

---

# DeepBlip: Estimating Conditional Average Treatment Effects Over Time

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Estimating the conditional average treatment effect (CATE) over time is crucial for  
2 making personalized decisions in medicine. Yet, existing neural methods for this  
3 task have limitations: they either (1) do not adjust for time-varying confounding  
4 and are thus biased (e.g., causal transformer), or (2) become unstable over long  
5 time horizons because the method has to learn the full counterfactual outcome  
6 trajectories (e.g., MSNs, G-computation). To address these limitations, we propose  
7 DeepBlip, the first neural framework that leverages the blip function from structural  
8 nested mean models to break the joint effect of treatment sequences over time into  
9 localized, time-specific “blip effects”. As a result, we learn a simpler estimand that  
10 does not require learning full counterfactual outcome trajectories, which is thus  
11 more stable over long horizons. Further, our DeepBlip adjusts for time-varying  
12 confounding and is thus unbiased. Our DeepBlip seamlessly integrates sequential  
13 models like LSTMs or transformers to capture complex temporal dependencies.  
14 Our DeepBlip has two further strengths for medical practice: (i) The loss is Neyman-  
15 orthogonal, meaning it is robust against model misspecification. (ii) The blip effects  
16 can be used to predict treatment effects for new treatment sequences without re-  
17 computation, which allows to identify optimal treatment sequences through offline  
18 evaluation. Finally, we evaluate our DeepBlip across various clinical datasets,  
19 where it achieves state-of-the-art performance.

## 20 1 Introduction

21 Predicting the effects of treatment sequences is crucial for personalized medicine to choose the  
22 best therapeutic strategy for a patient based on their history [9]. Methodologically, the **conditional**  
23 **average treatment effect (CATE) over time** captures the combined effect of multiple treatments  
24 in the next  $\tau$  time steps (see Fig. 1). Nowadays, CATEs over time are frequently estimated from  
25 observational data with patient histories, such as the electronic health records [2, 4].

26 Several works aim to estimate CATE over time from observational data, but they suffer from two key  
27 limitations: ① **No proper adjustment for confounding and thus bias:** There are methods that do  
28 *not* properly adjust for time-varying confounding (e.g., CRN [3], causal transformer [20]) and that are  
29 thus *biased*. This leads to unreliable estimates, which is particularly problematic for safety-critical  
30 applications such as personalized medicine. ② **Unstable for long time horizons:** Other methods  
31 require modeling the *full* counterfactual outcome trajectories. This is the case in MSNs, which must  
32 learn long-range treatment-response mappings (e.g., as RMSNs [19]), or g-computation, which relies  
33 on modeling the full data-generating process of covariates and outcomes (e.g., G-transformer [14]).  
34 To the best of our knowledge, there is **no** method for estimating CATE over time that has addressed  
35 both challenges ① and ②.

In principle, one way to address the above limitations is through the theoretical framework of **structural nested mean models** (SNMMs) [28, 29]. SNMMs provide a principled foundation for estimating CATEs over time in an *unbiased* way. For this, SNMMs decompose the time-varying CATE into a sequence of *incremental treatment effects*, formalized through so-called **blip functions**. This decomposition yields several important advantages: (ii) It enables a divide-and-conquer approach that breaks the CATE over time into localized, time-specific causal effects. As a result, SNMMs define an estimand that is easier to learn than in many other methods (e.g., MSNs) and thus avoid the need to model full counterfactual trajectories. (ii) Because blip functions are conditionally independent across time given a patient’s history, estimation errors do *not* propagate, which makes them more stable for long time horizons. *However*, SNMMs are *only* a theoretical foundation (and, therefore, *not* a model that can be directly applied). So far, one study by Lewis et al. [18] has employed SNMMs, yet only instantiated using linear models. To the best of our knowledge, no prior work has developed a neural version of SNMMs.

Here, we propose **DeepBlip**, the first *neural* framework to estimate CATE over time by leveraging the blip function from SNMMs. DeepBlip decomposes the joint effect of treatment sequences over time into localized, time-specific blip effects, which enables more tractable and stable learning. This allows our DeepBlip to overcome both of the two limitations from above: (1) Our DeepBlip adjusts for time-varying confounding and is thus *unbiased*. (2) Our DeepBlip targets a simpler estimand than many of the above methods, thereby avoiding the need to learn the full counterfactual outcome trajectories and which improves the stability of DeepBlip over long time horizons. Our DeepBlip is built on top of sequential neural networks (e.g., LSTMs, transformers) to capture complex temporal dependencies. For this, it employs a two-stage architecture: Stage 1 models the probability of time-varying treatments and mean outcomes conditioned on a patient’s history, while Stage 2 reformulates g-estimation [28, 29] as a risk minimization task to directly learn the blip functions.

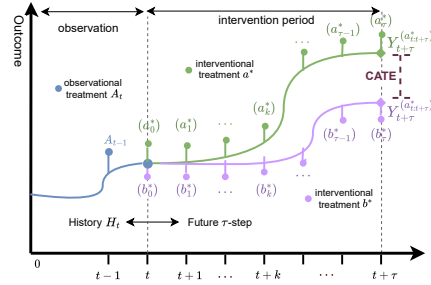


Figure 1: **CATE over time.** Trajectories of potential outcomes under two interventional sequences  $a_{t:t+\tau}^*$ ,  $b_{t:t+\tau}^*$  given the shared observed history  $H_t$ . The difference between the two curves is the CATE over time.

Our DeepBlip has two further strengths for medical practice: (i) The loss *Neyman-orthogonal*, meaning it is *robust against model misspecification*. (ii) The learned blip effects can be reused to predict treatment outcomes for new treatment sequences *without* re-computation. This enables an efficient approach for *offline evaluation* of different therapeutic strategies. Formally, at inference time, DeepBlip can identify the optimal treatment sequence within just one forward pass. This is unlike other methods, which typically require either re-training [32, 14] or multiple forward passes due to exhaustive search [19, 3, 20]. This property is especially beneficial for clinicians when searching for optimal treatment sequences.

Our **contributions** are three-fold:<sup>1</sup> (1) We introduce the first neural framework to predict CATE over time via the SNMM framework. (2) Our framework is carefully tailored to medical practice by benefiting from robustness due to Neyman-orthogonality and efficient offline evaluation. (3) We conduct extensive experiments across multiple medical datasets to demonstrate that our DeepBlip is effective and also robust across long time horizons.

## 2 Related Work<sup>2</sup>

**Estimating CATE in the static setting:** There has been extensive research in estimating ATE/CATE in the static setting (e.g., [1, 7, 17, 34, 39, 43, 17]). Recently, deep learning has been used to improve the non-parametric estimation of ATE/CATE [43]. However, these methods are aimed at static settings and thus struggle with medical datasets such as electronic health records, where patient histories are recorded *over time*.

<sup>1</sup>Code for review is available at <https://github.com/anonymous> Upon acceptance, we will move our code to a public GitHub repository.

<sup>2</sup>We provide an extended related work in Appendix E

**Estimating ATE over time:** One line of work has developed methods for the ATE over time [10, 35]. However, the ATE captures only population-level effects and thus overlooks differences in treatment effectiveness across patients. In contrast, we focus on the CATE, which provides a more granular, individualized estimate of treatment outcomes, which is highly relevant for personalized medicine [9].

**Estimating CATE over time:** There are several *neural* methods for this task<sup>3</sup>, which can be broadly categorized into two streams but with notable limitations (see Table 1):

**Limitation ①:** *No proper adjustment for confounding and thus bias:* Some methods for CATE over time fails to properly adjust for time-varying confounding properly, which leads to estimates that are *biased*. Here, prominent examples are counterfactual recurrent network (**CRN**) [3] and the causal transformer (**CT**) [20]. These methods attempt to alleviate time-varying confounding via balanced representations. However, balancing was originally designed for reducing finite-sample estimation variance and *not* for mitigating confounding bias [34]. Hence, such methods act as heuristics without a theoretical justification. The difficulty of enforcing balanced representations may even introduce further confounding bias [21]. Unlike these methods, our DeepBlip allows for proper adjustments and is thus unbiased.

**Limitation ②:** *Unstable for long time horizons:* Other neural methods build on frameworks from statistics such as marginal structural models (MSMs) [26, 22] and G-computation [27, 30]. Examples are: (i) **G-Net** [32] and G-transformer (**GT**) [14], which are both based on G-computation and thus compute nested conditional expectations over time. Hence, these methods require the entire counterfactual outcome trajectories and thus model the full data-generating process of covariates and outcomes, which becomes exponentially more complex as time horizons grow. (ii) MSMs-based methods like **R-MSNs** [19] and the **DR-learner** (for time-varying settings) [11] use inverse propensity weighting (IPW) to re-weight outcomes as in a randomized control trial. In time-varying settings, the propensity is a multiplication of a sequence of probabilities, which has known to become instable when the horizon is large.<sup>4</sup> In sum, while these methods are unbiased, they require modeling *entire* counterfactual outcome trajectories, which makes learning *unstable* over long time horizons. As a remedy to this, our DeepBlip breaks the CATE into localized effects at each time step via blip functions.

**Structural nested mean models (SNMMs):** SNMMs *refs* offer a principled framework for estimating CATE over time by directly estimating incremental treatment effects via so-called blip functions. In principle, SNMMs could address both of the above limitations; however, SNMMs are only an abstract theoretical foundation, *not* an off-the-shelf model that can be directly applied. Early implementations [28, 29] relied on strong parametric assumptions (e.g., linearity) and were limited to short time horizons. Recently, Lewis et al. [18] employed SNMMs with linear models, yet which thus fails to condition on the rich, complex information in patient histories and which makes it unsuitable for personalized medicine in realistic, high-dimensional settings for medicine. So far, to the best of our knowledge, a neural instantiation of SNMMs are missing.

**Research gap:** To the best of our knowledge, there is no neural implementation of SNMMs. We thus extend upon the work of Lewis et al. [18] and develop the first neural framework for SNMMs. For this, we introduce a new, flexible architecture and a tailored two-stage learning algorithm that allows us to leverage the blip function from SNMMs for CATE estimation over time. As a result, ours is the first neural method to address both limitations ① and ② from above.

Method	Methodological limitations		Benefits for medical practice	
	① Unbiased	② Stable wrt long horizons	Orthogonal	Offline efficiency
CRN [3]	✗	✓	✗	✓
CT [20]	✗	✓	✗	✓
G-Net [32]	✓	✗	✗	✗
GT [14]	✓	✗	✗	✗
R-MSNs [19]	✓	✗	✗	✗
DR-learner [11]	✓	✗	✓	✗
DeepBlip (ours)	✓	✓	✓	✓

Table 1: Neural methods for learning CATE over time.

**Research gap:** To the best of our knowledge, there is no neural implementation of SNMMs. We thus extend upon the work of Lewis et al. [18] and develop the first neural framework for SNMMs. For this, we introduce a new, flexible architecture and a tailored two-stage learning algorithm that allows us to leverage the blip function from SNMMs for CATE estimation over time. As a result, ours is the first neural method to address both limitations ① and ② from above.

<sup>3</sup>There have been some attempts to use non-parametric models [33, 36, 42], yet these approaches impose strong assumptions on the outcomes and their scalability is limited. For these reasons, we focus on neural methods, which offer better flexibility and scalability for complex, high-dimensional medical data.

<sup>4</sup>This is due to overlap violations, which is especially challenging in time-varying settings, since the number of possible treatment combinations grows exponentially with the horizon [11]. Hence, IPW often leads to extreme weights, and thus instabilities due to division by values close to zero.

### 3 Problem Formulation

**Setup:** We follow the standard setup [3, 20, 14, 11] for estimating CATE over time  $t \in \{1, 2, \dots, T\} \subset \mathbb{N}$  given by: (1) the target outcome  $Y_t \in \mathbb{R}$ ; (2) time-varying covariates  $X_t \in \mathcal{X} \subset \mathbb{R}^{d_x}$ ; and (3) treatment  $A_t \in \mathcal{A} \subset \mathbb{R}^{d_a}$  that can be either discrete or continuous. We assume w.l.o.g. that the static features (e.g., age, sex) are included in the covariate.

**Notation:** To simplify notation, we use overlines to denote the full sequence of a variable (e.g.,  $\bar{X}_t = (X_1, \dots, X_t)$ ). We refer to a sequence of variables that starts at  $t$  and ends at  $t + \tau$  via  $A_{t:t+\tau} = (A_t, A_{t+1}, \dots, A_{t+\tau})$ . We use the lowercase letter to denote a realization of a random variable (e.g.,  $X_t = x_t$ ). We use an asterisk  $*$  to indicate a constant quantity (e.g., a fixed treatment  $a_t^*$ ). We denote the patient history by  $H_t = (\bar{X}_t, \bar{A}_{t-1}, \bar{Y}_{t-1})$ .

**CATE estimation over time:** We build upon the potential outcome framework [31] for the time-varying setting [30]. We aim to estimate the CATE over time between two treatment sequences for a given patient history, i.e.,

$$\mathbb{E} \left[ Y_{t+\tau}^{(a_{t:t+\tau}^*)} - Y_{t+\tau}^{(b_{t:t+\tau}^*)} \mid H_t = h_t \right], \quad 0 \leq t \leq T - \tau, \quad (1)$$

where  $Y_{t+\tau}^{(a_{t:t+\tau}^*)}$  and  $Y_{t+\tau}^{(b_{t:t+\tau}^*)}$  represents the  $\tau$ -step-ahead potential outcomes under interventions  $\text{do}(A_{t:t+\tau} = a_{t:t+\tau}^*)$  and  $\text{do}(A_{t:t+\tau} = b_{t:t+\tau}^*)$ , respectively (see Appendix A for a formal definition of potential outcomes and interventions).

**Identifiability:** We make the following identifiability assumptions [26, 25] that are standard in the time-varying setting [3, 20, 19, 32]: (1) *Consistency*: The potential outcome under the intervention by the observed treatment equals the observed outcome, namely,  $Y_t^{(A_t)} = Y_t$ . (2) *Overlap*: Given an observed history  $H_t = h_t$ , if  $\Pr(H_t = h_t) > 0$ , then any possible treatment has a positive probability of being received:  $\forall a \in \mathcal{A}_t, \Pr(A_t = a \mid H_t = h_t) > 0$ . (3) *Sequential ignorability*: The potential outcome under an arbitrary intervention is independent of the treatment assignment conditioned on the history, i.e.,  $Y_t^{(a_t^*)} \perp A_t \mid H_t = h_t$ .

However, estimating the CATE over time is non-trivial due to *time-varying confounding* [6, 14]. In the time-varying setting, covariates act as confounders because they are influenced by earlier treatments and affect later treatments. However, these time-varying confounders are *unobserved*, because of which naïve adjustments as in the static setting are impossible (see Appendix A.2). Here, we adjust for time-varying confounding through the use of SNMMs.

**Blip function:** SNMMs model the incremental effect of treatments (which are called “blips”) at time  $t + k$  on the mean outcome at  $t + \tau$ , given observed patient history  $H_t$  [37]. These “blips” accumulate over time into the total treatment effect and thus allow to rigorously adjust for time-varying confounding [29] (see Appendix B.1 for details). Formally, the blips are defined via a **blip function** [28, 29]:

$$\gamma_{t,k}(\bar{x}_{t+1:t+k}, \bar{a}_{t:t+k}; h_t) = \mathbb{E} \left[ Y_{t+\tau}^{(a_{t:t+k}, d_{t+k+1:t+\tau})} - Y_{t+\tau}^{(a_{t:t+k-1}, 0, d_{t+k+1:t+\tau})} \mid A_{t:t+k} = a_{t:t+k}, X_{t+1:t+k} = x_{t+1:t+k}, H_t = h_t \right]. \quad (2)$$

Intuitively, the blip function  $\gamma_{t,k}$  isolates the causal effect of each treatment decision *locally*. This breaks the sequential dependencies over long temporal dependencies and thus is more stable over long horizons ( $\rightarrow$  thus addressing limitation ②).

Nevertheless, SNMMs offer *only* a theoretical framework for identifying CATEs – they are *not* ready-to-use algorithms or models. Hence, implementing SNMMs with neural networks in particular is non-trivial: this requires a tailored learning objective that allows for neural parameterization and that supports efficient, end-to-end training and inference, which is the contribution of our DeepBlip.

### 4 Our DeepBlip Framework

In this section, we present DeepBlip. First, we introduce how we learn the CATE via blip functions using a neural parameterization (Sec. 4.1), then introduce our  $L^1$ -moment loss (Sec. 4.2), our model architecture (Sec. 4.3), and the training and inference procedure (Sec. 4.4).

#### 4.1 Learning the CATE via blip functions

**Overview:** Our DeepBlip leverages Eq. (3) to adjust for time-varying confounding ( $\rightarrow$  thus addressing limitation ②). Our task thus reduces to estimating the blip functions – in particular, so-called *blip coefficients* that parametrize the blip functions. However, we do **not** attempt to estimate the coefficients directly. Instead, we optimize a  $L^1$ -moment loss that directly predicts the blip coefficients and which allows us to estimate Eq. (3) more efficiently.

**Parameterization trick:** We first explain how we estimate the CATE via the blip function. For this, we adopt a similar parametrization for the blip function as in [18], namely,  $\gamma_{t,k}(\bar{x}_{t+1:t+k}, \bar{a}_{t:t+k}; h_t) = \psi_{t,k}(h_t)' a_{t+k}$ , but where  $\psi_{t,k}$  is a **neural network**. Under identifiability assumptions and the parametrization for  $\gamma_{t,k}$  defined above, the CATE of  $a^*$  against  $b^*$  for any two treatment sequences  $a^*, b^* \in \mathbb{R}^{(\tau+1) \cdot d_a}$  is (see [18] for a formal derivation):

$$\mathbb{E}[Y_{t+\tau}^{(a^*)} - Y_{t+\tau}^{(b^*)} | H_t = h_t] = \sum_{k=0}^{\tau} \psi_{t,k}(h_t)' (a_{t+k}^* - b_{t+k}^*). \quad (3)$$

We refer to  $\psi_{t,k}(h_t)$  as the conditional *blip coefficients* of the blip function  $\gamma_{t,k}$ .

*Why do we need a tailored architecture and learning algorithm?* A key component of our framework is that the function  $\psi_{t,k}(h_t)$  is parameterized by a sequential neural network (e.g., LSTM or transformer). This is a crucial difference from traditional SNMMs, which were developed for estimating the ATE over a fixed number of time steps (i.e.  $\mathcal{H}_t = \emptyset \wedge t \equiv 0 \wedge T \equiv \tau$ ) and where, as a result, blip coefficients are constants. These constants are typically estimated through *iteratively* solving a set of moment equations via g-estimation [28, 29, 37] (see Appendix B.1). However, such an approach is not compatible with neural network-based learning. In contrast, DeepBlip introduces a tailored neural architecture (Sec. 4.3) that we can train via gradient-based optimization (Sec. 4.2).

#### 4.2 $L^1$ -moment loss

We reformulate the moment-based linear equations from [18] as an equivalent iterative minimization problem for  $k = \tau, \dots, k = 0$ . At each time step  $k$ , we aim to find the minimizer  $\psi_{t,k}^*(\cdot)$ , which is a function that maps the history  $h_t$  to the blip coefficients, via

$$\psi_{t,k}^* = \arg \min_{\hat{\psi}_{t,k}(\cdot) \in \Phi_{t,k}} \mathbb{E} \left[ \left\| \mathbb{E} \left[ (\tilde{Y}_{t,k} - \sum_{j=k+1}^{\tau} \psi_{t,j}^*(h_t)' \tilde{A}_{t,j,k} - \hat{\psi}_{t,k}(h_t)' \tilde{A}_{t,k,k}) \tilde{A}_{t,k,k} | H_t \right] \right\|_1 \right], \quad (4)$$

where  $\Phi_{t,k}$  is the function space for the blip coefficient predictors and  $\|\cdot\|_1$  is the  $L^1$ -norm operator. The expectation outside is taken over all random variables  $H_t$ . We name the target as the  **$L^1$ -moment loss**.

