1388–1389.
- [42] Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. 2018. Breaking the Curse of Horizon: Infinite-Horizon Off-Policy Estimation. In *NeurIPS*.
- [43] Henry B Mann and Donald R Whitney. 1947. On a test of whether one of two random variables is stochastically larger than the other. *The annals of mathematical statistics* (1947), 50–60.
- [44] C Marras, JC Beck, et al. 2018. Prevalence of Parkinson's disease across North America. *NPJ Parkinson's disease* 4, 1 (2018), 21.
- [45] Volodymyr Mnih, Adria Puigdomenech Badia, et al. 2016. Asynchronous methods for deep reinforcement learning. In *ICML*. 1928–1937.
- [46] Volodymyr Mnih, Koray Kavukcuoglu, et al. 2015. Human-level control through deep reinforcement learning. *Nature* 518, 7540 (2015), 529.
- [47] Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. 2019. DualDICE: Behavior-agnostic estimation of discounted stationary distribution corrections. *NeurIPS* 32 (2019).
- [48] Vivek Nagaraj, Andrew Lamperski, and Theoden I Netoff. 2017. Seizure control in a computational model using a reinforcement learning stimulation paradigm. *International J. of Neural Sys.* 27, 07 (2017), 1750012.
- [49] Michael S Okun. 2012. Deep-brain stimulation for Parkinson's disease. *New England Journal of Medicine* 367, 16 (2012), 1529–1538.
- [50] Enrico Opri, Stephanie Cernera, et al. 2020. Chronic embedded cortico-thalamic closed-loop deep brain stimulation for the treatment of essential tremor. *Science translational medicine* 12, 572 (2020), eaay7680.
- [51] Bahram Parvini, Christopher Scully, et al. 2018. Regulatory considerations for physiological closed-loop controlled medical devices used for automated critical care: food and drug administration workshop discussion topics. *Anesthesia and analgesia* 126, 6 (2018), 1916.
- [52] J. Pineau, A. Guez, et al. 2009. Treating epilepsy via adaptive neurostimulation: a reinforcement learning approach. *Int. J. of Neural Sys.* 19, 04 (2009), 227–240.
- [53] Rob Powers, Maryam Etezadi-Amoli, et al. 2021. Smartwatch inertial sensors continuously monitor real-world motor fluctuations in Parkinson's disease. *Science translational medicine* 13, 579 (2021), eabd7865.
- [54] Doina Precup. 2000. Eligibility traces for off-policy policy evaluation. *Computer Science Department Faculty Publication Series* (2000), 80.
- [55] Claudia Ramaker, Johan Marinus, Anne Margarethe Stiggelbout, and Bob Johannes Van Hilten. 2002. Systematic evaluation of rating scales for impairment and disability in Parkinson's disease. *Movement disorders* 17, 5 (2002), 867–876.
- [56] Andrei A Rusu, Sergio G Colmenarejo, et al. 2016. Policy Distillation. In *ICLR*.
- [57] David Silver, Guy Lever, et al. 2014. Deterministic policy gradient algorithms.
- [58] Rosa Q So, Alexander R Kent, and Warren M Grill. 2012. Relative contributions of local cell and passing fiber activation and silencing to changes in thalamic fidelity during deep brain stimulation and lesioning: a computational modeling study. *Journal of computational neuroscience* 32, 3 (2012), 499–519.
- [59] Scott Stanslaski, Jeffrey Herron, et al. 2018. A chronically implantable neural coprocessor for investigating the treatment of neurological disorders. *IEEE transactions on biomedical circuits and systems* 12, 6 (2018), 1230–1245.
- [60] Nicole C Swann, Coralie de Hemptinne, et al. 2016. Gamma oscillations in the hyperkinetic state detected with chronic human brain recordings in Parkinson's disease. *Journal of Neuroscience* 36, 24 (2016), 6445–6458.
- [61] Ziyang Tang, Yihao Feng, et al. 2019. Doubly Robust Bias Reduction in Infinite Horizon Off-Policy Estimation. In *ICLR*.
- [62] Philip Thomas and Emma Brunskill. 2016. Data-efficient off-policy policy evaluation for reinforcement learning. In *ICML*. PMLR, 2139–2148.
- [63] Joshua K Wong, Günther Deuschl, et al. 2022. Proc. the 9th Annual Deep Brain Stimulation Think Tank: Advances in Cutting Edge Technologies, Artificial Intelligence, Neuromodulation, Neuroethics, Pain, Interventional Psychiatry, Epilepsy, and Traumatic Brain Injury. *Frontiers in Human Neuroscience* (2022), 25.
- [64] Yuhuai Wu, Elman Mansimov, et al. 2017. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. In *NeurIPS*.
- [65] Mengjiao Yang, Ofir Nachum, et al. 2020. Off-Policy Evaluation via the Regularized Lagrangian. In *NeurIPS*, Vol. 33.