Here, we employ an  $L^1$  loss for empirical reasons (see our ablation studies in Appendix C.2). The reason is that the moment has a high variance, especially with growing time horizon  $\tau$  and due to the mini-batch sampling, which introduces another source of variance later. As a result, the  $L^1$  loss is beneficial since it is more robust to such variance.

**Theoretical properties:** Below, we first show that the loss recovers the ground-truth blip coefficients. Then, we show that our loss is *Neyman-orthogonality* (see [5] for formal definition), which ensures double robustness. This means that the target loss is *robust* against perturbations of the nuisance functions [5, 16].

**Remark 1.** If  $\forall 0 \leq k \leq \tau, \psi_{t,k} \in \Phi_{t,k}$ , then the solution of the risk minimization scheme in Eq. (4) given by  $(\psi_{t,0}^*, \dots, \psi_{t,\tau}^*)$ , yields the ground-truth blip coefficients. That is,  $\psi_{t,k}^* = \psi_{t,k}$ .

**Remark 2.** The moment loss is *Neyman-orthogonal*.

The above remarks follow from the theory in [18], which is easy to extend to our setting (see Appendix B.2 and Appendix B.3, respectively).

**Double optimization trick for our  $L^1$ -moment loss:** In order to find  $\psi_{t,k}^*$ , all the previous blip predictors  $\psi_{t,j}^*, j \geq k$  are required. However, the ground-truth predictors are generally not available at the beginning. To avoid solving  $\psi_{t,k}^*$  sequentially, we propose a *double optimization trick* that allows *simultaneous* training of all the blip predictors: During each iteration, first, the blip predictor



224  $\hat{\psi}_t$  makes two forward passes to generate two sets of the blip coefficients  $\hat{\psi}_t^1(h_t)$  and  $\hat{\psi}_t^2(h_t)$ . Then,  
 225  $\hat{\psi}_{t,j}^2(h_t)$  is treated as the pseudo blip effects that replaces  $\psi_{t,j}^*(h_t)$  in Eq. (4). For  $k = 0, \dots, \tau$ , the  
 226 adapted  $L^1$ -moment loss at step  $k$  is then given empirically by

$$\mathcal{L}_{\text{blip}}^k = \frac{1}{T - \tau} \sum_{t=1}^{T-\tau} \left\| \sum_{i=1}^n (\tilde{Y}_{t+k}^i - \sum_{j=k+1}^{\tau} \hat{\psi}_j^2(H_t^i)' \tilde{A}_{t,j,k}^i - \hat{\psi}_k^1(H_t^i)' \tilde{A}_{t,k,k}^i \cdot \tilde{A}_{t,k,k}^i) \right\|_1. \quad (5)$$

### 227 4.3 Model architecture

228 DeepBlip works in two stages (see  
 229 Fig. 2): **•Stage ① (nuisance network)**: models the nuisance functions  
 230 to estimate the residuals in Eq. (6).  
 231 **•Stage ② (blip prediction network)**:  
 232 estimates the blip coefficients given the observed history  $h_t$ . The neural  
 233 networks in both stages have a similar  
 234 structure: (i) a *sequential encoder* that  
 235 encodes the observed history  $H_t$ , and  
 236 (ii) multiple *prediction heads* that take  
 237 the encoded history as input to predict  
 238 the targets.  
 239  
 240

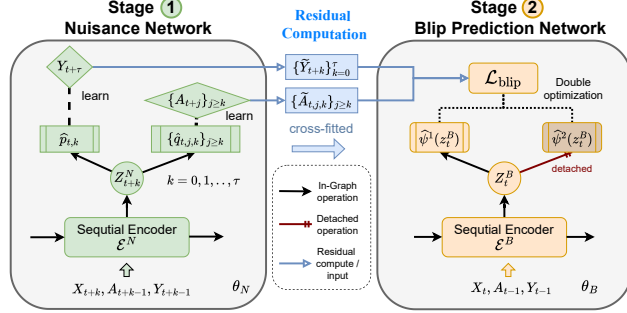


Figure 2: Neural architecture of the two-stage DeepBlip framework.

#### 241 Why we need a two-stage design:

242 To construct the  $L^1$ -moment loss defined Eq. (4), we need the variables  $\tilde{Y}_{t,k}$ ,  $\tilde{A}_{t,j,k}$ , defined as  
 243 (see Appendix B.1 for details):

$$\tilde{Y}_{t,k} = Y_{t+\tau} - \mathbb{E}[Y_{t+\tau} | H_{t+k} = h_{t+k}], \quad \tilde{A}_{t,j,k} = A_{t+j} - \mathbb{E}[A_{t+j} | H_{t+k} = h_{t+k}]. \quad (6)$$

244 As we see here, we must compute the *residuals* between the outcome variable and their regressed  
 245 means *before* we optimize  $\mathcal{L}_{\text{blip}}^k$ . We thus follow previous literature [18, 11, 16] and treat the  
 246 conditional expectations  $\mathbb{E}[Y_{t+\tau} | H_{t+k} = h_{t+k}]$  and  $\mathbb{E}[A_{t+j} | H_{t+k} = h_{t+k}]$  as *nuisance*  
 247 *functions*, so that we first estimate the nuisance function (Stage ①) and then train the blip prediction  
 248 network to minimize the  $L^1$ -moment loss (Stage ②).

249 **Neural backbone:** Our DeepBlip is flexible and allows for different neural backbones (e.g., LSTM  
 250 or transformer). These are necessary to capture the patient history: We notice that the networks at both  
 251 Stage ① and ② take the history variable  $H_t = (\bar{X}_t, \bar{A}_{t-1}, \bar{Y}_{t-1}) \in \mathcal{H}_t$  as input. Hence, both stages  
 252 can be written as a function  $f: \cup_{t=1}^T \mathcal{H}_t \rightarrow \mathbb{R}^c$ , where  $c$  is the number of outputs. However,  $\dim(H_t)$   
 253 varies over time, which makes  $H_t$  not suitable as a direct input to a neural network. A standard  
 254 way to handle this is by using a sequential model to iteratively take the inputs  $(X_t, A_{t-1}, Y_{t-1})$  and  
 255 then maintain a vector  $Z_t \in \mathbb{R}^{d_z}$  with fixed dimension that encodes all the necessary information  
 256 [19, 32, 20, 14]. Here, we thus use LSTMs and transformers (see details of the architectures in  
 257 Appendix G). Finally, we stress that each stage uses a *separate* encoder:  $\mathcal{E}_{\theta_N}^N$  for Stage ①, and  $\mathcal{E}_{\theta_B}^B$   
 258 for Stage ② ( $N$  for Nuisance and  $B$  for Blip) with different model weights  $\theta_N$  and  $\theta_B$ .

259 **Stage ①: nuisance network.** The nuisance network  $(\mathcal{E}_{\theta_N}^N, \{\text{gp}_{\theta_N}^k\}_{k=0}^{\tau}, \{\text{gq}_{\theta_N}^{j,k}\}_{0 \leq k \leq j \leq \tau})$  consists  
 260 of a sequential encoder  $\mathcal{E}_{\theta_N}^N$  and a collection of prediction heads  $\{\text{gp}_{\theta_N}^k\}_{k=0}^{\tau}, \{\text{gq}_{\theta_N}^{j,k}\}_{0 \leq k \leq j \leq \tau}$ . The  
 261 nuisance networks are responsible for computing the following nuisance functions:

$$p_{t,k}(h_{t+k}) := \mathbb{E}[Y_{t+\tau} | H_{t+k} = h_{t+k}], \quad 1 \leq t \leq T - \tau, 0 \leq k \leq \tau \quad (7)$$

$$q_{t,j,k}(h_{t+k}) := \mathbb{E}[Q_{t,j} | H_{t+k} = h_{t+k}], \quad 1 \leq t \leq T - \tau, 0 \leq k \leq j \leq \tau \quad (8)$$

262 For a patient with history  $H_t$  and subsequent covariates  $X_{t+1:t+k}, A_{t:t+k-1}$ , we proceed as follows:  
 263 First, the encoder  $\mathcal{E}_{\theta_N}^N$  learns the representation at time  $t+k$  (note that  $H_{t+k} = H_t \cup X_{t+1:t+k} \cup$   
 264  $A_{t:t+k-1}$ ), which is given by  $Z_{t+k}^N = \mathcal{E}_{\theta_N}^N(H_{t+k})$ . Second, the prediction heads receive  $Z_{t+k}^N$  to  
 265 compute the regressed outcomes for the nuisance functions via:

$$\text{gp}_{\theta_N}^k(Z_{t+k}^N) = \hat{p}_{t,k}(H_{t+k}), \quad \text{gq}_{\theta_N}^{j,k}(Z_{t+k}^N) = \hat{q}_{t,j,k}(H_{t+k}) \quad \text{for } k = 0, \dots, \tau \text{ and } k \leq j \leq \tau \quad (9)$$

where  $Z_{t+k}^N = \mathcal{E}_{\theta_N}^N(H_{t+k})$ . Third, the residuals are computed via

$$\tilde{Y}_{t,k} \approx Y_{t+\tau} - \text{gp}_{\theta_N}^k(Z_{t+k}^N), \quad \tilde{A}_{t,j,k} \approx A_{t+j} - \text{gq}_{\theta_N}^{j,k}(Z_{t+k}^N) \quad (10)$$

**Stage ②: blip prediction network.** The blip prediction network ( $\mathcal{E}_{\theta_B}^B, \{\text{gb}_{\theta_B}^k\}_{k=0}^\tau$ ) is responsible for predicting the blip coefficients  $\psi_t(h_t) = (\psi_{t,0}, \dots, \psi_{t,\tau}) \in \mathbb{R}^{r(\tau+1)}$  as described in Eq. (4). Here, we proceed as follows. First, the sequential encoder  $\mathcal{E}_{\theta}^B$  (B for **Blip**) processes the patient's history  $H_t$  into a representation  $Z_t^B = \mathcal{E}_{\theta}^B(h_t)$ . Then, for each horizon  $k \in \{0, 1, \dots, \tau\}$ , the prediction head  $\text{gb}_{\theta_B}^k$  maps  $Z_t^B$  onto the corresponding blip coefficient:

$$\hat{\psi}_{t,k}(H_t) = \text{gb}_{\theta_B}^k(Z_t^B) \sim \psi_{t,k}(H_t) \in \mathbb{R}^r \quad \text{where } Z_t^B = \mathcal{E}_{\theta_B}^B(H_t). \quad (11)$$

#### 4.4 Training and Inference

Taken together, the training procedure of DeepBlip now follows two steps (see Fig. 2): (1) train the nuisance networks and compute the residuals, and (2) train the blip prediction network. In contrast, inference with DeepBlip is highly efficient as it involves *only* the second-stage blip prediction network. Details are below. We provide the pseudocode in Alg. 1 and Alg. 2 in the appendix.

**Step ①: Train nuisance network.** The nuisance network is trained to predict nuisance functions  $p_{t,k}(h_{t+k})$  and  $q_{t,j,k}(h_{t+k})$  simultaneously. Since  $p_{t,k}(h_{t+k})$  is the conditional expectation of real outcome  $Y_{t+\tau} \in \mathbb{R}$ , we use the squared error loss  $\mathcal{L}_p = \frac{1}{(T-\tau)(\tau+1)} \sum_{t=1}^{T-\tau} \sum_{k=0}^\tau (\text{gp}_{\theta_N}^k(Z_{t+k}^N) - Y_{t+\tau})^2$ . For  $q_{t,j,k}(h_{t+k})$ , which denotes the treatment response, we proceed for the  $i$ -th treatment in  $A_{t+j} \in \mathbb{R}^{d_a}$  as follows. If  $(A_{t+j})_i$  is a continuous variable, then we apply the squared loss:  $\mathcal{L}_{q,i} = \frac{2}{(T-\tau)(\tau+1)(\tau+2)} \sum_{t=1}^{T-\tau} \sum_{0 \leq k \leq j \leq \tau} (\text{gq}_{\theta_N}^{k,j}(Z_{t+k}^N)_i - (A_{t+j})_i)^2$ . If  $(A_{t+j})_i$  is a binary variable, then we apply the binary cross entropy loss  $\mathcal{L}_{q,i} = \frac{2}{(T-\tau)(\tau+1)(\tau+2)} \sum_{t=1}^{T-\tau} \sum_{0 \leq k \leq j \leq \tau} \text{BCE}((A_{t+j})_i, \text{gq}_{\theta_N}^{k,j}(Z_{t+k}^N)_i)$ . For categorical variables with more than 2 classes, we preprocess the variable into a one-hot vector of binary variables. Since the network predicts these targets simultaneously, we update the parameter  $\theta_N$  by backpropagating the sum of all the losses discussed above, i.e.,  $\mathcal{L}_N = \mathcal{L}_p + \frac{1}{d_a} \sum_{i=1}^{d_a} \mathcal{L}_{q,i}$

**Step ②: Train blip prediction network.** After having trained the nuisance network, we freeze its parameters and then compute the residuals of each sample as in Eq. (6) for the  $L^1$ -moment loss. To accelerate the training process, we adopt the double optimization trick from above: For  $k = 0, \dots, \tau$ , we perform two forward passes that create two predictions  $\hat{\psi}_{t,k}^1(H_t)$  and  $\hat{\psi}_{t,k}^2(H_t)$ . The latter is then *detached* from the computation graph before feeding into the adapted  $L^1$ -moment loss at step  $k$ . The final loss target is then given by:  $\mathcal{L}_{\text{blip}} = \sum_{k=0}^\tau \mathcal{L}_{\text{blip}}^k$ . We note that, for  $k = \tau$ , there is **no** detached blip coefficients term in Eq. (5). Hence,  $\hat{\psi}_{t,\tau}^1(H_t)$  is directly supervised by the true  $L^1$ -moment loss and can directly learn the ground-truth. As such,  $\mathcal{L}_{\text{blip}}^{\tau-1}$  gradually approximates the true  $L^1$ -moment loss, which then supervises  $\hat{\psi}_{t,\tau-1}^1(H_t)$ , and so on. As a result, all prediction heads will gradually learn to predict the blip coefficients from  $t = \tau$  to  $t = 0$ .

*Remark 3. Under standard assumptions, the output of our DeepBlip has a mean squared error guarantee:*

$$\max_{t \leq T-\tau} \max_{k \in \{0, \dots, \tau\}} \mathbb{E} \left[ \left\| \hat{\psi}_{t,k} - \psi_{t,k} \right\|_{2,2}^2 \right] = O(r^2 \delta_n^2), \quad \delta_n^2 \propto \frac{\log \log(n)}{n} \quad (12)$$

The above is adopted from SNMM methods [18] that were originally developed for linear models, yet we offer a neural instantiation. Details are in Appendix B.4).

**Inference at runtime:** Once trained, our DeepBlip predicts the CATE over time (i.e.,  $\mathbb{E}[Y_{t+\tau}^{(a^*)} - Y_{t+\tau}^{(b^*)} \mid H_t = h_t]$ ) through only the blip prediction network:

$$\sum_{k=0}^\tau \text{gb}_{\theta_B}^k(z_t^B)'(a_k^* - b_k^*), \quad \text{where } z_t^B = \mathcal{E}_{\theta_B}^B(h_t) \quad (13)$$

**Efficient offline evaluation:** Once we have estimated the blip coefficients, then we can instantly identify the treatment sequence  $a^*$  with the best effect compared to the baseline  $b^*$  (e.g., a treatment sequence with no interventions). The reason is that the blip coefficients do *not* depend on treatments. Hence, our DeepBlip is much more efficient for evaluating the personalized effects of different treatment sequences compared to existing methods that require re-computation [19, 20, 3] or even re-training [32, 14]. This is highly relevant in personalized medicine where clinicians and patients jointly reason about different treatment strategies [9].

**Implementation details.** We instantiate our DeepBlip with a transformer architecture (see Appendix G). We also provide a variant based on an LSTM, which, despite the simpler architecture, is still highly competitive (see Appendix C.1).

## 5 Experiments

**Baselines:** We demonstrate the performance of our DeepBlip against key baselines from the literature (see Table 1) for the task of estimating CATE (or conditional average potential outcomes) on medical datasets. Descriptions of the baseline methods are available in Appendix E. We further select the HA-PI-learner from [11] instantiated by transformer (named **HA-TRM**) as a naïve baseline. We provide additional implementation details – including architecture choices, training procedures, and hyperparameter tuning – in Appendix G. To ensure a fair comparison, all methods – including baselines – use the **same** neural backbone architecture, so any performance differences must be attributed solely to that our learning objective is better (i.e., unbiased and stable over longer time horizons). All results are averaged over 5 runs.

**Ablations:** We include ablation studies in Appendix C.1, where we validate our component-wise blip coefficient estimates in Appendix C.2. We also provide an instantiation of our DeepBlip with an LSTM instead of a transformer. Importantly, even our ablation is highly competitive and outperforms the majority of transformer-based baselines (see Appendix C.1).

### 5.1 Tumor growth dataset

**Setting:** We use the pharmacokinetic-pharmacodynamic tumor growth dataset [12], which is commonly used for benchmarking CATE methods over time [19, 32, 3, 14, 20]. The dataset describes the time-varying effects of chemotherapies and radiotherapies, for which treatment assignments depend on previous outcomes, subject to time-varying confounding. The amount of confounding is controlled by the simulation parameter  $\gamma_{\text{conf}}$ . Details are in Appendix D.1.

**Results:** Figure 3 shows the average RMSE of CATE against increasing confounding  $\gamma_{\text{conf}}$  and under  $\tau = 2$ . Our **DeepBlip** outperforms all baselines for  $\gamma_{\text{conf}} \geq 2$ . This matches the purpose of our method to deal with time-varying confounding. More importantly, our DeepBlip achieves large performance gains under strong confounding levels ( $\gamma_{\text{conf}} > 6$ ). This highlights that our DeepBlip is robust against time-varying confounding by providing adequate adjustment for time-varying confounding.

We could further make the following observations: ① The MSM-based **R-MSN** performs poorly across all confounding levels and even has a higher RMSE than HA-LSTM for  $\gamma_{\text{conf}} \leq 6$ . This aligns with the inverse propensity weighting in MSMs is highly unstable, which was the motivation for our method. ② Baselines like **CT** and **CRN** that use balanced representation (in orange) are ineffective. This is expected as balanced representations were originally developed for reducing finite-sample estimation variance and *not* for proper adjustment (see the original work [34] on balanced representations for a discus-

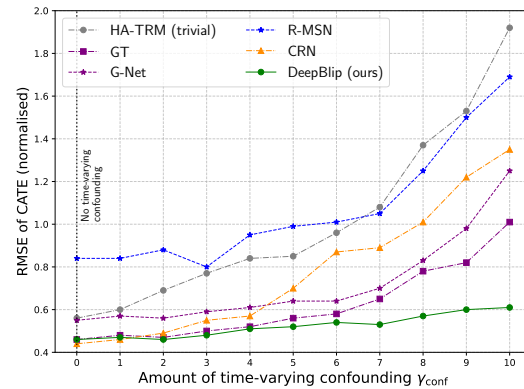


Figure 3: **Results for tumor growth dataset.** Normalized RMSE (averaged over 5 runs) of CATE predictions against ground-truth over growing confounding. Here:  $\tau = 2$



sion), because of which these baselines are known to be biased. ③ G-computation-based methods like **G-Net** and **GT** show slightly lower RMSE for  $\gamma_{\text{conf}} \leq 1$  but still perform significantly worse than our DeepBlip for  $\gamma_{\text{conf}} \geq 6$ . We attribute this to the fact that the learning is unstable, which we empirically verify in the following by varying the prediction horizon  $\tau$ .

## 5.2 MIMIC-III dataset

**Setting:** Next, we evaluate the performance for longer prediction horizons  $\tau$ . We build upon MIMIC-III [15], a widely used benchmark for evaluating CATE over time [3, 20, 14]. Following previous literature [33, 20, 14], we extract patient vitals from MIMIC-III and then simulate the patient outcome over time with the mixed dynamics of exogenous dependency, endogenous dependency, and treatment effects combined. Treatments are assigned based on previous outcomes and patient vitals, which again, introduces time-varying confounding (see Appendix D.2).

**Results:** Table 2 shows the average RMSE (with std. dev.) over five different runs with  $\gamma_{\text{conf}} = 1$  and varying prediction horizon  $\tau$ . First, our **DeepBlip** consistently achieves lower RMSE compared to other baselines for  $\tau \geq 2$ . Of note, the performance gain from DeepBlip becomes larger as  $\gamma_{\text{conf}}$  increases. For  $\tau = 10$ , DeepBlip achieves  $\sim 38\%$  performance gain compared to the second-best model (here: GT). This highlights that our DeepBlip is temporally stable over longer horizons.

We further make the observations that all baselines either struggle with high-dimensional covariates or become unstable as  $\tau$  increases. ① The MSM-based method (**R-MSNs**) exhibits the highest variance across all  $\tau$ , indicating that it struggles with high-dimensional propensity modeling and becomes unstable over time with increasing standard deviation. The reason is that inverse propensity weighting produces unstable weights. ② Methods like **CRN** and **CT** perform better than baselines with high variance like R-MSNs due to the way they handle high-dimensional covariates. However, both **CRN** and **CT** are known to be biased and thus inferior to GT and our DeepBlip. ④ G-computation-based methods (i.e., **G-Net** and **GT**) achieve a lower RMSE than the other baselines due to proper adjustments, but still are not as stable as our method. This is because G-computation accumulates error over time due to modeling nested expectations.

	$\tau = 2$	$\tau = 3$	$\tau = 4$	$\tau = 5$	$\tau = 6$	$\tau = 7$	$\tau = 8$	$\tau = 9$	$\tau = 10$
HA-LSTM (naïve) [11]	$0.68 \pm 0.02$	$0.89 \pm 0.03$	$0.97 \pm 0.04$	$1.02 \pm 0.10$	$1.42 \pm 0.20$	$1.92 \pm 0.40$	$2.57 \pm 0.44$	$2.58 \pm 0.56$	$3.11 \pm 0.72$
R-MSNs [19]	$0.73 \pm 0.14$	$0.98 \pm 0.17$	$1.12 \pm 0.21$	$1.25 \pm 0.28$	$1.65 \pm 0.57$	$2.25 \pm 1.02$	$2.85 \pm 1.18$	$3.20 \pm 1.42$	$3.55 \pm 1.50$
CRN [3]	$0.49 \pm 0.05$	$0.66 \pm 0.11$	$0.82 \pm 0.12$	$1.05 \pm 0.22$	$1.22 \pm 0.35$	$1.43 \pm 0.33$	$1.62 \pm 0.42$	$1.83 \pm 0.43$	$2.04 \pm 0.54$
CT [20]	$0.52 \pm 0.07$	$0.64 \pm 0.12$	$0.79 \pm 0.11$	$1.01 \pm 0.18$	$1.18 \pm 0.33$	$1.77 \pm 0.52$	$1.85 \pm 0.49$	$1.99 \pm 0.63$	$1.98 \pm 0.60$
G-Net [32]	$0.42 \pm 0.05$	$0.58 \pm 0.08$	$0.73 \pm 0.12$	$1.05 \pm 0.25$	$1.38 \pm 0.40$	$1.75 \pm 0.60$	$2.15 \pm 0.80$	$2.55 \pm 0.90$	$3.12 \pm 1.05$
GT [14]	$0.40 \pm 0.01$	$0.52 \pm 0.02$	$0.63 \pm 0.08$	$0.75 \pm 0.17$	$0.85 \pm 0.13$	$0.95 \pm 0.26$	$1.10 \pm 0.34$	$1.25 \pm 0.37$	$1.50 \pm 0.45$
DeepBlip (ours)	<b><math>0.39 \pm 0.11</math></b>	<b><math>0.48 \pm 0.12</math></b>	<b><math>0.56 \pm 0.16</math></b>	<b><math>0.64 \pm 0.19</math></b>	<b><math>0.70 \pm 0.21</math></b>	<b><math>0.79 \pm 0.24</math></b>	<b><math>0.82 \pm 0.27</math></b>	<b><math>0.88 \pm 0.28</math></b>	<b><math>0.93 \pm 0.32</math></b>
Improvement	2.5%	7.6%	11.1%	14.7%	17.6%	16.8%	25.5%	29.6%	38.0%

Table 2: **MIMIC-III with longer time horizons  $\tau$ .** Normalized RMSE (mean  $\pm$  std. dev. over 5 runs) for  $\tau$ -step-ahead CATE estimation on the MIMIC-III dataset. We highlight the relative improvement over the best-performing baseline.  $\Rightarrow$  Our DeepBlip consistently outperforms the baselines for  $\tau \geq 2$ .

## 6 Discussion

**Limitations:** (1) Our work is subject to the standard assumptions for treatment effect estimation, which are standard in the literature [3, 11, 14, 19, 20, 28, 29, 32]. (2) Our work is further subject to the characteristics of how the blip function is parameterized. (3) The overall training cost is comparable to that of the baseline. Importantly, the runtime ( $\sim 30$  min, see Appendix G.2) is similar across all baselines, and in practice, we often observe faster convergence due to the more stable learning of our approach. Importantly, the computational cost is typically not a major concern in medical applications, as the models are trained only once, and all baselines scale efficiently to all real-world medical datasets from practice.

**Broader impact:** We expect our contribution to have a significant impact on *reliable* decision-making in personalized medicine. DeepBlip provides a *stable* learning framework for *efficient* offline evaluation of personalized treatment strategies over long time horizons.

**Conclusion:** We are the first to build a neural framework using blip functions to estimate CATEs over time.

## References

- [1] Ahmed M. Alaa and Mihaela van der Schaar. Bayesian inference of individualized treatment effects using multi-task Gaussian processes.
- [2] Ahmed Allam, Stefan Feuerriegel, Michael Rebhan, and Michael Krauthammer. Analyzing patient trajectories with artificial intelligence. *Journal of Medical Internet Research*, page e29812, 2021.
- [3] Ioana Bica, Ahmed M. Alaa, James Jordon, and Mihaela van der Schaar. Estimating Counterfactual Treatment Outcomes over Time Through Adversarially Balanced Representations. In *ICLR*, 2020.
- [4] Ioana Bica, Ahmed M. Alaa, Christoph Lambert, and Mihaela van der Schaar. From real-world patient data to individualized treatment effects using machine learning: Current and future methods to address underlying challenges. *Clinical Pharmacology & Therapeutics*, 2021.
- [5] Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/Debiased Machine Learning for Treatment and Causal Parameters. *The Econometrics Journal*, 2018.
- [6] Amanda Coston, Edward Kennedy, and Alexandra Chouldechova. Counterfactual predictions under runtime confounding. In *Advances in Neural Information Processing Systems*, 2020.
- [7] Alicia Curth and Mihaela van der Schaar. Nonparametric estimation of heterogeneous treatment effects: From theory to learning algorithms. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, Proceedings of Machine Learning Research, 13–15 Apr 2021.
- [8] William Falcon, Jirka Borovec, Adrian Wälchli, Katharina Eggersperger, Roman Schiele, Massimo Patacchiola, Ryutaro Tanno, Alexander J. Landgraf, Chad Baskin, Jeff Jordan, Raghuvardhan Sinha, et al. PyTorch Lightning. In *NeurIPS 2019 Reproducibility Challenge*, 2019.
- [9] Stefan Feuerriegel, Dennis Frauen, Valentyn Melnychuk, Jonas Schweisthal, Konstantin Hess, Alicia Curth, Stefan Bauer, Niki Kilbertus, Isaac S. Kohane, and Mihaela van der Schaar. Causal machine learning for predicting treatment outcomes. *Nature Medicine*, 04 2024.
- [10] Dennis Frauen, Tobias Hatt, Valentyn Melnychuk, and Stefan Feuerriegel. Estimating average causal effects from patient trajectories. 37:7586–7594, Jun. 2023.
- [11] Dennis Frauen, Valentyn Melnychuk, and Stefan Feuerriegel. Model-agnostic meta-learners for estimating heterogeneous treatment effects over time. In *International Conference on Learning Representations (ICLR)*, 2025.
- [12] Changran Geng, Harald Paganetti, and Clemens Grassberger. Prediction of Treatment Response for Combined Chemo- and Radiation Therapy for Non-Small Cell Lung Cancer Patients Using a Bio-Mathematical Model. *Scientific Reports*, 2017. Author correction published in Scientific Reports, 2018, 8(1):12631, doi:10.1038/s41598-018-30761-7.
- [13] Miguel A. Hernán, Babette Brumback, and James M. Robins. Marginal structural models to estimate the joint causal effect of nonrandomized treatments. *Journal of the American Statistical Association*, 2001.
- [14] Konstantin Hess, Dennis Frauen, Valentyn Melnychuk, and Stefan Feuerriegel. G-transformer for conditional average potential outcome estimation over time. In *NeurIPS*, 2025. arXiv preprint arXiv:2405.21012.
- [15] Alistair E. W. Johnson, Tom J. Pollard, Lu Shen, Li wei H. Lehman, Mengling Feng, Mohammad Mahdi Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. MIMIC-III, a freely accessible critical care database. *Scientific Data*, 2016.
- [16] Edward H. Kennedy. Towards Optimal Doubly Robust Estimation of Heterogeneous Causal Effects. *Electronic Journal of Statistics*, 2023.

- [17] Sören R. Künzel, Jasjeet S. Sekhon, Peter J. Bickel, and Bin Yu. Meta-learners for estimating heterogeneous treatment effects using machine learning. *Proceedings of the National Academy of Sciences*, 116:4156–4165, 2019.
- [18] Greg Lewis and Vasilis Syrgkanis. Double/Debiased Machine Learning for Dynamic Treatment Effects via g-Estimation. In *Advances in Neural Information Processing Systems*, Virtual, Dec 2021.
- [19] Bryan Lim. Forecasting treatment responses over time using recurrent marginal structural networks. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- [20] Valentyn Melnychuk, Dennis Frauen, and Stefan Feuerriegel. Causal Transformer for Estimating Counterfactual Outcomes. In *ICML*, 2022.
- [21] Valentyn Melnychuk, Dennis Frauen, and Stefan Feuerriegel. Bounds on representation-induced confounding bias for treatment effect estimation. In *ICLR*, 2024. arXiv:2311.11321.
- [22] Liliana Orellana, Andrea Rotnitzky, and James M. Robins. Dynamic regime marginal structural mean models for estimation of optimal dynamic treatment regimes, part i: main content. *The International Journal of Biostatistics*, 2010.
- [23] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019.
- [24] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of Causal Inference: Foundations and Learning Algorithms*. MIT Press, 2017.
- [25] Maya L. Petersen and Mark J. van der Laan. Causal models and learning from data: integrating causal modeling and statistical estimation. *Epidemiology*, 2014.
- [26] J. M. Robins, M. A. Hernan, and B. Brumback. Marginal structural models and causal inference in epidemiology. *Epidemiology*, 2000.
- [27] James Robins. A new approach to causal inference in mortality studies with a sustained exposure period. *Mathematical Modelling*, January 1986.
- [28] James M. Robins. Correcting for non-compliance in randomized trials using structural nested mean models. *Communications in Statistics - Theory and Methods*, 1994.
- [29] James M. Robins. Optimal structural nested models for optimal sequential decisions. In *Proceedings of the Second Seattle Symposium in Biostatistics*. 2004.
- [30] James M. Robins, Sander Greenland, and Fu-Chang Hu. Estimation of the causal effect of a time-varying exposure on the marginal mean of a repeated binary outcome. *Journal of the American Statistical Association*, 1999.
- [31] Donald Rubin. Estimating causal effects of treatments in experimental and observational studies. *ETS Research Bulletin Series*, 1972.
- [32] Rui Li, Stephanie Hu, Mingyu Lu, Yuria Utsumi, Prithwish Chakraborty, Daby M. Sow, Piyush Madan, Jun Li, Mohamed Ghalwash, Zach Shahn, and Li-wei Lehman. G-Net: a recurrent network approach to G-Computation for counterfactual prediction under a dynamic treatment regime. *Machine Learning for Health*, 2021.
- [33] Peter Schulam and Suchi Saria. Reliable Decision Support using Counterfactual Models. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 6076–6086, Long Beach, CA, USA, 04–09 Dec 2017. Curran Associates, Inc.
- [34] Uri Shalit, Fredrik D. Johansson, and David Sontag. Estimating Individual Treatment Effect: Generalization Bounds and Algorithms. In Doina Precup and Yee Whye Teh, editors, *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 3076–3085. PMLR, 06–11 Aug 2017.

- 491 [35] Toru Shirakawa, Yi Li, Yulun Wu, Sky Qiu, Yuxuan Li, Mingduo Zhao, Hiroyasu Iso, and Mark  
492 van der Laan. Longitudinal targeted minimum loss-based estimation with temporal-difference  
493 heterogeneous transformer. *arXiv preprint arXiv:2404.04399*, apr 2024.
- 494 [36] Hossein Soleimani, Adarsh Subbaswamy, and Suchi Saria. Treatment-response models for  
495 counterfactual reasoning with continuous-time, continuous-valued interventions. In *UAI*, Sydney,  
496 Australia, 2017. AUAI Press.
- 497 [37] Stijn Vansteelandt and Marshall Joffe. Structural nested models and g-estimation: The partially  
498 realized promise. *Statistical Science*, 29, November 2014.
- 499 [38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
500 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information*  
501 *Processing Systems*, 30, 2017.
- 502 [39] Stefan Wager and Susan Athey and. Estimation and inference of heterogeneous treatment effects  
503 using random forests. *Journal of the American Statistical Association*, 2018.
- 504 [40] Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*, volume 48  
505 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press,  
506 Cambridge, UK, 2019.
- 507 [41] Shirly Wang, Matthew B. A. McDermott, Geeticka Chauhan, Marzyeh Ghassemi, Michael C.  
508 Hughes, and Tristan Naumann. MIMIC-Extract: a data extraction, preprocessing, and represen-  
509 tation pipeline for MIMIC-III. In *Proceedings of the ACM Conference on Health, Inference,*  
510 *and Learning*, ACM CHIL '20, April 2020.
- 511 [42] Yanbo Xu, Yanxun Xu, and Suchi Saria. A non-parametric bayesian approach for estimating  
512 treatment-response curves from sparse time series. In *ML4H*, Proceedings of Machine Learning  
513 Research, Northeastern University, Boston, MA, USA, aug 2016.
- 514 [43] Jinsung Yoon, James Jordon, and Mihaela van der Schaar. GANITE: Estimation of individual-  
515 ized treatment effects using generative adversarial nets. In *ICLR*, 2018.

## A Preliminaries

In this section we briefly introduce the preliminaries of causal inference.

### A.1 Causal modeling and treatment effect

**Structural causal model** A structural causal model is a tuple  $\mathcal{M} = (\mathcal{U}, \mathcal{V}, \mathcal{F}, P)$  where:

1.  $\mathcal{U}$  is a set of exogenous variables, which are not caused by any other variables, namely the unknown background factors or simply noises. Each  $U_i \in \mathcal{U}$  has a value domain denoted  $\text{dom}(U_i)$ .  $\mathcal{U}$  is defined on a probability space  $(\Omega, \mathcal{B}, P)$ .
2.  $\mathcal{V}$  is a set of endogenous variables that are determined within the model from other variables. The domain of  $V \in \mathcal{V}$  is  $\text{dom}(V)$ .
3.  $\mathcal{F}$  is a set of structural functions:  $\mathcal{F} = \{f_i : \text{dom}(\text{Pa}_i) \times \text{dom}(U_i) \rightarrow \text{dom}(V_i) \mid V_i \in \mathcal{V}\}$ , where  $\text{Pa}_i$  is the set of endogenous variable that directly influences  $V_i$ , and  $U_i \in \mathcal{U}$  is the exogenous variable that directly influences  $V_i$ . Here,  $f_i$  is a deterministic function that describes the structural causal mechanism for generating  $V_i$ :  $V_i = f_i(\text{Pa}_i, U_i)$ .

Furthermore, the causal relationships specified by  $\mathcal{F}$ 's  $\{\text{Pa}_i \mid V_i \in \mathcal{V}\}$  induce a directed graph  $\mathcal{G} = (\mathcal{V}, E)$ , for which nodes correspond to the endogenous variables and  $(V_i, V_j) \in E \iff V_i \in \text{Pa}_j$ . It is important to note that the causal graph  $\mathcal{G}$  must be a **directed acyclic graph (DAG)**, meaning that there should not be a causal cycle, which is contradictory to both real world and mathematics.

For SCM, the joint probability distribution  $P$  determines induces the probability on the whole variable sets. Once the set  $\mathcal{U}$ 's variables takes a fixed value, the whole endogenous set's values are determined through  $\mathcal{F}$ . Let  $P_M$  be the joint probability distribution over  $\mathcal{V}$ , then for any value  $v \in \prod_{V_i \in \mathcal{V}} \text{dom}(V_i)$ :

$$P_M(V = v) = \Pr(\{u \in \prod_{U_i \in \mathcal{U}} \text{dom}(U_i) \mid \mathcal{F}(u) = v\}) \quad (14)$$

**Intervention:** Intervention could be defined under the context of SCM as forcing a subset of variables (here we assume endogenous)  $X \subset \mathcal{V}$  to take a specific value  $x \in \prod_{V_i \in X} \text{dom}(V_i)$ <sup>5</sup>, independent of the original induced probability. This action simulates the process of controlled experiment or a hypothetical scenario.

An intervention, denoted as  $\text{do}(X = x)$ , creates a **new** SCM  $\mathcal{M}_x = (\mathcal{U}, \mathcal{V}, \mathcal{F}_x, P)$ , which is a mutation from the original SCM  $\mathcal{M}$ , where:

1. The set of endogenous  $\mathcal{V}$  and exogenous  $\mathcal{U}$  remain unchanged with the same probability space  $(\Omega, \mathcal{B}, P)$ .
2. For each  $V_i \in X$ , the old structural equation  $V_i = f_i(\text{Pa}_i, U_i)$  is replaced by a degenerated constant equation  $V_i \equiv x_i$ .
3. For all the other endogenous variables, the mapping  $f_i$  stays unchanged, but their distributions may adjust.

On the whole, the intervention  $\text{do}(X = x)$  induces a new probability distribution  $P_{M_x}$ . The interventional distribution of a variable, say  $Y \in \mathcal{V}$ , is denoted as  $\Pr(Y \mid \text{do}(X = x))$  or  $\Pr(Y^{\text{do}(X=x)})$ . Under certain context for simplicity, we could directly use  $\Pr(Y^{(x)})$  without ambiguity.