## A NOTATION

Here, we define the notation used in the paper. The sets of reals, integers, and positive integers, are denoted by  $\mathbb{R}$ ,  $\mathbb{Z}$ , and  $\mathbb{Z}^+$ , respectively. Further,  $x \sim p(x)$  denotes that random variable  $x$  is sampled from distribution  $p(x)$ . We also use  $\mathcal{N}(x; \mu, \Sigma)$  to denote Gaussian distributions with mean  $\mu$  and covariance matrix  $\Sigma$  over variable  $x$ . For simplicity, we write  $x \sim \mathcal{N}(\mu, \Sigma)$  during sampling. The KL-divergence between distributions  $p(x)$  and  $q(x)$  is defined as

$$KL(p||q) = \mathbb{E}_p \left[ \log \frac{q(x)}{p(x)} \right]. \quad (18)$$

## B AVAILABILITY OF DATA AND CODE

To facilitate research in developing RL-based DBS controllers, we plan to open-source the data collected from the clinical testing, as well as the implementation for training RL policies. The RC+S system as well as its Summit code base are considered proprietary, which may not be published online. The implementation of our OPE method, DLMM, is built on our previous works [17, 20], with code published at <https://github.com/gaoqitong/vlbm>.

## C ADDITIONAL PRELIMINARIES

Below, we introduce additional details on preliminaries discussed Section 2.

### C.1 Deep Actor-Critic RL

We now briefly introduce the deep actor-critic algorithm [39] and refer the readers to [19, 20, 39] for more details. First, the state-action value functions can be defined as follows.

**DEFINITION C.1 (STATE-ACTION VALUE FUNCTION).** *Given an MDP  $\mathcal{M}$  and policy  $\pi$ , the state-action value function  $Q^\pi(s, a)$ , where  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ , is defined as the expected return for taking action  $a$  when at state  $s$  following policy  $\pi$  at stage  $t$ , i.e.,*

$$Q^\pi(s, a) = \mathbb{E}_{s \sim \mathcal{S}, a \sim \mathcal{A}} [G_t | s_t = s, a_t = a]. \quad (19)$$

Two neural networks, with weights  $\theta_a$  and  $\theta_c$ , can be used to parameterize the policy (actor)  $\pi_{\theta_a}(s) : \mathcal{S} \rightarrow \mathcal{A}$  and the Q-functions (critic)  $Q_{\theta_c}(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , respectively. Finally, the target policy  $\pi^* = \pi_{\theta_a^*}$  can be obtained by optimizing over

$$\max_{\theta_a, \theta_c} \mathbb{E}_{s, a, r, s' \sim \mathcal{E}^\mu} [Q_{\theta_c}(s, \pi_{\theta_a}(s))]; \quad (20)$$

this can be achieved using gradient descent, over all the training samples in the experience replay buffer  $\mathcal{E}^\mu$  [39].

### C.2 Deep Latent MDP Model (DMLL)

The DLMM is trained to fit the transitions of the MDP  $p(s_{t+1}|s_t, a_t)$  and rewards  $r_t = R(s_t, a_t)$ , which consist of three components, i.e., a latent prior, posterior and sampling distribution. The prior  $p_\psi(z_t)$  is parameterized by  $\psi$ , over the latent variable space (LVS)  $\mathcal{Z} \subset \mathbb{R}^d$ , where  $d \in \mathbb{Z}^+$  is the dimension. The prior represents one's belief over the latent distribution of the states (i.e., probability density function over the latent variables) which is considered unknown; thus, it is usually chosen to be a multivariate Gaussian with zero mean and identity covariance. Then, the encoder (or approximated posterior)  $q_\phi(z_t|s_t)$  is parameterized by  $\phi$ , which is responsible for encoding the MDP state  $s_t \in \mathcal{S}$  into the LVS,  $\mathcal{Z}$ .