**Potential Outcome and Treatment effect:** From an intuition perspective, intervention is like cutting the causal links to  $X$  from the original diagram, setting the value to  $x$  and then observe how the system responds. We emphasize that  $\Pr(Y \mid \text{do}(X = x)) \neq \Pr(Y \mid X = x)$ , which is just a conditional distribution under the original SCM.

A **potential outcome (PO)** of a variable  $Y \in \mathcal{V}$  is the value under a hypothetical scenario, where we make an intervention  $\text{do}(X = x)$ , given a specific realization of  $U = u$ . We denote this as

<sup>5</sup>General intervention could be defined as setting some variables to follow a specific **distribution**[24], however we follow a common practice to assume the intervention means forcing a fixed value



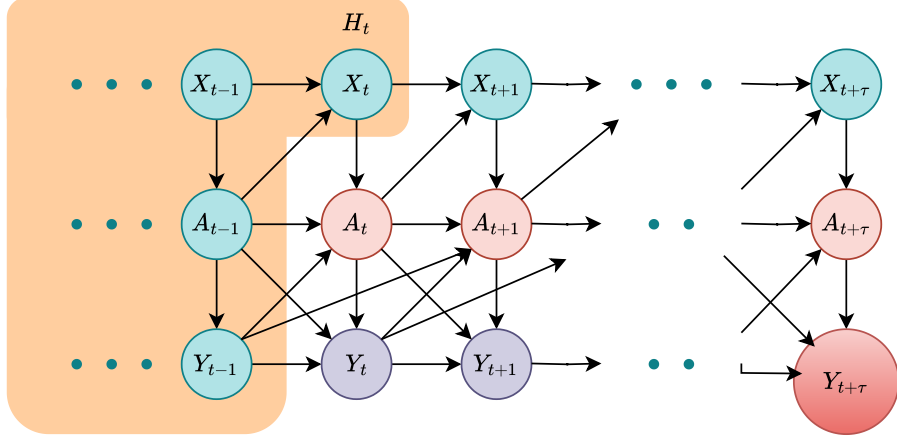


Figure 4: An example of causal diagram for the SCM in the setting of estimating CATE over time. The variable sequences are truncated at  $t + \tau$  (maximal length is  $T$ ). Exogenous variables  $\mathcal{U}$  are hidden.

558  $Y_x(u)$ . When  $U \sim P$ . This random variable  $Y_x(U)$  is just  $Y^{\text{do}(X=x)}$ , which follows the distribution  
559  $\Pr(Y \mid \text{do}(X = x))$ . The average potential outcome (APO) is the expectation of PO:  $\mathbb{E}[Y^{\text{do}(X=x)}]$ .  
560 A **conditional average potential outcome (CAPO)** is nothing but the expectation of the random  
561 variable  $Y_x(U)$  conditioned on some subgroups  $W = w$ , where  $W \subset \mathcal{V}$ . This conditioning on  $W = w$   
562 could be traced back to a posterior distribution on  $U$ , namely  $\Pr(U \mid W = w) := P_w$ , then the CAPO  
563 is just  $\mathbb{E}_{U_w \sim P_w}[Y_x(U_w)]$  or equivalently  $\mathbb{E}[Y^{(x)} \mid W = w]$ .

564 An **average treatment effect (ATE)** is the expectation difference between two POs under intervention  
565  $x_1$  and  $x_2$  respectively:  $\mathbb{E}[Y^{(x_1)} - Y^{(x_2)}]$ . A **conditional average treatment effect (CATE)** is then  
566 the difference between two CAPOs:  $\mathbb{E}[Y^{(x_1)} - Y^{(x_2)} \mid W = w]$ .

567 So far we have strictly defined these foundational concepts. Overall, the potential outcomes are  
568 counterfactual variables in intervened SCM while treatment effects are the contrast between these  
569 potential outcomes. Within this rigorous framework, we are able to clearly define the goal for our  
570 work on estimating CATE over time conditioned on patient history  $H_t$  as in Sec. 3.

## 571 A.2 SCM for CATE estimation over time

572 We formalize the setting for estimating CATE over time with the following structural causal model  
573  $\mathcal{M}$ :

- 574 1. Exogenous variables:  $\mathcal{U} = \{U_{X_t}, U_{A_t}, U_{Y_t}\}_{t=1}^T$ , where  $U_{X_t}, U_{A_t}, U_{Y_t}$  are mutually indepen-  
575 dent noise.
- 576 2. Endogenous variables are  $\mathcal{V} = \{X_t, A_t, Y_t\}_{t=1}^T$ . Denote  $H_t = \{X_s\}_{s=1}^t \cup \{A_j\}_{j=1}^T \cup$   
577  $\{Y_i\}_{i=1}^T$  as the history up until time  $t$ .
- 578 3. Structural functions  $\mathcal{F}$ : for  $t \in \{1, \dots, T\}$

$$\begin{aligned} X_t &= f_{X_t}(H_{t-1} \cup \{A_{t-1}\}, U_{X_t}) \\ A_t &= f_{A_t}(H_t, U_{A_t}) \\ Y_t &= f_{Y_t}(H_t, A_t, U_{Y_t}) \end{aligned}$$

579 These equations generate the full observed patient trajectory in the following temporal order:  $(X_1 \rightarrow$   
580  $A_1 \rightarrow Y_1 \rightarrow \dots \rightarrow X_T \rightarrow A_T \rightarrow Y_T)$ .

581 A possible instantiation of temporal graph of this given structural causal model is presented in Figure.4  
582 When talking about SCMs, it is important to point out that SCM represents perfect knowledge of the  
583 data generating process. In real-world settings, there are possibly omitted variables (e.g. unmeasured  
584 confounders).

585 **Time-varying treatment and potential outcomes:** A **time-varying intervention** is intervention  
 586 on a continuous sequence of treatment variables  $\text{do}(A_{t:t+\tau} = a_{t:t+\tau})$ , which modifies the structural  
 587 equations for  $A_t, \dots, A_{t+\tau}$ :

$$\forall k \in \{t, \dots, t + \tau\} : \quad A_k \equiv a_k.$$

588 This induces a new SCM  $\mathcal{M}^{(a_{t:t+\tau})}$ , generating the potential outcome:

$$Y_{t+\tau}^{(a_{t:t+\tau})} = f_{Y_{t+\tau}} \left( H_{t+\tau}^{(a_{t:t+\tau})}, X_{t+\tau}^{(a_{t:t+\tau})}, a_{t+\tau}, U_{Y_{t+\tau}} \right),$$

589 where  $H_{t+\tau}^{(a_{t:t+\tau})}, X_{t+\tau}^{(a_{t:t+\tau})}$  are recursively computed under the intervention using structural equations  
 590 defined in 3.

591 **Identification assumptions:** The structural functions, although might be unspecified, already  
 592 contains rich information about the data generating process. It is easy to verify that under the  
 593 specified data generating process and the mutual independence between the exogenous variables, the  
 594 following two conditions are satisfied:

- 595 1. **Consistency:** If  $A_{t:t+\tau} = a_{t:t+\tau}$  is the observed treatment sequence for a given patient  
 596  $u$ , then  $Y_{t+\tau}^{(a_{t:t+\tau})}(u) = Y(u)$ . Or, equivalently  $Y_{t+\tau}^{(A_{t:t+\tau})} = Y_{t+\tau}$ . This means when we  
 597 intervene with the observed treatment, the potential outcome equals the observed outcome.
- 598 2. **Sequential ignorability:** The potential outcome under a fixed intervention is conditionally  
 599 independent with the treatment assignment:

$$Y_t^{(a_t)} \perp A_t \mid H_t, \quad \forall a_t, 1 \leq t \leq T$$

600 This ignorability assumption implies that there are no hidden confounders that affect both  
 601 outcome and the treatment, which is already decided with the SCM.

- 602 3. **Sequential overlap:** If  $\Pr(H_t = h_t) > 0$ , then:

$$\Pr(A_t = a_t \mid H_t = h_t) > 0 \quad \forall a_t \in \mathcal{A}_t, 1 \leq t \leq T.$$

603 These three assumptions together allow the use of observational data to compute causal quantities,  
 604 such as potential outcome over time and treatment effect over time [26]. Without these assumptions,  
 605 exact identification would be impossible. While these assumptions themselves are untestable, they  
 606 could be justified under appropriate study design and domain knowledge.

607 **Time-varying confounding:** Estimating treatment effect over time poses a significant challenge  
 608 due the existence of time-varying confounding [6], which differentiate the task from the static setting  
 609 ( $\tau = 0$ ). The underlying reason is that, the future time-varying confounders, such as  $X_{t+1:t+\tau}$   
 610 or  $Y_{t:t+\tau-1}$ , are *unobserved* during the inference phase. As a result, the naïve adjustment with the  
 611 history variable  $H_t$  is not sufficient for the backdoor criteria [24] and, thus, leads to biased estimation:

$$\text{When } \tau > 0 : \quad \mathbb{E} \left[ Y_{t+\tau}^{(a_{t:t+\tau})} \mid H_t = h_t \right] \neq \mathbb{E} [Y_{t+\tau} \mid A_{t:t+\tau} = a_{t:t+\tau}, H_t = h_t]. \quad (15)$$

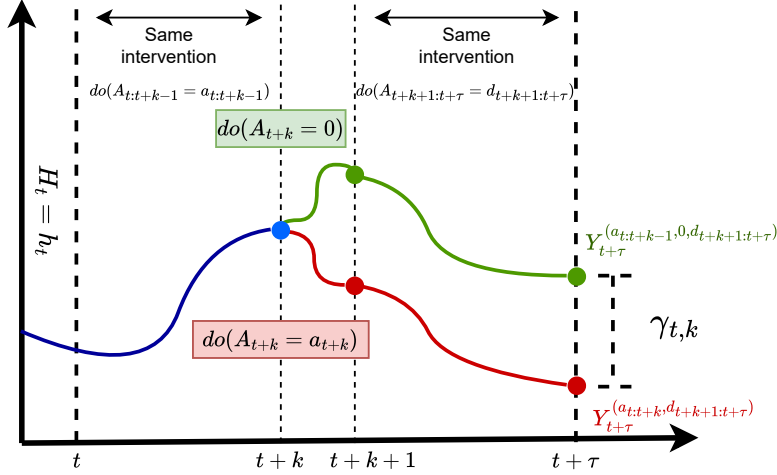


Figure 5: Visualization of the blip function  $\gamma_{t,k}$ . The blip function quantifies the localized treatment effect at each time step.

## 612 B Theory of structural nested mean models

613 In this section, we show by strict causal reasoning how to use the structural nested mean models  
 614 to estimate CATE over time. The theory consists of the foundations of SNMM [29] and a specific  
 615 realization [18], which we adopt for our method. We show the mathematical derivations for complete-  
 616 ness. Then in Subsec. B.2 we show that our  $L^1$ -moment loss is equivalent with the exact identification  
 617 process of SNMM [29, 18]. Next, in Subsec. B.3 we show that our loss is Neyman orthogonal.  
 618 Finally, we briefly explain the mean squared error rate for the estimated blip coefficients which is  
 619 proved by previous works in SNMM [18].

620 To simplify the notation, we denote  $Y = Y_{t+\tau}$ . And we use  $\underline{W}_{t+k}$  to denote  $\underline{W}_{t+k:t+\tau}$  for any  
 621 variable  $W \in \mathcal{V}$ .

### 622 B.1 SNMM

623 **Identification:** The blip functions can be accumulated to identify the potential outcomes, as  
 624 summarized by the lemma [29, 18] below. We show the proof for completeness.

625 **Lemma 1. Identification via blip functions [29, 18]** Given a policy  $d$  between time  $t$  and  $t + \tau$ , the  
 626 following identification holds under the sequential ignorability assumption in 2:

$$\mathbb{E}[Y^{(d)} \mid H_t = h_t] = \mathbb{E}\left[Y + \sum_{k=0}^{\tau} \rho_{t,k}(X_{t+1:t+k}, A_{t:t+k}; h_t) \mid H_t = h_t\right] \quad (16)$$

627 where  $\rho_{t,k}$  is defined as:

$$\rho_{t,k}(X_{t+1:t+k}, A_{t:t+k}; h_t) = \gamma_{t,k}(X_{t+1:t+k}, (A_{t:t+k-1}, d_{t+k}); h_t) - \gamma_{t,k}(X_{t+1:t+k}, A_{t:t+k}; h_t).$$

628 *Proof.* By the definition of  $\gamma_{t,k}$ , we could write  $\gamma_{t,k}$  into the following form:

$$\begin{aligned} & \gamma_{t,k}(x_{t+1:t+k}, (a_{t:t+k-1}, d_{t+k}); h_t) \\ & \stackrel{\text{Def.}}{=} \mathbb{E}\left[Y^{(a_{t:t+k-1}, d_{t+k})} - Y^{(a_{t:t+k-1}, 0, d_{t+k+1})} \mid X_{t+1:t+k} = x_{t+1:t+k}, A_{t:t+k-1} = a_{t:t+k-1}, A_{t+k} = d_{t+k}, H_t = h_t\right] \\ & \stackrel{\text{Ignor.}}{=} \mathbb{E}\left[Y^{(a_{t:t+k-1}, d_{t+k})} - Y^{(a_{t:t+k-1}, 0, d_{t+k+1})} \mid X_{t+1:t+k} = x_{t+1:t+k}, A_{t:t+k} = a_{t:t+k}, H_t = h_t\right] \end{aligned} \quad (17)$$

629 Similarly:

$$\begin{aligned} & \gamma_{t,k}(x_{t+1:t+k}, a_{t:t+k}; h_t) \\ \stackrel{\text{Def.}}{=} & \mathbb{E} \left[ Y^{(a_{t:t+k}, \underline{d}_{t+k+1})} - Y^{(a_{t:t+k-1}, 0, \underline{d}_{t+k+1})} \mid X_{t+1:t+k} = x_{t+1:t+k}, A_{t:t+k} = a_{t:t+k}, H_t = h_t \right] \end{aligned} \quad (18)$$

630 This implies that:

$$\begin{aligned} & \rho_{t,k}(X_{t+1:t+k}, A_{t:t+k}; h_t) \\ &= \gamma_{t,k}(X_{t+1:t+k}, (A_{t:t+k-1}, d_{t+k}); h_t) - \gamma_{t,k}(X_{t+1:t+k}, A_{t:t+k}; h_t) \\ & \stackrel{\text{By 17,18}}{=} \mathbb{E} \left[ Y^{(A_{t:t+k-1}, \underline{d}_{t+k})} - Y^{(A_{t:t+k}, \underline{d}_{t+k+1})} \mid X_{t+1:t+k}, A_{t:t+k}, H_t = h_t \right] \end{aligned} \quad (19)$$

631 By consistency 1,  $Y^{(A_{t:t+\tau})} = Y$ , and that by towering law of the conditional expectation:

$$\begin{aligned} & \mathbb{E} \left[ \rho_{t,k}(X_{t+1:t+k}, A_{t:t+k}; h_t) \mid H_t = h_t \right] \\ &= \mathbb{E} \left[ \mathbb{E} \left[ Y^{(A_{t:t+k-1}, \underline{d}_{t+k})} - Y^{(A_{t:t+k}, \underline{d}_{t+k+1})} \mid X_{t+1:t+k}, A_{t:t+k}, H_t = h_t \right] \mid H_t = h_t \right] \\ &= \mathbb{E} \left[ Y^{(A_{t:t+k-1}, \underline{d}_{t+k})} - Y^{(A_{t:t+k}, \underline{d}_{t+k+1})} \mid H_t = h_t \right] \end{aligned} \quad (20)$$

632 Summing up (20) for  $k = 0, 1, \dots, \tau$ , we observe that all the middle terms cancel off:

$$\begin{aligned} & \mathbb{E} \left[ \sum_{k=0}^{\tau} \rho_{t,k}(X_{t+1:t+k}, A_{t:t+k}; h_t) \mid H_t = h_t \right] \\ &= \mathbb{E} \left[ \sum_{k=0}^{\tau} Y^{(A_{t:t+k-1}, \underline{d}_{t+k})} - Y^{(A_{t:t+k}, \underline{d}_{t+k+1})} \mid H_t = h_t \right] \\ &= \mathbb{E} \left[ Y^{(d)} - Y \mid H_t = h_t \right] \end{aligned} \quad (21)$$

633 This is equivalent with the Eq. 16. □

634 Lemma 1 shows that, in order to estimate the potential outcomes, it suffices to estimate the blip  
635 functions. Now the identification problem is reduced into correctly estimating these blip func-  
636 tions. Suppose  $\gamma_t^\dagger(\cdot, \cdot; h_t) = (\gamma_{t,0}^\dagger, \gamma_{t,1}^\dagger, \dots, \gamma_{t,\tau}^\dagger)$  is an **arbitrary** candidate function that takes  
637  $x_{t+1:t+k}, a_{t:t+k}$  as input and is subject to the following conditions:

$$\forall k \in [0, \tau], \quad \gamma_{t,k}^\dagger(x_{t+1:t+k}, (a_{t:t+k-1}, 0); h_t) = 0. \quad (22)$$

638 Let:

$$H_{t,k}(\gamma_t^\dagger) = Y + \sum_{j=k}^T \left( \gamma_{t,k}^\dagger(X_{t+1:t+j}, (A_{t:t+j-1}, d_{t+j})) - \gamma_{t,k}^\dagger(X_{t+1:t+j}, A_{t:t+j}) \right) \quad (23)$$

639 Then, with Lemma 1 and the sequential conditional ignorability, theorem 3.2 from [29] states that the  
640 moment restriction below is satisfied. Here we present it as a crucial lemma and show the proof for  
641 completeness.

642 **Lemma 2. Moment restriction for blip functions** When (22) is satisfied, then for arbitrary function  
643  $S_k$ , the following two conditions are equivalent:

- 644 1.  $\gamma_{t,k}^\dagger(\cdot, \cdot; h_t)$  is the true  $\gamma_{t,k}(\cdot, \cdot; h_t)$ .
- 645 2. The following moment restriction is satisfied:

$$\mathbb{E} \left[ H_{t,k}(\gamma_t^\dagger)(S_k(X_{t+1:t+k}, A_{t:t+k}) - \mathbb{E}[S_k(X_{t+1:t+k}, A_{t:t+k}) \mid X_{t+1:t+k}, A_{t:t+k-1}]) \mid H_t = h_t \right] = 0 \quad (24)$$

646 *Proof.* We only prove the direction of  $1 \rightarrow 2$ , which already conveys the core thought of the moment  
 647 restriction.