## Algorithm 1 Train DLMM.

**Input:** Model weights  $\psi, \phi$ , experience replay buffer  $\mathcal{E}^\mu$ , and learning rate  $\alpha$ .

**Begin:**

- 1: Initialize  $\psi, \phi$
- 2: **for**  $iter$  in  $1 : max\_iter$  **do**
- 3:   Sample a trajectory  $[(s_0, a_0, r_0, s_1), \dots, (s_{T-1}, a_{T-1}, r_{T-1}, s_T)] \sim \mathcal{E}^\mu$
- 4:    $z_0^\phi \sim q_\phi(z_0|s_0)$
- 5:    $z_0^\psi \sim p_\psi(z_0)$
- 6:   Run forward pass of DLMM following (13) and (15) for  $t = 1 : T$ , and collect all variables needed to evaluate the all terms within the expectation in  $\mathcal{L}_{ELBO}$ , which is denoted as  $\hat{\mathcal{L}}_{ELBO}$ .
- 7:    $\psi \leftarrow \psi + \alpha \nabla_\psi \hat{\mathcal{L}}_{ELBO}$
- 8:    $\phi \leftarrow \phi + \alpha \nabla_\phi \hat{\mathcal{L}}_{ELBO}$
- 9: **end for**

Note that the true posterior  $p_\psi(z_t|s_t)$  is intractable, since its density function contains integration over  $\mathcal{Z}$  which is deemed unknown; we refer to [20] for more details. Lastly, the decoder (or sampling distribution)  $p_\psi(s_{t+1}, r_t | z_t, a_t)$  enforces the MDP transition from  $t$  to  $t+1$ .

Hence, the DLMM can be used to interact with the policy  $\pi$  and generate simulated trajectories via

$$z_t \sim q_\phi(z_t|\hat{s}_t), \quad \hat{s}_{t+1}, \hat{r}_t \sim p_\psi(s_{t+1}, r_t | z_t, a_t); \quad (21)$$

here,  $\hat{s}_t, \hat{s}_{t+1}$  and  $\hat{r}_t$  represent the states and rewards predicted from DLMM. Consequently, the expected return of  $\pi$  can be estimated as  $\frac{1}{M} \sum_{i=1}^M \sum_{t=0}^T \gamma^t \hat{r}_t^{(i)}$ , where  $M$  is the total number of simulated trajectories generated following the process above, and  $\hat{r}_t^{(i)}$  is the predicted reward at step  $t$  in the  $i$ -th simulated trajectory. To train DLMM, one can maximize the evidence lower bound (ELBO) of the joint log-likelihood  $\sum_{t=0}^T \log p_\psi(s_{t+1}, r_t)$ , where the derivation of the ELBO for DLMM can be found in [20].

From (21), it can be observed that the predicted  $\hat{s}_{t+1}, \hat{r}_t$  are conditioned on  $z_t$ 's, which are dependent on the predicted state  $\hat{s}_t$  from the last step. As a result, such an iterative process may not scale well to environments with longer horizons and more complicated dynamics, as the prediction error from all earlier steps are propagated into the future steps. Moreover, this DLMM is not capable of predicting the QoC metrics that are only evaluated *once at the end of each session*, including the bradykinesia results, patient ratings and tremor severity as discussed in Sec. 3. To address such limitations, we introduce a new latent modeling method in Sec. 5, which decouples the dependencies between  $z_t$ 's and  $\hat{s}_t$ 's by directly enforcing the temporal transitions over latent variables, i.e.,  $p(z_{t+1}|z_t, a_t)$ .

## D ALGORITHM TO TRAIN DLMM

Here we introduce how to use gradient descent to maximize the ELBO (16), resulting in Algorithm 1. For simplicity, we first illustrate with the case where the training batch only contains a single trajectory, and then extend to the cases where each batch contain  $n$  trajectories. In each iteration, a trajectory is sampled from the experience replay buffer  $\mathcal{E}^\mu$ . Then, the initial latent state in the encoder is obtained following  $z_0^\phi \sim q_\phi(z_0|s_0)$ , while the initial latent state for the sampling distribution is generated following the latent

prior  $z_0^\psi \sim p_\psi(z_0)$ . The processes introduced in (13) can be used to generate  $z_t^\psi$ 's iteratively given  $z_0^\psi$ . Similarly,  $s_t^\psi, r_t^\psi, z_t^\psi$  can be generated iteratively following (15). As a result, the log-likelihoods and KL-divergence terms within the expectation in  $\mathcal{L}_{ELBO}$ , defined in (16), can be evaluated using the variables above, after which  $\psi, \phi$  can be updated using the gradients  $\nabla_\psi \tilde{\mathcal{L}}_{ELBO}, \nabla_\phi \tilde{\mathcal{L}}_{ELBO}$ , respectively, where  $\tilde{\mathcal{L}}_{ELBO}$  refers to all the terms within the expectation in  $\mathcal{L}_{ELBO}$ . This algorithm is summarized in Alg. 1.

To extend to batch gradient descent, in line 3 of Algorithm 1, a batch of  $n$  trajectories,  $\mathcal{B}(n)$ , will be sampled, i.e.,

$$\mathcal{B}(n) = \left[ \left[ (s_0^{(1)}, a_0^{(1)}, r_0^{(1)}, s_1^{(1)}), \dots, (s_{T-1}^{(1)}, a_{T-1}^{(1)}, r_{T-1}^{(1)}, s_T^{(1)}) \right], \dots, \left[ (s_0^{(n)}, a_0^{(n)}, r_0^{(n)}, s_1^{(n)}), \dots, (s_{T-1}^{(n)}, a_{T-1}^{(n)}, r_{T-1}^{(n)}, s_T^{(n)}) \right] \right] \sim \mathcal{E}^\mu. \quad (22)$$

Then, the processes illustrated in lines 4-6 in Algorithm 1 can be executed in  $n$  parallel threads, with each corresponding to a unique trajectory in  $\mathcal{B}(n)$ . Further, then,  $\tilde{\mathcal{L}}_{ELBO}$  can be evaluated as

$$\begin{aligned} \tilde{\mathcal{L}}_{ELBO}(\psi, \phi) = & \frac{1}{n} \sum_{i=1}^n \left[ \sum_{t=0}^T \log p_\psi(s_t^{(i)} | z_t^{(i)}) \right. \\ & + \sum_{t=1}^T \log p_\psi(r_{t-1}^{(i)} | z_t^{(i)}) \\ & + \log p_\psi(r_{end}^{(i)} | z_T^{(i)}) - KL(q_\phi(z_0^{(i)} | s_0^{(i)}) || p(z_0^{(i)})) \\ & \left. - \sum_{t=1}^T KL(q_\phi(z_t^{(i)} | z_{t-1}^{(i)}, a_{t-1}^{(i)}, s_t^{(i)}) || p_\psi(z_t^{(i)} | z_{t-1}^{(i)}, a_{t-1}^{(i)})) \right], \end{aligned} \quad (23)$$

with  $s_t^{(i)}, a_t^{(i)}, z_t^{(i)}, r_t^{(i)}, r_{end}^{(i)}$  being the variables involved in one of the threads above processing the  $i$ -th trajectory in  $\mathcal{B}(n)$ .

## E PARTICIPANT CHARACTERISTICS

*Participant 1.* Episodes of tremor only. Bradykinesia in the left hand correlates with right STN beta amplitude. STN beta amplitudes in both hemispheres correlate with stimulation amplitude. This participant takes 2 tablets of Sinemet 25mg/100mg for each clinical visit, one at 2 hrs before the testing begins (~7am) and another in the middle of the visit (around noon), respectively.

*Participant 2.* Very large amplitude tremor returns within seconds of low amplitude DBS. Bradykinesia in right hand highly correlated with left STN beta amplitude. Left STN beta amplitude also correlated with DBS amplitude. This participant takes Flexeril 10mg, Selegiline 5mg and Pramipexole 0.125mg (1 tablet for each) at 2 hrs before the testing begins.

*Participant 3.* Pronounced tremor (hands and jaw) returns within seconds of low amplitude closed-loop DBS. Bradykinesia in right hand correlates with left STN beta amplitude. Left STN beta amplitude is the most responsive to stimulation of the cohort. This participant takes 2 tablets of Sinemet 25mg/100mg for each clinical visit, one at 2 hrs before the testing begins (~7am) and another in the middle of the visit (around noon), respectively.

*Participant 4.* Pronounced tremor. Monopolar left STN DBS can produce a dyskinesia in neck, so clinical settings have used a bipolar configuration. Right hand bradykinesia correlated to left STN

beta power. This participant takes 1 tablets of Sinemet 25mg/100mg at 2 hrs before the testing begins (~7am).