We denote:

$$R_k(X_{t+1:t+k}, A_{t:t+k}) = S_k(X_{t+1:t+k}, A_{t:t+k}) - \mathbb{E}[S_k(X_{t+1:t+k}, A_{t:t+k}) \mid X_{t+1:t+k}, A_{t:t+k-1}, H_t = h_t]$$

648 Observing the LHS of (24):

$$\begin{aligned} \text{L.H.S.} &= \mathbb{E} \left[ H_{t,k}(\gamma_t^\dagger) R_k(X_{t+1:t+k}, A_{t:t+k}) \mid H_t = h_t \right] \\ &\stackrel{\text{Cond.}}{=} \mathbb{E} \left[ \mathbb{E} [H_{t,k}(\gamma_t) R_k(X_{t+1:t+k}, A_{t:t+k}) \mid X_{t+1:t+k}, A_{t:t+k}] \mid H_t = h_t \right] \\ &= \mathbb{E} \left[ \mathbb{E} [H_{t,k}(\gamma_t) \mid X_{t+1:t+k}, A_{t:t+k}] \cdot R_k(X_{t+1:t+k}, A_{t:t+k}) \mid H_t = h_t \right] \\ &\stackrel{\text{Lem. 1}}{=} \mathbb{E} \left[ \mathbb{E} [Y^{(A_{t:t+k-1}, \underline{d}_{t+k})} \mid X_{t+1:t+k}, A_{t:t+k}] \cdot R_k(X_{t+1:t+k}, A_{t:t+k}) \mid H_t = h_t \right] \\ &\stackrel{\text{Ignor.}}{=} \mathbb{E} \left[ \mathbb{E} [Y^{(A_{t:t+k-1}, \underline{d}_{t+k})} \mid X_{t+1:t+k}, A_{t:t+k-1}] \cdot R_k(X_{t+1:t+k}, A_{t:t+k}) \mid H_t = h_t \right] \\ &\stackrel{\text{Tower}}{=} \mathbb{E} \left[ \mathbb{E} [Y^{(A_{t:t+k-1}, \underline{d}_{t+k})} \mid X_{t+1:t+k}, A_{t:t+k-1}] \cdot \mathbb{E} [R_k(X_{t+1:t+k}, A_{t:t+k}) \mid X_{t+1:t+k}, A_{t:t+k-1}] \mid H_t = h_t \right] \\ &\stackrel{\text{Def.}}{=} \mathbb{E} \left[ \mathbb{E} [Y^{(A_{t:t+k-1}, \underline{d}_{t+k})} \mid X_{t+1:t+k}, A_{t:t+k-1}] \cdot 0 \mid H_t = h_t \right] = 0 \end{aligned} \quad (25)$$

649

□

650 To further reduce the complexity of the moment restriction to achieve locally efficient doubly robust  
 651 estimator (see Section 3.3 from [29]), we could subtract  $H_{t,k}(\gamma^*)$ (24) by its conditional expectation  
 652  $\mathbb{E}[H_{t,k}(\gamma_t^\dagger) \mid X_{t+1:t+k}, A_{t:t+k-1}]$ :

$$\mathbb{E} \left[ (H_{t,k}(\gamma_t^\dagger) - \mathbb{E}[H_{t,k}(\gamma_t^\dagger) \mid X_{t+1:t+k}, A_{t:t+k-1}]) R_k(X_{t+1:t+k}, A_{t:t+k}) \mid H_t = h_t \right] = 0 \quad (26)$$

653 **Parametrization of the blip function:** To achieve parametric rates of the estimation, we use the  
 654 common semi-parametrization of the blip function<sup>6</sup> [29, 18]:

$$\gamma_{t,k}(x_{t+1:t+k}, a_{t:t+k}; h_t) = \psi_{t,k}(h_t)' a_{t+k} \quad (27)$$

655 Here  $\psi_{t,k}$  is a function that maps a given  $H_t = h_t$  to  $\psi_{t,k}(h_t)$ , which is the parameter of the blip  
 656 function  $\gamma_{t,k}(\cdot, \cdot; h_t)$ . Intuitively, these  $\psi_{t,k}(h_t)$  works as the coefficients in a linear parametrization.  
 657 Hence, we call them in our method *blip coefficients* and the mapping  $\psi_{t,k}(\cdot)$  *blip coefficient predictor*  
 658 or *blip coefficient estimator*.

659 **Deriving the moment restriction of SNMM:** We abuse the notation of  $H_{t,k}(\gamma_t^*)$  with using the  
 660 parameter  $\psi_t$  to represent the  $H_{t,k}(\gamma_t^{\psi_t})$  as  $H_{t,k}(\psi_t)$ . For convenience, we introduce several notation  
 661 for residuals:

$$\tilde{A}_{t,j,k} = A_{t+j} - \mathbb{E}[A_{t+j} \mid X_{t+1:t+k}, A_{t:t+k}, H_t = h_t], \quad 0 \leq k \leq j \leq \tau \quad (28)$$

$$\tilde{Y}_{t,k} = Y - \mathbb{E}[Y \mid X_{t+1:t+k}, A_{t:t+k}, H_t = h_t], \quad 0 \leq k \leq \tau \quad (29)$$

662 Then we could derive the following (with definition of  $H_{t,k}(\cdot)$  in (23)), which is an adaptation from  
 663 [18]:

$$\mathbb{E} \left[ H_{t,k}(\psi_t) - \mathbb{E}[H_{t,k}(\psi_t) \mid X_{t+1:t+k}, A_{t:t+k-1}] \mid H_t = h_t \right] = \tilde{Y}_{t,k} - \sum_{j=k}^{\tau} \psi'_{t,j} \tilde{A}_{t,j,k} \quad (30)$$

664 Moreover, we specify  $S_k(x_{t+1:t+k}, a_{t:t+k-1}) = a_{t+k}$  from Lemma 2, then:

$$\begin{aligned} &R_k(X_{t+1:t+k}, A_{t:t+k-1}) \\ &= S_k(X_{t+1:t+k}, A_{t:t+k}) - \mathbb{E}[S_k(X_{t+1:t+k}, A_{t:t+k}) \mid X_{t+1:t+k}, A_{t:t+k-1}, H_t = h_t] = \tilde{A}_{t,j,k} \end{aligned} \quad (31)$$

<sup>6</sup>In our paper, we choose the simplest parametrization for clarity. In fact, the blip functions could have more sophisticated parametrizations, such as  $\gamma_{t,k}(x_{t+1:t+k}, a_{t:t+k}; h_t) = \psi_{t,k}(h_t)' \phi_{t,k}(x_{t+1:t+k}, a_{t:t+k}; h_t)$ , where  $\phi$  is a nonlinear, him-dimensional feature map [29, 18].



Therefore, the moment restriction in (26) could be simplified as:

$$\forall k \in \{0, \dots, \tau\} : \mathbb{E} \left[ \left( \tilde{Y}_{t,k} - \sum_{j=k+1}^{\tau} \psi_{t,j}(h_t)' \tilde{A}_{t,j,k} - \psi_{t,k}(h_t)' \tilde{A}_{t,k,k} \right) \tilde{A}_{t,k,k} \mid H_t = h_t \right] = 0 \quad (32)$$

which is the primary moment restriction for our setting, adapted from [18]. To estimate  $\psi_t(h_t)$  we need to solve the set of equations defined in (32) in an iterative way (See Algorithm 3 from [18]):

1. Step 0: for  $k = \tau$ , solve  $\hat{\psi}_{t,\tau}(h_t) \in \mathbb{R}^r$
2. Step  $\tau - k$ , ( $k \geq 1$ ), With the given  $\hat{\psi}_{t,j}(h_t)$ ,  $j > k$ , solve  $\hat{\psi}_{t,k}(h_t)$ .

Each step we need to solve a linear equation, which has a unique solution under standard assumptions (see theorem 8 from [18]). And since the solution is unique, the solution equals the true blip coefficients. However, we argue that this classical method has limitations due to its sequential nature which prevents us from batch training. Noticing that solving (32) is equivalent to solving this minimization problem, we successfully establish the estimation of  $\psi_t$  under risk minimization paradigm, suitable for training with neural networks:

$$\psi_{t,k}^* = \arg \min_{\hat{\psi}_{t,k}(\cdot) \in \Phi_{t,k}} \mathbb{E} \left[ \left\| \mathbb{E} \left[ \left( \tilde{Y}_{t,k} - \sum_{j=k+1}^{\tau} \psi_{t,j}^*(h_t)' \tilde{A}_{t,j,k} - \hat{\psi}_{t,k}(h_t)' \tilde{A}_{t,k,k} \right) \tilde{A}_{t,k,k} \mid H_t \right] \right\|_1 \right] \quad (33)$$

**Estimating CATE over time:** Once the parameters  $\psi_t(h_t) = (\psi_{t,0}(h_t), \dots, \psi_{t,\tau}(h_t))$  are estimated, we could proceed to estimate the potential outcome from Eq. 16:

$$\begin{aligned} & \mathbb{E} \left[ Y^{(d)} \mid H_t = h_t \right] \\ &= \mathbb{E} \left[ Y + \sum_{k=0}^{\tau} \rho_{t,k} (X_{t+1:t+k}, A_{t:t+k}; h_t) \mid H_t = h_t \right] \\ &= \mathbb{E} \left[ Y + \sum_{k=0}^{\tau} \hat{\psi}_{t,k}(h_t)' (d_{t+k} - A_{t+k}) \mid H_t = h_t \right]. \end{aligned} \quad (34)$$

Then the CATE over time under treatment sequences  $a^*$  and  $b^*$  is computed by taking difference between the potential outcomes respectively [28, 29, 18]:

$$\begin{aligned} & \mathbb{E} \left[ Y^{(a^*)} - Y^{(b^*)} \mid H_t = h_t \right] \\ &= \mathbb{E} \left[ \sum_{k=0}^{\tau} \psi_{t,k}(h_t)' [(a_{t+k}^* - A_{t+k}) - (b_{t+k}^* - A_{t+k})] \mid H_t = h_t \right] \\ &= \mathbb{E} \left[ \sum_{k=0}^{\tau} \psi_{t,k}(h_t)' (a_{t+k}^* - b_{t+k}^*) \mid H_t = h_t \right] \\ &= \sum_{k=0}^{\tau} \psi_{t,k}(h_t)' (a_{t+k}^* - b_{t+k}^*) \end{aligned} \quad (35)$$

The expectation term vanishes since all the terms left are constant in the conditional expectation, enabling us to directly estimate CATE over time.

## B.2 Equivalence of the minimization scheme

Our proposed  $L^1$ -moment loss  $\mathcal{L}_D$  is:

$$\mathcal{L}_D(\theta) = \sum_{k=0}^{\tau} \mathcal{L}_{D,k}(\theta), \quad (36)$$

where  $\mathcal{L}_{D,k}(\theta) = \mathbb{E} \left[ \left\| \mathbb{E} \left[ \left( \tilde{Y}_{t,k} - \sum_{j=k+1}^{\tau} \psi_{t,j}(h_t; \theta)' \tilde{A}_{t,j,k} - \hat{\psi}_{t,k}(h_t; \theta)' \tilde{A}_{t,k,k} \right) \tilde{A}_{t,k,k} \mid H_t \right] \right\|_1 \right]$ , and  $\theta$  is the parameter of the blip coefficient predictor  $\psi_t(\cdot; \theta) = (\psi_{t,0}(\cdot; \theta), \dots, \psi_{t,\tau}(\cdot; \theta))$ . Under the parametrization assumption in Eq. 27, we can assume that  $\psi_{t,k} \in \Phi_{t,k}$ .

687 **Note:** We point out that  $\psi_{t,k}(\cdot)$  is not  $\psi_{t,k}(\cdot; \theta)$ . The former is the blip coefficient mapping of the  
 688 **true** blip function defined from the distribution and expectation, while the latter is just a function  
 689 with parameter  $\theta$  from the functional space  $\Phi_{t,k}$ . Moreover, the minimizer of an optimization has  
 690 subscript  $*$ .

691 **Our objective:** We want to show that the blip coefficient estimator parametrized by the minimizer  
 692 of the problem:

$$\theta^* = \arg \min_{\theta \in \Theta} \mathcal{L}_D(\theta) \quad (37)$$

693 equals the true blip coefficient mapping  $\psi_{t,k}$  almost surely:

$$\psi_{t,k}(h_t; \theta^*) = \psi_{t,k}(h_t) \quad \text{a.s.} \quad (38)$$

694 Since we are working in the field of machine learning, we treat "almost surely equal" and "constantly  
 695 equal" as the same thing and do not differentiate these two concepts from now on.

696 *Proof.* The proof is straightforward. We assume the true blip coefficient mappings  $\{\psi_{t,k}\}_{k=0}^\tau$  are  
 697 parametrized by  $\theta^*$ . Since  $\theta^* \in \Theta$  is the minimizer, then:

$$\mathcal{L}_D(\theta^*) = \sum_k \mathcal{L}_{D,k}(\theta^*) \leq \mathcal{L}_D(\theta^t) \stackrel{\text{Eq. 32}}{=} 0, \quad (39)$$

698 which means  $\forall k \in [0, \tau]$  :

$$\mathcal{L}_{D,k}(\theta^*) = \mathbb{E} \left[ \left\| \mathbb{E} \left[ (\tilde{Y}_{t,k} - \sum_{j=k+1}^\tau \psi_{t,j}(h_t; \theta^*)' \tilde{A}_{t,j,k} - \psi_{t,k}(h_t; \theta^*)' \tilde{A}_{t,k,k}) \tilde{A}_{t,k,k} \mid H_t \right] \right\|_1 \right] = 0. \quad (40)$$

699 Since the  $L^1$ -norm  $\|\cdot\|_1$  always produces non-negative values, then it must be that  $\forall h_t \in \mathcal{H}_t, k \in$   
 700  $[0, \tau]$ :

$$\mathbb{E} \left[ (\tilde{Y}_{t,k} - \sum_{j=k+1}^\tau \psi_{t,j}(h_t; \theta^*)' \tilde{A}_{t,j,k} - \psi_{t,k}(h_t; \theta^*)' \tilde{A}_{t,k,k}) \tilde{A}_{t,k,k} \mid H_t \right] = 0. \quad (41)$$

701 This means the minimizer of our  $L^1$  loss also produces a solution to the same set of moment restrions  
 702 as the true blip coefficients in Eq. 32. By the uniqueness of the solution for Eq. 32 under certain  
 703 regularity conditions (see theorem 8 from [18]), we reach the conclusion that:

$$\forall h_t \in \mathcal{H}_t, \quad \psi_{t,k}(h_t; \theta^*) = \psi_{t,k}(h_t) \quad (42)$$

704

□

### 705 B.3 Neyman orthogonality of the $L^1$ -moment loss

706 We show the orthogonality of the defined  $L^1$ -moment loss. Neyman orthogonality of the moment  
 707 equations for solving  $\psi$  as constants is already proved in previous SNMM theory (Lemma 14 from  
 708 [18]). Here we show, with a step forward, that the gradient of our  $L^1$ -moment is Neyman orthogonal  
 709 w.r.t. the nuisance function  $h = (p_{t,k}, \{q_{t,j,k}\}_{j=k}^\tau)$ . We note that some parts of the proof are  
 710 straightforward adaptations from [18].

711 The Fréchet derivative of any functional  $\mathcal{L}(f)$  is defined as:

$$\forall v, \quad D_f \mathcal{L}(f)[v] = \left. \frac{\partial}{\partial t} \mathcal{L}(f + tv) \right|_{t=0}. \quad (43)$$

712 *Proof.* Recall that our  $L^1$ -moment loss is:

$$\begin{aligned} \mathcal{L}_{D,k}(\theta; h) = \mathbb{E} \left[ \left\| \mathbb{E} \left[ (Y - p_{t,k}(H_{t+k}) - \sum_{j=k+1}^\tau \psi_{t,j}(H_t; \theta)' (A_{t+j} - q_{t,j,k}(H_{t+k})) \right. \right. \right. \\ \left. \left. \left. - \psi_{t,k}(H_t; \theta)' (A_{t+k} - q_{t,k,k}(H_{t+k})) \right) (A_{t+k} - p_{t,k,k}(H_{t+k})) \mid H_t \right] \right\|_1 \right], \end{aligned} \quad (44)$$

713 where  $h = (p, q)$  is the nuisance parameter,  $\theta$  is the trainable weight of the blip prediction model  
 714 from Stage 2. Our objective is to show that:

$$D_h(\partial_\theta \mathcal{L}_{D,k}(\theta; h)) = 0 \quad (45)$$

715 where  $D_h$  is the Fréchet derivative w.r.t. the nuisance function  $h$ .

716 We start by fixing  $H_t = h_t$  and denote the inner conditional moment as

$$\begin{aligned} \phi_k(h_t; \theta, h) := & \mathbb{E} \left[ (Y - p_{t,k}(H_{t+k}) - \sum_{j=k+1}^{\tau} \psi_{t,j}(h_t; \theta)' (A_{t+j} - q_{t,j,k}(H_{t+k})) \right. \\ & \left. - \psi_{t,k}(h_t; \theta)' (A_{t+k} - q_{t,k,k}(H_{t+k})) (A_{t+k} - p_{t,k,k}(H_{t+k})) \mid H_t = h_t \right], \end{aligned} \quad (46)$$

717 so that  $\mathcal{L}_{D,k}(\theta; h) = \mathbb{E}[\|\phi_k(H_t; \theta, h)\|_1]$ .

718 First, we derive the gradient of  $\phi_k$  over  $\theta$ . Note that the nuisance functions  $p_{t,k}, q_{t,j,k}$  does not depend  
 719 on the model weight  $\theta$ , the gradient is then:

$$\begin{aligned} \partial_\theta \phi_k(h_t; \theta, h) = & \mathbb{E} \left[ - \sum_{j=k+1}^{\tau} \partial'_\theta \psi_{t,j}(h_t; \theta)' (A_{t+j} - q_{t,j,k}(H_{t+k})) \right. \\ & \left. - \partial'_\theta \psi_{t,k}(h_t; \theta)' (A_{t+k} - q_{t,k,k}(H_{t+k})) (A_{t+k} - p_{t,k,k}(H_{t+k})) \mid H_t = h_t \right]. \end{aligned} \quad (47)$$

720 Then we proceed to show that the Fréchet derivative of  $\phi_k$  over  $p_{t,k}, q_{t,j,k}$  (where  $j > k$ , and  $q_{t,k,k}$   
 721 are all zero. hence its Fréchet derivative along any direction  $v_p$  is 0. Let  $h^* = (p^*, q^*)$  denote the true  
 722 nuisance functions with  $p_{t,k}^*(H_{t+k}) = \mathbb{E}[Y_{t+\tau} \mid H_{t+k}]$  and  $q_{t,j,k}^*(H_{t+k}) = \mathbb{E}[A_{t+j} \mid H_{t+k}]$ . (Note  
 723  $H_{t+k} = (\bar{X}_{t+k}, \bar{A}_{t+k-1}, \bar{Y}_{t+k-1})$ .)