## F PROOF OF THEOREM 5.1

We now derive the evidence lower bound (ELBO) for the joint log-likelihood distribution, i.e.,

$$\log p_\psi(s_{0:T}, r_{0:T-1}, r_{end}) \quad (24)$$

$$= \log \int_{z_{1:T} \in \mathcal{Z}} p_\psi(s_{0:T}, z_{1:T}, r_{0:T-1}, r_{end}) dz \quad (25)$$

$$= \log \int_{z_{1:T} \in \mathcal{Z}} \frac{p_\psi(s_{0:T}, z_{1:T}, r_{0:T-1}, r_{end})}{q_\phi(z_{0:T} | s_{0:T}, a_{0:T-1})} q_\phi(z_{0:T} | s_{0:T}, a_{0:T-1}) dz \quad (26)$$

$$\geq \mathbb{E}_{z_t \sim q_\phi} [\log p(z_0) + \log p_\psi(s_{0:T}, z_{1:T}, r_{0:T-1} | z_0) + \log p_\psi(r_{end} | z_T) - \log q_\phi(z_{0:T} | s_{0:T}, a_{0:T-1})] \quad (27)$$

$$\begin{aligned} = & \mathbb{E}_{z_t \sim q_\phi} \left[ \log p(z_0) + \log p_\psi(s_0 | z_0) + \log p_\psi(r_{end} | z_T) \right. \\ & + \sum_{t=1}^T \log p_\psi(s_t, z_t, r_{t-1} | z_{t-1}, a_{t-1}) \\ & \left. - \log q_\phi(z_0 | s_0) - \sum_{t=1}^T \log q_\phi(z_t | z_{t-1}, a_{t-1}, s_t) \right] \end{aligned} \quad (28)$$

$$\begin{aligned} = & \mathbb{E}_{z_t \sim q_\phi} \left[ \log p(z_0) - \log q_\phi(z_0 | s_0) + \log p_\psi(s_0 | z_0) + \log p_\psi(r_{end} | z_T) \right. \\ & + \sum_{t=1}^T \log (p_\psi(s_t | z_t) p_\psi(r_{t-1} | z_t) p_\psi(z_t | z_{t-1}, a_{t-1})) \\ & \left. - \sum_{t=1}^T \log q_\phi(z_t | z_{t-1}, a_{t-1}, s_t) \right] \end{aligned} \quad (29)$$

$$\begin{aligned} = & \mathbb{E}_{z_t \sim q_\phi} \left[ \sum_{t=0}^T \log p_\psi(s_t | z_t) + \log p_\psi(r_{end} | z_T) \right. \\ & + \sum_{t=1}^T \log p_\psi(r_{t-1} | z_t) - KL(q_\phi(z_0 | s_0) || p(z_0)) \\ & \left. - \sum_{t=1}^T KL(q_\phi(z_t | z_{t-1}, a_{t-1}, s_t) || p_\psi(z_t | z_{t-1}, a_{t-1})) \right]. \end{aligned} \quad (30)$$

Note that the transition from (26) to (27) follows Jensen's inequality.

## G OPE METRICS

*Rank correlation.* Rank correlation measures the Spearman's rank correlation coefficient between the ordinal rankings of the estimated returns and actual returns across policies, i.e.,

$$\rho = \frac{\text{Cov}(\text{rank}(V_{1:P}^\pi), \text{rank}(\hat{V}_{1:P}^\pi))}{\sigma(\text{rank}(V_{1:P}^\pi)) \sigma(\text{rank}(\hat{V}_{1:P}^\pi))}, \text{ where } \text{rank}(V_{1:P}^\pi) \text{ is the ordinal rankings of the actual returns, and } \text{rank}(\hat{V}_{1:P}^\pi) \text{ is the ordinal rankings of the OPE-estimated returns.}$$

*Regret@1.* Regret@1 is the (normalized) difference between value of the actual best policy, against value of the policy associated with the best OPE-estimated return, which is defined as  $(\max_{i \in 1:P} V_i^\pi - \max_{j \in \text{best}(1:P)} V_j^\pi) / \max_{i \in 1:P} V_i^\pi$  where  $\text{best}(1:P)$  denotes the index of the best policy over the set of  $P$  policies as measured by estimated values  $\hat{V}^\pi$ .

*Mean Absolute error (MAE).* MAE is defined as the absolute difference between the actual return and estimated return of a policy:  $MAE = |V^\pi - \hat{V}^\pi|$ ; here,  $V^\pi$  is the actual value of the policy  $\pi$ , and  $\hat{V}^\pi$  is the estimated value of  $\pi$ .