724 (i) **Orthogonality to  $p_{t,k}$ .** Let  $\Delta_k = p_{t,k} - p_{t,k}^*$ , then the Fréchet derivative over the single nuisance  
 725 term  $p_{t,k}$  is:

$$D_{p_{t,k}} \partial_\theta \phi_k(h_t; \theta, h^*)[\Delta_k] = -\mathbb{E}[\Delta_k(H_{t+k}) \cdot (A_{t+k} - q_{t,k,k}^*(H_{t+k})) \mid H_t = h_t] = 0 \quad (48)$$

726 (ii) **Orthogonality to  $q_{t,j,k}$  for  $j > k$ .** Let  $\delta_{j,k} = q_{t,j,k} - q_{t,j,k}^*$ . Then the directional derivative at  
 727  $t = 0$  equals

$$\begin{aligned} D_{q_{t,j,k}} \partial_\theta \phi_k(h_t; \theta, h^*)[\delta_{j,k}] = & \mathbb{E} \left[ \partial'_\theta \psi_{t,j}(h_t; \theta)' \delta_{j,k}(H_{t+k}) (A_{t+k} - p_{t,k,k}(H_{t+k})) \mid H_t = h_t \right] \\ = & \partial'_\theta \psi_{t,j}(h_t; \theta)' \mathbb{E}[\delta_{j,k}(H_{t+k}) \cdot (A_{t+k} - p_{t,k,k}(H_{t+k})) \mid H_t = h_t] = 0. \end{aligned} \quad (49)$$

728 (iii) **Orthogonality to  $q_{t,k,k}$ .** Let  $\delta_k = q_{t,k,k} - q_{t,k,k}^*$ . Denote:

$$R_k(q_{t,k,k}) = Y - p_{t,k}(H_{t+k}) - \sum_{j=k+1}^{\tau} \partial_\theta \psi'_{t,j}(h_t; \theta)' (A_{t+j} - q_{t,j,k}^*(H_{t+k})) - \partial_\theta \psi_{t,k}(h_t; \theta)' (A_{t+k} - q_{t,k,k}(H_{t+k})) \quad (50)$$

729 We substitute  $q_{t,k,k}$  by  $q_{t,k,k}^* + t\delta_k$  into  $R_k$  and then take derivative over  $t$  at  $t = 0$ :

$$\frac{\partial}{\partial t} R_k(q_{t,k,k}^* + t\delta_k) \Big|_{t=0} = \partial'_\theta \psi_{t,k}(h_t; \theta)' \delta_k(H_{t+k}). \quad (51)$$

730 Now we take the directional Fréchet derivative over  $q_{t,k,k}$ :

$$D_{q_{t,k,k}} \partial_\theta \phi_k(h_t; \theta; h^*)[\delta_k] = \mathbb{E} \left[ \left( \frac{\partial}{\partial t} R_k \right) (A_{t+k} - q_{t,k,k}^*(H_{t+k})) - R_k(q_{t,k,k}^*) \delta_k(H_{t+k}) \mid H_t = h_t \right] \quad (52)$$

731 Substituting Eq. 51 into Eq. 52 and note the fact that  $\mathbb{E}[R_k(q_{t,k,k}^* \mid H_{t+k})] = 0$ , we have  
732  $D_{q_{t,k,k}} \partial_\theta \phi_k(h_t, \theta; h^*)[\delta_k] = 0$ .

733 Putting (i),(ii), and (iii) together, we prove that the overall Fréchet derivative of  $\phi_k$  over nuisance  
734 parameter  $h$  is zero for any given condition  $H_t = h_t$ .

735 Because  $\|\cdot\|_1$  is almost-surely differentiable, and for every  $h_t$  the derivative of  $\partial_\theta \phi_k$  in any nuisance  
736 direction is 0, integrating over  $H_t$  preserves the result:

$$\begin{aligned} D_h \frac{\partial \mathcal{L}_{D,k}}{\partial \theta}(\theta^*; h^*)[h - h^*] &= \mathbb{E} \left[ D_h \left\| \frac{\partial \phi_k}{\partial \theta}(\theta^*; h^* | H_t) \right\|_1 [h - h^*] \right] \\ &= \mathbb{E} \left[ \text{sign} \left( \frac{\partial \phi_k}{\partial \theta} \right) \cdot \underbrace{D_h \frac{\partial \phi_k}{\partial \theta} [h - h^*]}_{=0} \right] = 0 \end{aligned}$$

737 Hence  $\partial_\theta \mathcal{L}_{D,k}$  is Fréchet-orthogonal to all directions in the nuisance space at  $h = h^*$ , completing  
738 the proof.  $\square$

#### 739 B.4 The MSE rate theorem

740 A mean squared error (MSE) guarantee under finite samples for the estimator  $\hat{\psi}_t$  in a primary  
741 heterogeneous setting, where the conditioning variable is the static feature  $X_0$  (see theorem 10  
742 from [18]). We point out that the conclusion of the theorem could be extended to our setting, where  
743 conditioning variable is the history  $H_t$ . Before we show the theorem, we clarify the prerequisite  
744 notations and regular assumptions:

##### 745 Mathematical notations:

746 1. **Norm of a function in function space** Define the norm for a vector-valued function  
747  $\psi : \mathcal{X} \rightarrow \mathbb{R}^d$ , then the  $u, v$ -norm is defined as:

$$\|\psi\|_{u,v} = \mathbb{E}_{X \sim P_X} \left[ \left( \sum_{i=1}^d \psi_i(X)^u \right)^{\frac{v}{u}} \right]^{\frac{1}{v}}. \quad (53)$$

748 We point out that when  $u = v$ , the  $u, v$ -norm degenerate into the normal norm:  $\|\cdot\|_{u,u} =$   
749  $\|\cdot\|_u$ . Further, we could define the product norm between two functions  $f, g$  as:

$$\|f \circ g\|_{u,v} = \mathbb{E} \left[ \|f(\hat{X})\|_u^v \cdot \|g(\hat{Y})\|_u^v \right]^{1/v} = \mathbb{E} \left[ \left( \sum_{i=1}^r f_i(X)^u \right)^{\frac{v}{u}} \cdot \left( \sum_{i=1}^r g_i(X)^u \right)^{\frac{v}{u}} \right]^{\frac{1}{v}} \quad (54)$$

750 2. **Localized Rademacher complexity:** The localized Rademacher complexity is an extended  
751 concept of the standard Radmacher complexity. It measures the the fitting capacity of a  
752 function space  $\mathcal{F}$  where functions have bounded second moment. Formally it is defined as:

$$R_{P_X, n}(\mathcal{F}; \delta) = \mathbb{E}_{\epsilon_{1:n}, X_{1:n} \sim P_X} \left[ \sup_{f \in \mathcal{F}: \|f\|_2 \leq \delta} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(X_i) \right] \quad (55)$$

753 **Assumptions:** We adapt the assumptions from the original heterogeneous setting [18] to conditioning  
754 on the history variable  $H_t$ :

- 755 1. All the variables and functions in our problem setting are bounded.
- 756 2. A regularity assumption regarding the treatment variables: The following quantities are also  
757 bounded:

$$c_{k,j} := \sup_{t \leq T-\tau} \sup_{h_t \in \mathcal{H}_t} \mathbb{E} \left[ \|C_{t,k,j}\|_{u,u}^v \mid H_t = h_t \right]^{2/v} < +\infty \quad (56)$$

758 where  $C_{t,k,j} := \text{Cov}(A_{t+k}, A_{t+j} \mid H_{t+k})$ , and:

$$\forall h_t \in \mathcal{H}_t, \quad \mathbb{E}[C_{t,k,k} \mid H_t = h_t] \succeq \lambda I \quad (57)$$

759 3. The  $u, \infty$ -norms of  $\psi_{t,k}$  are bounded by:

$$M := \max_{t \leq T-\tau} \max_{0 \leq k \leq \tau, \psi_{t,k} \in \Psi_{t,k}} \|\psi_{t,k}\|_{u,\infty} < +\infty \quad (58)$$

760 **Theorem 1 (Adapted from Heterogeneous R-learner[18]):** Suppose the assumptions above are  
 761 satisfied. Let  $\delta_n$  be an upper bound on the critical radius of the star hull of all the function spaces

$$\{\Psi_{t,k,i} - \psi_{t,k,i}\}_{t \in [SL-\tau], k \in [\tau], i \in [r]}$$

762 and that  $\delta_n = \Omega\left(\sqrt{\frac{r \log \log(n)}{n}}\right)$ . Suppose:

$$\forall 1 \leq t \leq T - \tau, \quad \max_{0 \leq k \leq j \leq \tau} \mathbb{E}_{\hat{q}_{t,k,k}, \hat{q}_{t,j,k}} \left[ \|(\hat{q}_{t,k,k} - q_{t,k,k}) \circ (\hat{q}_{t,j,k} - q_{t,j,k})\|_{2,2}^2 \right] = O(r^2 \delta_{n/2}^2) \quad (59)$$

$$\forall 1 \leq t \leq T - \tau, \quad \max_{0 \leq k \leq \tau} \mathbb{E}_{\hat{q}_{t,k,k}, \hat{q}_{t,k}} \left[ \|(\hat{q}_{t,k,k} - q_{t,k,k}) \circ (\hat{p}_{t,k} - p_{t,k})\|_{2,2}^2 \right] = O(r^2 \delta_{n/2}^2), \quad (60)$$

763 then the blip coefficient predictor  $\hat{\psi}_{t,k}$  trained under the  $L^1$ -moment loss satisfies:

$$\max_{t \leq T-\tau} \max_{k \in \{0, \dots, \tau\}} \mathbb{E} \left[ \|\hat{\psi}_{t,k} - \psi_{t,k}\|_{2,2}^2 \right] = O(r^2 \delta_n^2), \quad \delta_n^2 \propto \frac{\log \log(n)}{n}. \quad (61)$$

764 The theorem uses the notion of critical radius which is often applied to describe the statistical  
 765 complexity of functional spaces[40]. Intuitively, the theorem states that when:

- 766 • our sequential neural architecture is correctly specified and trained with empirical risk  
 767 minimization,
- 768 • the product norm between nuisance networks has a error small enough,
- 769 • the boundness assumption stated above are satisfied,

770 then our trained neural network for predicting the sequential blip effect will only have an  $\Omega(\frac{\log \log(n)}{n})$   
 771 root mean squared error rate.



## C Additional Results

### C.1 Ablation studies

We report the ablation studies on (1) the neural backbone of the sequential encoder and (2) the  $L^1$ -moment loss with double optimization trick. To this, we further introduce the following baselines:

1. **DeepBlip-LSTM:** We replace the transformer architecture in the sequential encoders  $\mathcal{E}_{\theta_N}^N, \mathcal{E}_{\theta_B}^B$  of nuisance network and blip prediction network by the LSTM network. From the experimental results below, we find that DeepBlip-LSTM is still highly effective: it outperforms other baselines, which proves the effectiveness of our propose DeepBlip framework regardless of the instantiations. From the examples of DeepBlip instantiated by transformer and LSTM, we demonstrate that our DeepBlip framework can seamlessly integrate popular sequential networks.
2. **DeepBlip-WDO:** We remove the double optimization trick (WDO stands for **Without Double Optimization**) and only make a single forward pass on  $\hat{\psi}_t$  to construct the  $L^1$ -moment loss. In this scenario, we discard the mandate on the order of solving the blip coefficients estimators, which is required for a correct estimation in SNMMs. This way, we are able to identify the contribution of applying the double optimization trick. Our results show that DeepBlip-WDO suffers from high estimation error and thus, demonstrate the necessity of the double optimization trick.
3. **DeepBlip- $L^2$ :** We use the squared loss adopted by heterogeneous dynamic R-learner [18] instead of our  $L^1$ -moment loss and then optimize it with the double optimization trick:

$$\mathcal{L}_2^k = \frac{1}{T - \tau} \sum_{t=1}^{T-\tau} \sum_{i=1}^n (\tilde{Y}_{t+k}^i - \sum_{j=k+1}^{\tau} \hat{\psi}_j^2(H_t^i)' \tilde{A}_{t,j,k}^i - \hat{\psi}_k^1(H_t^i)' \tilde{A}_{t,k,k}^i)^2 \quad (62)$$

This design enables us to compare the performance between two losses. The  $L^2$  squared loss also works for our setting. However, as shown in the results, it has slower convergence and slightly higher bias. Hence, we highlight the advantage of using our  $L^1$ -moment loss which ensures stable optimization.

We report the RMSE on both synthetic dataset and semi-synthetic dataset as in the experiment section. We adopt the similar test setting in experiment section, where we test the performance against growing time-varying confounding  $\gamma_{\text{conf}}$  in the synthetic data experiment and test the performance against increasing prediction horizons  $\tau$  in the semi-synthetic dataset.

### C.2 Granular analysis of the blip coefficients prediction

So far, we only demonstrate the overall RMSE on the CATE over time estimation. Our DeepBlip works by predicting the blip coefficients  $\hat{\psi}_t(H_t) = (\hat{\psi}_{t,0}(H_t), \dots, \hat{\psi}_{t,\tau}(H_t)) \in \mathbb{R}^{(\tau+1)d_a}$ . In this subsection, we report a granular analysis on the prediction of the blip coefficients, bringing more transparency to our DeepBlip.

In fact, we could directly derive the ground truth personalized blip coefficients  $\psi_t = (\psi_{t,0}(H_t), \dots, \psi_{t,\tau}(H_t))$  for both the synthetic dataset in D.1 and semi-synthetic dataset in D.2. This enables us to accurately evaluate the performances of all the blip coefficient predictors. Specifically, we visualize the distributions of all the components of the difference between the true blip coefficients and the prediction  $\psi_t - \hat{\psi}_t \in \mathbb{R}^{(\tau+1)d_a}$ . For synthetic dataset, we set  $\gamma = 5$  and  $\tau = 2$ . For semi-synthetic dataset, we set  $\gamma = 1$  and  $\tau = 4$ . In both settings, we use DeepBlip models instantiated by transformer that achieve the lowest RMSE and DeepBlip models trained without the double optimization trick to further demonstrate its necessity. The results from Figure 8 and Figure 9 show that our DeepBlip model is indeed capable of predicting each blip coefficient with high accuracy while with low variance. This supports our claim that **DeepBlip achieves robust estimation of CATE by decomposing the total effect into incremental blip effects with controlled error**. Further, by comparing the blip prediction of DeepBlip with and without double optimization, we see that

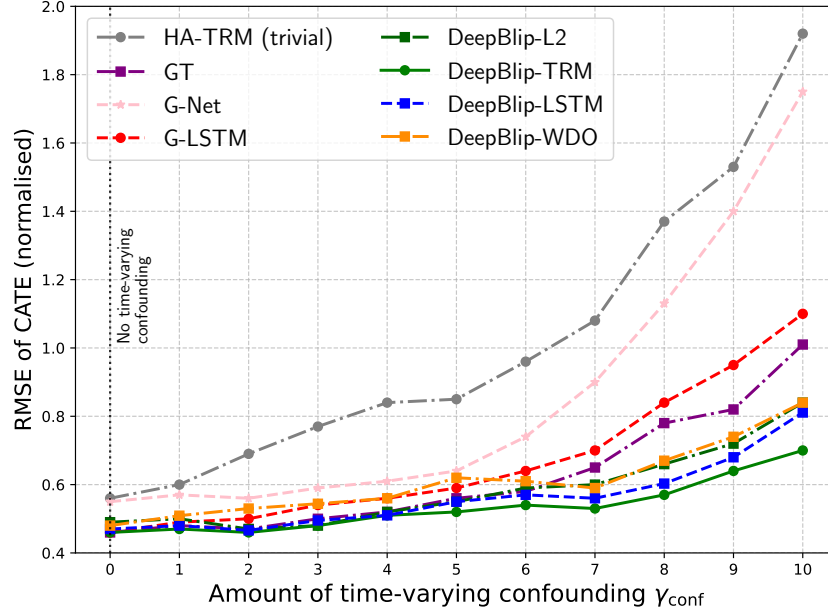


Figure 6: Ablation study on the tumor growth synthetic dataset. We explore two ablations: (1) We replace the transformer sequential encoder of DeepBlip and GT by LSTMs to create **DeepBlip-LSTM** and **G-LSTM**. (2) We remove the double optimization trick to create **DeepBlip-WDO** or replace  $L^1$ -moment loss with the  $L^2$  loss to create **DeepBlip-L<sup>2</sup>**. As in previous experiment, we report the average RMSE against increasing confounding  $\gamma_{conf}$ . We note that DeepBlip-LSTM is still more competitive than other baselines, which indicates the *effectiveness of the model-agnostic framework of DeepBlip*. DeepBlip-TRM, as our originally proposed model, remain the strongest.

## 818 D Dataset

819 In this section we will give a more detailed elaboration of the data generating process as well as the  
 820 key simulation of the counterfactual outcome (namely the potential outcome). The treatment effect is  
 821 then achieved by subtracting two counterfactual outcomes.

### 822 D.1 Tumor growth dataset

823 We specify the dynamic of the tumor growth model to be:

$$Y_{t+1} = Y_t + \underbrace{\left( \rho \log\left(\frac{K}{m(\bar{Y}_{t-l})}\right) + \epsilon_{t+1} \right) m(\bar{Y}_{t-l})}_{\text{Tumor Growth}} - \underbrace{m(\bar{Y}_{t-l}) \alpha_c c_{t+1}}_{\text{Chemotherapy}} + \underbrace{m(\bar{Y}_{t-l}) (\alpha_r d_{t+1} + \beta_r d_{t+1}^2)}_{\text{Radiotherapy}} \quad (63)$$

824  $\alpha_c, \alpha_r, \beta_r$  are the coefficients on the treatment effect of chemotherapy and radiotherapy,  $\rho, K$  together  
 825 control the growth model.  $m(\bar{Y}_{t-l})$  represents the mean of tumor volumes  $Y_1, \dots, Y_{t-l}$ , where  $l$  is  
 826 a lag parameter. The averaging of the volumes avoids abrupt changes for the tumor volume. The  
 827 chemotherapy drug dosage  $c_t$  and the radiation dosage  $d_t$  are applied with probability:

$$A_t^c, A_t^d \sim \text{Ber}\left(\sigma \cdot \left(\frac{\gamma_{conf}}{D_{max}} (\bar{D}_{15}(\bar{Y}_t) - \frac{D_{max}}{2})\right)\right) \quad (64)$$

828 Since the probability distribution is determined on the mean of previous 15 steps' tumor volumes  
 829  $\bar{D}_{15}(\bar{Y}_t)$ , time-varying confounding exists through future outcomes with its influence controlled by  
 830 the strength parameter  $\gamma_{conf}$ .

831 To generate the observational dataset, we simulate the trajectories of  $N = 1000$  patients with maximal  
 832 sequence length  $T = 30$  under confounding strength  $\gamma_{conf} \in \{0, 1, \dots, 10\}$ , resulting in 11 datasets

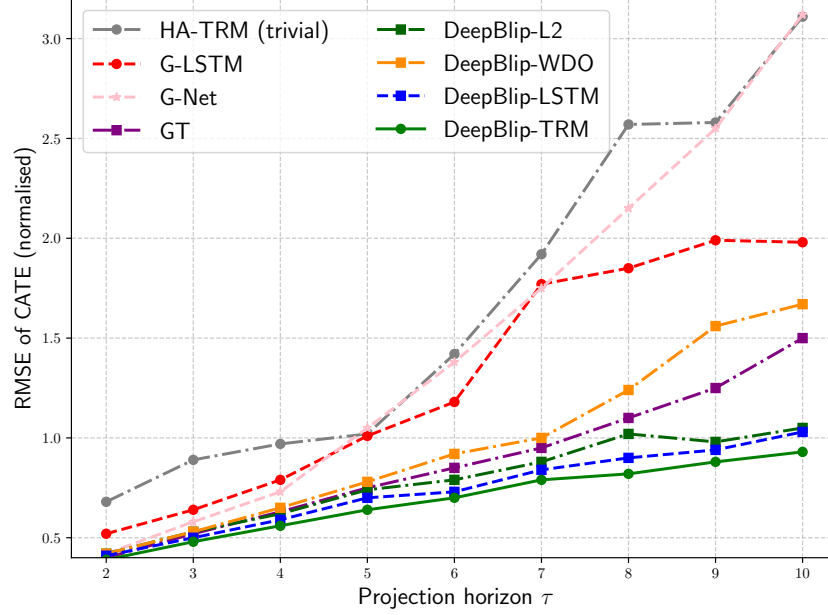


Figure 7: Ablation study on the MIMIC semi-synthetic dataset. Same ablations are investigated using the average RMSE against increasing projection horizon  $\tau$ . Even though DeepBlip-LSTM slightly underperforms the original model DeepBlip-TRM, it is still superior over other baselines. Hence, our proposed DeepBlip framework is *robust* over longer prediction horizon regardless of the backbone. We also observe that DeepBlip-WDO fail to predict CATE as accurately as the proposed DeepBlip-TRM. We also find that DeepBlip- $L^2$  has higher RMSE compared to  $L^1$ . This shows that the adoption of an stable  $L^1$ -moment loss, together with the double optimization trick, are crucial for learning the blip coefficients unbiasedly.

833  $\mathcal{D}_{N,T,\gamma_{conf}}$ . Then we split the simulated dataset into training, validation and test set by a split of  
834 (70%, 15%, 15%).

835 **Simulating counterfactuals:** For the test dataset, we compute the ground truth CATE by simulating  
836 the counterfactual outcomes  $Y_{t+\tau}^{a^*}$  and  $Y_{t+\tau}^{b^*}$  via a sliding window treatment along the sequence. To  
837 simulate the counterfactual outcome  $Y_{t+\tau}^{(a^*)}$  under an intervention  $a^* \in \mathbb{R}^{(\tau+1) \times 2}$  during  $[t, t + \tau]$   
838 for the tumor growth model, we iteratively compute the counterfactual trajectory starting from the  
839 observed history  $H_t$ . The DGP of outcome is consisted by the lagged dependency on tumor  
840 volume averages and the cumulative effects of treatments. Here we provide a strict formulation of the  
841 counterfactual outcome generation process, including the mathematical formulation and the steps for  
842 simulation.

843 Given the observed history  $H_t = (\bar{X}_t, \bar{A}_{t-1}, \bar{Y}_{t-1})$ , the counterfactual trajectory  $\{Y_s^{(a^*)}\}_{s=t}^{t+\tau}$  under  
844 intervention  $a^* = [a_t^*, a_{t+1}^*, \dots, a_{t+\tau}^*]$  (where  $a_s^* = (c_s^*, d_s^*)$ ) is generated with the following steps:

- 845 1. **Initialization:** Get observed history:  $Y_s$  for  $s < t$ . Note  $Y_{t-1}^{(a^*)} = Y_{t-1}$ .
- 846 2. **Iterative simulation:** For each time  $s \in [t, t + \tau - 1]$ , we simulate  $Y_{s+1}^{(a^*)}$  based on previous  
847 simulated outcomes. Notice that since  $l > \tau$ ,  $s - l < t$ , then  $m(\bar{Y}_{s-l})$  is unaffected by the  
848 interventions and therefore could be treated as constant when conditioning on  $H_t$ . Therefore

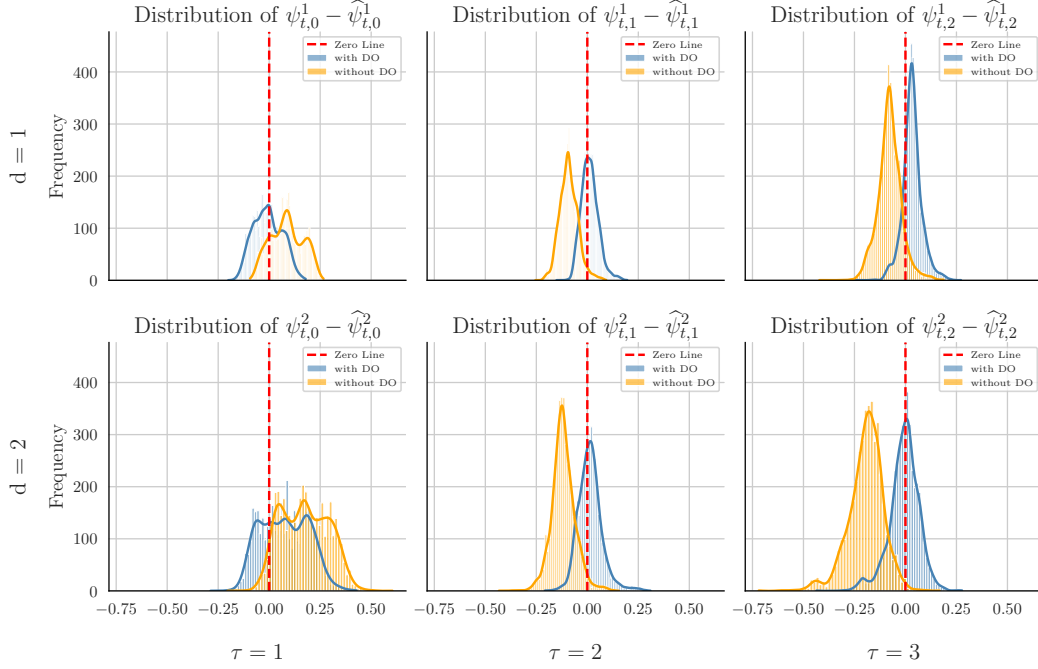


Figure 8: Tumor synthetic dataset ( $\gamma_{\text{conf}} = 5$ ,  $\tau = 2$ ): We visualize the distribution of the difference between the ground truth blip coefficients and the predictions made by our blip prediction network. Observing the histogram marked in blue represented by our DeepBlip, we find that the blip prediction network unbiasedly predict the blip coefficients, which offers proper adjustment for the confounding.

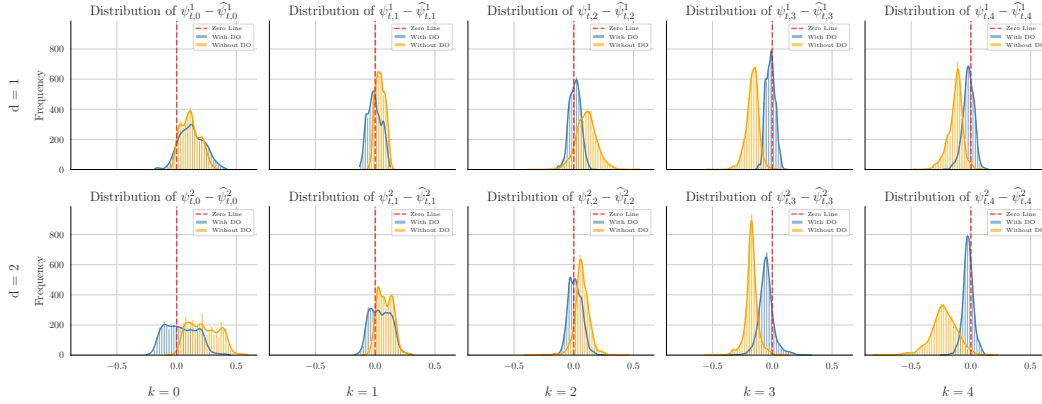


Figure 9: MIMIC semi-synthetic dataset ( $\gamma_{\text{conf}} = 1$ ,  $\tau = 4$ ): We conduct the same visualization. When  $k$  grows, the prediction gets more accurate. We expect this to happen, since  $\hat{\psi}_{t,\tau}$  is the first to optimize by the double optimization trick and then  $\hat{\psi}_{t,\tau-1}$  and so on. In contrast, the blip prediction network trained without double optimization exhibits significant bias and generally higher variance. This implies that the double optimization trick is crucial for the training DeepBlip.

849

we have:

$$\begin{aligned}
Y_{s+1}^{(\mathbf{a}^*)} = & Y_s^{(\mathbf{a}^*)} + \underbrace{\rho \log \left( \frac{K}{m(\bar{Y}_{s-l})} \right) m(\bar{Y}_{s-l})}_{\text{Growth Term}} + \epsilon_{s+1} \\
& - \underbrace{m(\bar{Y}_{s-l}) \alpha_c c_{s+1}^*}_{\text{Chemotherapy Effect}} - \underbrace{m(\bar{Y}_{s-l}) (\alpha_r d_{s+1}^* + \beta_r (d_{s+1}^*)^2)}_{\text{Radiotherapy Effect}}, \quad (65)
\end{aligned}$$

850

where  $\epsilon_{s+1}$  is the **original** noise term during the generation of the observational data.

851

 $c_{s+1}^*, d_{s+1}^*$  are the intervened chemotherapy and radiotherapy dosages at  $s+1$ , drawn from  $\mathbf{a}^*$ 

852

853 Under this iterative scheme, we could compute the final counterfactual outcome  $Y_{t+\tau}^{(\mathbf{a}^*)}$ .

854

**Treatment effect:** We then proceed to compute the individual conditioned treatment effect for a patient under treatment  $\mathbf{a}^*$  and the null treatment  $\mathbf{0}$ . Observing the process in (65), we find that the growth term **cancels off** each other during the subtraction, leaving only the terms for chemotherapy effect and radiotherapy effect:

855

856

857

$$Y_{t+\tau}^{(\mathbf{a}^*)} - Y_{t+\tau}^{(\mathbf{0})} = \sum_{s=t}^{t+\tau} \{ m(\bar{Y}_{s-l}) \alpha_c c_{s+1}^* + m(\bar{Y}_{s-l}) (\alpha_r d_{s+1}^* + \beta_r (d_{s+1}^*)^2) \} \quad (66)$$

858

The individual treatment effect is calculated for all the units on the test dataset  $D_{test}$  in a sliding treatment window rolling over  $t \in \{1, \dots, T - \tau\}$ , producing  $N_{test} \times (T - \tau)$  targets of individual treatment effect. We use these individual treatment effects as realizations of the CATE target

860

861

$$\mathbb{E}[Y_{t+\tau}^{(\mathbf{a}^*)} - Y_{t+\tau}^{(\mathbf{0})} | H_t].$$

862

## D.2 MIMIC-III semi-synthetic dataset

863

**MIMIC Data and Semi-Synthetic Benchmark** We use the **MIMIC-Extract** pipeline [41] to preprocess the **MIMIC-III** ICU dataset [15], aggregating hourly clinical measurements. Missing values are addressed via forward and backward filling, and continuous time-varying covariates are normalised. Our semi-synthetic dataset extends [33] to achieve ground-truth treatment effects.

867

**Cohort & Features:** A cohort of 1,000 patients is selected, restricted to ICU stays of 50–100 hours ( $T^{(i)} \in [50, 100]$ ). The feature space includes 18 time-varying vital signs and 3 static features (gender, ethnicity, age), where categorical variables are one-hot encoded. These variables, dynamic and static, are concatenated together to form a covariate vector  $X \in \mathcal{X}$  with  $d_x = \dim(\mathcal{X}) = 18 + 2 + 3 + 1 = 25$ .

871

**Untreated Outcomes:** For each patient  $i$ , the untreated trajectory  $\mathbf{Z}_t^{(i)}$  are simulated as:

$$Z_t^{(i)} = \underbrace{\alpha_s \text{B-spline}(t) + \alpha_g g^{(i)}(t)}_{\text{Endogenous}} + \underbrace{\alpha_f f_Z(X_t^{(i)})}_{\text{Exogenous}} + \underbrace{\epsilon_t}_{\text{Noise}}, \quad \epsilon_t \sim \mathcal{N}(0, 0.005^2), \quad (67)$$

872

where  $\alpha_s^j, \alpha_g^j, \alpha_f^j \in [0, 1]$  control component weights. The **endogenous** term combines: (1) a global trend  $\text{B-spline}(t)$  sampled from three cubic splines; (2) patient-specific local variations  $g^{j,(i)}(t)$  from a Gaussian process with Matérn kernel (length-scale  $\ell = 2.0$ ). The **exogenous** term is modeled by the effect  $f_Z(X_t^{(i)})$  via random Fourier features (RFF), which approximates Gaussian process. By combining the three terms we are able to simulate a complex untreated outcome that has a intrinsic trend modeled by (1) endogenous dependencies on a different scale (B-spline for the global trend and RFF for the local perturbations), and (2) an exogenous dependency on the time-varying covariates  $X_t$ . The untreated outcome model presents a challenging task which requires complex deep neural networks.

881

**Treatment Assignment:** 2 binary treatments  $\mathbf{A}_t^l, l \in \{1, 2\}$  are assigned sequentially using:

$$p_t^l = \sigma(\gamma_A^l m(Y_{t-w:t-1}) + \gamma_X^l f_X^l(X_t) + b^l), \quad A_t^l \sim \text{Bernoulli}(p_t^l), \quad (68)$$

where  $(\gamma_A^1 = 2.5, \gamma_X^1 = 1.0, \gamma_A^2 = 2.0, \gamma_X^2 = 0.25)$  control confounding strength,  $b = (-3.5, -1.75)$  biases treatment probability, and  $f_X^l(\cdot)$  uses RFF approximations of Gaussian process similar to  $f_Y$ . The term  $m(Y_{t-w:t-1})$  computes the average of prior outcomes over a window  $[t-w, t-1]$ .

**Treatment Effects:** Treatments affect outcomes via cumulative effects that decay over time:

$$E(t) = \sum_{i=t-w}^t \sum_{l=1}^2 \frac{\beta_l A_i \kappa^l(X_i)}{\sqrt{t-i+1}}, \quad \text{where } \kappa^l(X_i) = \tanh((\omega^l)' X_i) + 1 \quad (69)$$

where  $w = 5$  defines the effect window, and the inverse square root decay ensures maximal effect  $\beta_l$  at application time  $t$ . Effects are summed across all treatments to model the local treatment effect at time  $t$ . Sparse dependencies are enforced: Outcomes are affected by at most 3 covariates. The function  $\kappa^l$  further provides heterogeneity to the treatment effect by scaling the full effect  $\beta$ , parametrized by the coefficients  $\omega^l \in \mathbb{R}^{d_x}$  (also sparse). These arrangements **increase** the challenge to estimate treatment effect compared to previous benchmarks [20, 14].

**Observed outcomes:** Finally, we combine the untreated outcome and the treatment effects together to form the observed outcome trajectories:

$$Y_t^i = Z_t^i + E(t)^i \quad (70)$$

**Simulating counterfactuals:** To generate the counterfactual outcome under an intervention  $\mathbf{a}^* \in \mathbb{R}^{(\tau+1) \times 2}$  during  $[t, t + \tau]$ , we intervene on the treatment effect  $E(t)$  by replacing the observed treatments  $A_t$  by the intervened treatments  $\mathbf{a}^*$  while keeping the untreated trajectory  $Z_t$  and covariates  $X_t$  unchanged (note that the covariates  $X_t$  are predefined vital signs which remain unaffected during the DGP). The counterfactual outcome  $Y_{t+\tau}^{\text{CF}}$  at time  $s \in [t, t + \tau]$  is:

$$Y_{t+\tau}^{(\mathbf{a}^*)} | H_t = Z_{t+\tau} + \underbrace{\sum_{i=t+\tau-w}^s \sum_{l=1}^2 \frac{\beta_l \cdot a_i^{*,l} \cdot \kappa^l(X_i)}{\sqrt{t+\tau-i+1}}}_{E^*(t+\tau)},$$

where:  $Z_{t+\tau}$  is the untreated outcome (identical to the observed scenario). The treatments are divided by intervention period and the observed period (history)

$$A_i^{*,l} = \begin{cases} a_i^{*,l} & \text{if } t \leq i \leq t + \tau \text{ (intervention period),} \\ A_i^l & \text{otherwise (observed treatment).} \end{cases}$$

And  $\kappa^l(X_i) = \tanh((\omega^l)^\top X_i) + 1$  encodes effect heterogeneity via covariates. The simulation is performed for the units only on the **test** set in a sliding window treatment pattern. This means we start by conditioning on  $H_1$  to simulate counterfactual outcome at time  $\tau + 1$  and end at  $t = T - \tau$  to simulate counterfactual outcome at time  $T$ . This results in  $N_{\text{test}} \times (T - \tau)$  counterfactual outcomes under intervention  $\mathbf{a}^*$ .

The full implementation is available in our GitHub repository, including parameter configurations for reproducibility.

## E Baseline methods

We select six baselines to demonstrate the performance of our DeepBlip for CATE estimation over time. They are (1) History adjusted plug-in learner with LSTM instantiation (**HA-TRM**) [11], (2) recurrent marginal structural networks (**R-MSNs**) [19], (3) counterfactual recurrent network (**CRN**), [3] (4) **G-Net** [32], (5) G-transformer (**GT**) [14], and (6) causal transformer (**CT**) [20]. We will provide details for each baseline and briefly introduce the theory framework of MSM and G-computation. The hyperparameter tuning is detailed in Appendix G.5 and Appendix G.6.

### E.1 History Adjusted with transformer (HA-TRM)

A naïve approach for estimating CATE over time is to create a regressor with the treatment as the condition:

$$\hat{\delta}_\theta(a_{t:t+\tau}, h_t) \approx \mathbb{E}[Y_{t+\tau} \mid A_{t:t+\tau} = a_{t:t+\tau}, H_t = h_t]$$

where  $\hat{\delta}_\theta$  is a non-parametric model (like neural networks). Then we estimate the CATE of  $\mathbf{a}^*$  and  $\mathbf{b}^*$  as:

$$\hat{\delta}_\theta(\mathbf{a}^*, h_t) - \hat{\delta}_\theta(\mathbf{b}^*, h_t) \approx \mathbb{E}[Y_{t+\tau} \mid A_{t:t+\tau} = \mathbf{a}^*, H_t = h_t] - \mathbb{E}[Y_{t+\tau} \mid A_{t:t+\tau} = \mathbf{b}^*, H_t = h_t]$$

This method is given the name **PI-HA learner** (plugged-in history adjusted meta-learner) from the work of [11]. We instantiate this method with the transformer architecture to encode the history into a latent representation (see Appendix G), which is then fed into a fully connected linear network to predict the outcome.

Since this approach does not adjust for time-varying confounding, it is biased [11]. While an advantage of the HA-TRM is low-variance, it is subject to the level of confounding within the dataset as we observed in the experiment.

### E.2 Marginal Structural Models (MSM) and its neural Extension

**MSM Framework:** Marginal Structural Models (MSMs) [26, 13] address time-varying confounding via inverse probability of treatment weighting (**IPTW**). IPTW re-weights the targets to create a pseudo population that approximates the randomized controlled trial. With projection horizon  $\tau$ , the stabilized weight (**SW**) at time  $t$  for each sample is defined as:

$$SW(t, \tau) = \prod_{n=t}^{t+\tau} \frac{f(\mathbf{A}_n | \bar{\mathbf{A}}_{n-1})}{f(\mathbf{A}_n | \bar{\mathbf{H}}_n)}, \quad (71)$$

where  $f(\mathbf{A}_n | \bar{\mathbf{A}}_{n-1})$  and  $f(\mathbf{A}_n | \bar{\mathbf{H}}_n)$  represent treatment probabilities conditioned on past treatments or full history. These treatments are assumed discrete, often binary. The probabilities are estimated via logistic regressions. In practice the weights are truncated at 1st and 99th percentiles to avoid numeric overflow and then renormalised [19]. However, this division of probabilities still creates significant instability due to the sequential multiplication of the propensities in Eq. 71.

**Recurrent Marginal Structural networks (R-MSNs):** R-MSNs [19] extend MSMs via LSTM networks. Four components constitute the whole R-MSNs: propensity treatment network, propensity history network, encoder, decoder. Each sub-network uses LSTM as the sequential modeling architecture and a fully connected network to produce the desired prediction. Here is a detailed description of each model:

- *Propensity Network (conditioned on past treatment)* estimates numerator of  $SW(t, \tau)$ .
- *Propensity Network (conditioned on history)* estimates denominator of  $SW(t, \tau)$ .
- *Encoder* maps history  $\bar{\mathbf{H}}_t$  to latent representation
- *Decoder* predicts  $\tau$ -step outcomes using the representation generated by the encoder from last step.

In our experiment, we replace the LSTM backbone of the encoder and decoder by transformers to ensure fairness. Training occurs in 3 phases: (1) First the two propensity networks are trained to predict the binary treatment probabilities using binary cross entropy losses. Then the stabilized

weights can be computed for subsequent training of encoder and decoder. (2) Encoder is trained to predict the 1-step ahead outcome with the computed  $SW(\cdot, 1)$ -weighted MSE. (3) The decoder takes the latent representation from the encoder and then transform the vector into the latent representation fit for the decoder via a memory adapter (normally a fully connected linear layer). The decoder then predict the  $\tau$ -step outcome with the  $SW(\cdot, \tau)$ -weighted MSE.

In each phase, the training loss is aggregated through averaging the local loss at each time step  $0 \leq t \leq T - \tau$  (identical to our DeepBlip). At each local time step  $t$ , a representation is built with the transformer for the tasks. For more details please refer to [19].

### E.3 Causal estimation via balanced representations

**Counterfactual recurrent network (CRN):** CRN [3] employs adversarial learning to create treatment-invariant representations. CRN heuristically addresses the time-varying confounding through creating temporal representations that are non-predictive of the treatment assignment. For this an LSTM encoder-decoder architecture with gradient reversal layers is built to enforce balanced representation:

- *Encoder:* LSTM produces hidden states  $\mathbf{h}_t$
- *Balanced Representation:*  $\Phi_t = \text{FC}(\mathbf{h}_t)$
- *Adversarial Loss:*

$$\mathcal{L} = \|\mathbf{Y}_{t+1} - G_Y(\Phi_t, \mathbf{A}_t)\|^2 - \lambda \sum_{j=1}^{d_a} \mathbf{1}_{[\mathbf{A}_t=a_j]} \log G_A(\Phi_t) \quad (72)$$

Here,  $G_Y$  predicts outcomes while  $G_A$  attempts (but fails due to gradient reversal) to predict treatments from  $\Phi_t$ . We fix  $\lambda = 1$  as in original implementations.

**Causal transformer (CT):** CT [20] is the first transformer-based approach that estimates potential outcomes over time. CT employs a multi-input transformer architecture specifically designed to create history representations  $\Phi_t$  under time-varying treatment effect estimation setting. On top of  $\Phi_t$ , two additional networks are built:

- $G_Y$ : The outcome prediction network, which predicts the 1-step-ahead outcome  $\mathbf{Y}_{t+1}$  using  $\Phi_t$  and the current treatment  $\mathbf{A}_t$ .
- $G_A$ : The treatment classifier network, which attempts to predict the current treatment  $\mathbf{A}_t$  from  $\Phi_t$ .

Similar to CRN, CT attempts to create representations that are **predictive of outcomes** but **non-predictive of treatment assignments**, thereby mitigating confounding bias. To this, CT proposes counterfactual domain loss (CDC). The loss consists of: (1) Factual Outcome Loss:  $\mathcal{L}_{G_Y} = \|\mathbf{Y}_{t+1} - G_Y(\Phi_t, \mathbf{A}_t)\|^2$  This mean squared error ensures that  $\Phi_t$  is useful for predicting  $\mathbf{Y}_{t+1}$ . (2): Confusion Loss:  $\mathcal{L}_{\text{conf}} = -\sum_{j=1}^{d_a} \frac{1}{d_a} \log G_A(\Phi_t)$  This cross-entropy loss encourages  $G_A$  to output a uniform distribution over treatments, making  $\Phi_t$  treatment-invariant.

The training of CT involves optimizing two adversarial objectives:

- Optimize representation parameters ( $\theta_R$ ) and  $G_Y$  parameters ( $\theta_Y$ ):  $(\hat{\theta}_Y, \hat{\theta}_R) = \arg \min_{\theta_Y, \theta_R} \mathcal{L}_{G_Y}(\theta_Y, \theta_R) + \alpha \mathcal{L}_{\text{conf}}(\hat{\theta}_A, \theta_R)$
- Optimize the  $G_A$  parameters ( $\theta_A$ ):  $\hat{\theta}_A = \arg \min_{\theta_A} \alpha \mathcal{L}_{G_A}(\theta_A, \hat{\theta}_R)$

During inference, CT autoregressively predicts outcomes  $\hat{Y}_{t+k}$  using the balanced representations  $\Phi_{t+k}$  based on the observed history  $H_t$ .

### E.4 G-Computation-based network (G-Net, G-LSTM)

Both methods are based on the G-computation formula [27, 30]:



$$\mathbb{E}[Y_{t+\tau}^{(a_{t:t+\tau})} | H_t = h_t] = \mathbb{E} \left\{ \mathbb{E} \left[ \cdots \mathbb{E} \left\{ \mathbb{E}[Y_{t+\tau} | H_{t+\tau}, A_{t:t+\tau} = a_{t:t+\tau}] \right. \right. \right. \\ \left. \left. \left. \left| H_{t+\tau-1}, A_{t:t+\tau-1} = a_{t:t+\tau-2} \right\} \right. \right. \right. \\ \left. \left. \left. \cdots \left| H_{t+1}, A_{t:t+1} = a_{t:t+1} \right| \right| H_t = \bar{h}_t, A_t = a_t \right\} \right\} \quad (73)$$

990 Due to the nested structure of G-computation in 73, the model misspecification error get easily  
991 propagated through the iterative steps, leading to potential bias and high variance.

992 **G-Net (Simulation-based)** G-Net[32] is the first neural network model that uses G-computation to  
993 estimate CAPO over time. It make the Monte-Carlo simulation based on the integration:

$$\int_{\mathbb{R}^{d_x \times (\tau)} \times \mathbb{R}^{d_y \times (\tau)}} \mathbb{E}[Y_{t+\tau} | H_{t+\tau} = h_{t+\tau-1}, A_{t:t+\tau} = a_{t:t+\tau}] \\ \times \prod_{\delta=1}^{\tau} p(x_{t+\delta}, y_{t+\delta} | H_t = h_t, x_{t+1:t+\delta-1}, y_{t:t+\delta-1}, a_{t:t+\delta-1}) \\ d(x_{t+1:t+\tau-1}, y_{t:t+\tau-1}) \quad (74)$$

994 G-Net predicts the potential outcome in two steps:

- 995 1. Estimate the conditional distribution  $p(x_{t+\delta}, y_{t+\delta} | H_t = h_t, x_{t+1:t+\delta-1}, y_{t:t+\delta-1}, a_{t:t+\delta-1})$
- 996 2. Compute the empirical sum of M (M=50) trajectories via Monte-Carlo simulation to estimate  
997 the integration in 74.

998 G-Net originally uses LSTM to encode the history into a hidden state, with an extra linear transfor-  
999 mation layer mapping the hidden state to the latent representation. We replace the LSTM architecture  
1000 by a transformer in our experiment to ensure fairness. The representation is fed into output com-  
1001 puting head (Fully connected networks with respective activation output) to predict the conditional  
1002 distribution. For fairness, we replace the LSTM encoder with the encoder of the transformer.

1003 **G-transformer (regression-based)** Unlike G-Net, which uses Monte-Carlo simulation, G-  
1004 transformer (GT) [14] is built via a iterative regression on the pseudo outcomes defined as:

- 1005 1.  $G_{t+\tau}^a = Y_{t+\tau}$
- 1006 2.  $G_{t+\delta}^a = \mathbb{E}[G_{t+\delta+1}^a | \bar{H}_{t+\delta}^t, A_{t:t+\delta} = a_{t:t+\delta}]$  for  $\delta = 0, \dots, \tau - 1$

1007 GT predicts the outcome in a masked transformer encoder. It uses transformer to encode the history  
1008 at  $t + \delta$  into a representation  $Z_{t+\delta} = z_\phi(h_{t+\delta})$ , and then predicts the pseudo outcome at step  $\delta$  as:

$$g_\phi^\delta(Z_{t+\delta}^a, A_{t+\delta}) \approx G_{t+\delta}^a$$

1009 And then the CAPO is predicted as  $g_\phi^0(Z_t, a_t)$ . The encoder  $z_\phi(\cdot)$  is instantiated by a transformer  
1010 architecture similar to CT (For more details see [14]).

## F Algorithm

---

**Algorithm 1** DeepBlip training algorithm
 

---

```

1: Input: Dataset  $\mathcal{D}$ , horizon  $\tau$ , feature map  $\phi$ , number of folds  $W$ , learning rates  $\eta_N, \eta_B$ 
2: Output: Blip prediction network  $(\mathcal{E}_{\theta_B}^B, \mathbf{gb}_{\theta_B}^k) \hat{\psi}_{\theta_B}$ 
3:
4: First Stage: Train Nuisance Network  $(\mathcal{E}_{\theta_N}^N, \mathbf{gp}_{\theta_N}^k, \mathbf{gq}_{\theta_N}^{j,k})$ 
5: Partition  $\mathcal{D}$  into  $W$  folds:  $S_1, S_2, \dots, S_W$ 
6: Partition  $\mathcal{D}$  into  $W$  folds:  $S_1, S_2, \dots, S_W$ 
7: for each fold  $w = 1$  to  $W$  do
8:   Train nuisance network on  $\mathcal{D} \setminus S_w$ 
9:   for each  $t = 1$  to  $T - \tau$ ,  $k = 0$  to  $\tau$  do
10:     $z_{t+k}^N \leftarrow \mathcal{E}_{\theta_N}^N(h_{t+k})$ 
11:     $\hat{p}_k(z_{t+k}^N) \leftarrow \mathbf{gp}_{\theta_N}^k(z_{t+k}^N)$ 
12:     $\mathcal{L}_{t,k}^p = (Y_{t+\tau} - \hat{p}_k(z_{t+k}^N))^2$ 
13:    for  $j = k$  to  $\tau$  do
14:       $\hat{q}_{j,k}(z_{t+k}^N) \leftarrow \mathbf{gq}_{\theta_N}^{j,k}(z_{t+k}^N)$ 
15:       $Q_{t,j} = \phi(X_{t+1:t+j}, A_{t:t+j}) - \phi(X_{t+1:t+j}, (A_{t:t+j-1}, 0))$ 
16:       $\mathcal{L}_{t,j,k}^q = BCE(\hat{q}_{j,k}(z_{t+k}^N), Q_{t,j})$  (if binary treatment)
17:    end for
18:  end for
19:   $\mathcal{L}_N \leftarrow \frac{1}{T-\tau} \sum_{t=1}^{T-\tau} \left( \frac{1}{\tau+1} \sum_{k=0}^{\tau} \mathcal{L}_{t,k}^p + \frac{2}{(\tau+1)(\tau+2)} \sum_{\tau \geq j \geq k \geq 0} \mathcal{L}_{t,j,k}^q \right)$ 
20:  {Compute the gradient update the nuisance network parameter  $\theta_N$ }
21:   $\theta_N \leftarrow \theta_N - \eta_N \nabla_{\theta_N} \mathcal{L}_N$ 
22:  for each sample in  $S_w$  do
23:    Compute the residuals:
24:     $\tilde{Y}_{t,k} = Y_{t+\tau} - \mathbf{gp}_{\theta_N}^k(z_{t+k}^N)$ 
25:     $\tilde{A}_{t,j,k} = A_j - \mathbf{gq}_{\theta_N}^{j,k}(z_{t+k}^N)$ 
26:  end for
27: end for
28:
29: Second Stage: Train Blip Prediction Network  $(\mathcal{E}_{\theta_B}^B, \mathbf{gb}_{\theta_B}^k)$ 
30: for each  $t = 1$  to  $T - \tau$  do
31:   Compute the hidden state of the encoder:  $z_t^B \leftarrow \mathcal{E}_{\theta_B}^B(h_t)$ 
32:   for  $k = 0$  to  $\tau$  do
33:      $\hat{\psi}_k^1(z_t^B) \leftarrow \mathbf{gb}_{\theta_B}^k(z_t^B)$ 
34:      $\hat{\psi}_k^2(z_t^B) \leftarrow \mathbf{gb}_{\theta_B}^k(z_t^B)$ 
35:      $\mathcal{L}_{\text{blip},t,k} = \left\| \left( \tilde{Y}_{t,k} - \sum_{j=k+1}^{\tau} \hat{\psi}_j^2(z_t^B)' \tilde{A}_{t,j,k} - \hat{\psi}_k^1(z_t^B)' \tilde{A}_{t,k,k} \right) \tilde{A}_{t,k,k} \right\|_1$ 
36:   end for
37: end for
38:  $\mathcal{L}_{\text{blip}} \leftarrow \frac{1}{T-\tau} \sum_{t=1}^{T-\tau} \left( \frac{1}{\tau+1} \sum_{k=0}^{\tau} \mathcal{L}_{\text{blip},t,k} \right)$ 
39: Compute the gradient update the blip prediction network parameter  $\theta_B$ 
40:  $\theta_B \leftarrow \theta_B - \eta_B \nabla_{\theta_B} \mathcal{L}_{\text{blip}}$ 
41:

```

---

*Note:* We use  $\leftarrow$  operation for value assignment in the computation graph while  $\leftarrow$  operations are detached from the computation graph.

---

---

**Algorithm 2** DeepBlip inference algorithm

---

- 1: **Input:** Patient history  $H_t = h_t$ , treatment sequence  $a^*$ , baseline treatment sequence  $b^*$
  - 2: **Output:** CATE estimate  $\mathbb{E}[Y_{t+\tau}^{(a^*)} - Y_{t+\tau}^{(b^*)} \mid H_t = h_t]$
  - 3: **Step 1: Compute the hidden state of the encoder**
  - 4:  $z_t^B \leftarrow \mathcal{E}_{\theta_B}^B(h_t)$
  - 5: **Step 2: Predict the blip coefficients**
  - 6:  $\hat{\psi}_k(z_t^B) \leftarrow \text{gb}_{\theta_B}^k(z_t^B)$
  - 7: **Step 3: Estimate CATE**
  - 8: (when  $\phi$  only depends the last treatment:  $\phi(x_{t+1:t+k}, a_{t:t+k}) = \phi(a_{t+k})$ )
  - 9:  $\text{CATE} = \sum_{k=0}^{\tau} \hat{\psi}_k(z_t^B)'(\phi(a_k^*) - \phi(b_k^*))$
- 

## G Implementation details

In subsection G.2, we introduce the runtime of our system. We provide a runtime comparison of the different methods in subsection G.2. Then in subsection G.4 and G.3, we introduce how to instantiate our DeepBlip with LSTM and transformers respectively. In subsection G.5 and subsection G.6, we report the details of hyperparameter tuning.

### G.1 Software & system

Our DeepBlip is implemented using **PyTorch** [23], combined with **PyTorch Lightning** [8] for streamlined code organization and model training. Our code is available at <https://anonymous.4open.science/r/DeepBlip-A39B>. Training and inference are performed on a single GeForce RTX 4080 with 18GB memory.

### G.2 Runtime comparison

In this section we report: (1) The runtime for training a single epoch on the training set, (2) The runtime for inference on the test set, and (3) The theoretical runtime required to identify the optimal treatment offline. The last metric is crucial for efficient offline treatment planning in practical medical scenarios. During the test, we choose transformer as the main backbone for all the selected methods. For two stage methods like RMSNs and DeepBlip, the first stage uses a simpler LSTM instantiation instead for higher efficiency. The results are reported in table 3.

Baselines like HA-TRM, RMSNs, CRN, CT estimate one conditional average potential outcome (CAPO) over time within a single inference. Therefore it takes  $N(\tau)$  inferences to identify the treatment sequence with best outcome, where  $N(\tau)$  denotes the number of possible treatment combinations within the time window  $[t, t + \tau]$ . G-Net computes the CAPO over time by performing Monte Carlo simulations, generating  $M$  trajectories for each possible treatment sequence. G-transformer only predicts the CAPO under a fixed intervention sequence, which means it requires additional  $N(\tau) - 1$  re-trainings to calculate all the possible outcomes. Our DeepBlip, however, estimates CATE over time via predicting blip coefficients to identify the treatment combination with the optimal effect. By Eq. 3, only one forward pass to predict the blip coefficients suffices to identify the optimal treatment. Hence, although our DeepBlip is not the most efficient method during training or inference, **it is significantly faster in offline treatment planning**, which is critical for personalized medicine.

### G.3 Transformer instantiation for DeepBlip

Our DeepBlip uses the transformer architecture to encode history into a latent representation. The transformer is built upon the multi-input transformer from causal transformer [20]. The transformer is specially designed for medical application where the inputs include the outcomes, covariates and treatments. Each type of variable has a corresponding sub-transformer. The three sub-transformers perform not only the classic multi-headed self-attention within themselves, but also the cross-attention in-between. This ensures that information is shared across these sub-transformers.

Table 3: Runtime report for the methods instantiated by transformers.  $t_T$  (in min) is the training time on the training dataset per epoch.  $t_I$  (in min) is the inference time on the test dataset. In the last column we report the time needed to identify the optimal treatment.  $N(\tau)$  is the number of all possible combination of treatments during the future horizons  $[t : t + \tau]$ .

Methods	Training ( $t_T$ )	Inference ( $t_I$ )	Offline optimal treatment identification
HA-TRM [11]	$2.5 \pm 0.4$	$0.16 \pm 0.04$	$t_I \cdot N(\tau)$
RMSNs [19]	$7.2 \pm 0.9$	$0.30 \pm 0.06$	$t_I \cdot N(\tau)$
CRN [3]	$5.4 \pm 0.6$	$0.28 \pm 0.06$	$t_I \cdot N(\tau)$
CT [20]	$3.4 \pm 0.5$	$0.20 \pm 0.04$	$t_I \cdot N(\tau)$
G-Net [32]	$2.8 \pm 0.4$	$0.24 \pm 0.03$	$t_I \cdot N(\tau) \cdot M$ <sup>†</sup>
GT [14]	$2.6 \pm 0.4$	$0.20 \pm 0.04$	$(N(\tau) - 1) \cdot t_T \cdot n_e + t_I \cdot N(\tau)$ <sup>††</sup>
<b>DeepBlip (ours)</b>	$3.6 \pm 0.5$	$0.17 \pm 0.06$	$t_I$

<sup>†</sup>: G-Net performs  $M$  simulations to infer the potential outcome.

<sup>††</sup>: G-transformer requires re-training additionally  $N(\tau) - 1$  times.

Let the three sub-transformers be  $z_\theta^k(\cdot)$ ,  $k = 1, 2, 3$ . Suppose for  $k = 1$ , the input sequence is the  $\mathbf{U}^k = (X_1, \dots, X_t)$ . Then the sub-transformer  $z_\phi^k(\cdot)$  generate the representation in the following steps.

**Input transformation:** The raw input first goes through a linear transformation:

$$H_0^k = \text{Linear}_k(\mathbf{U}) \in \mathbb{R}^{t \cdot d_{\text{um}}} \quad (75)$$

, where  $H_0^k$  is the input of the first transformer blocks.

**Transformer blocks:** The multi-input sub-transformer has  $B$  transformer blocks. Within each transformer block, the multi-headed self-attention and cross-attention are performed. Then the output goes through a feed forward network. Here are the details of a transformer block at layer  $b$ .

(1) Multi-headed self/cross-attentions First the keys, queries and values, namely  $K, Q, V$  for  $n_h$  attention heads, are computed from the output from last transformer block:

$$\begin{aligned} Q_b^{k,i} &= H_b^{k,i} W_Q^{k,i} + \mathbf{b}_Q^{k,i} \in \mathbb{R}^{t \cdot d_{qkv}}, \\ K_b^{k,i} &= H_b^{k,i} W_K^{k,i} + \mathbf{b}_W^{k,i} \in \mathbb{R}^{d_t \cdot qkv}, \\ V_b^{k,i} &= H_b^{k,i} W_V^{k,i} + \mathbf{b}_V^{k,i} \in \mathbb{R}^{d_t \cdot qkv}, \end{aligned} \quad (76)$$

where  $i \in [1, n_h]$  is the index of a single attention head. Then the  $i$ -th attention is computed via a softmax over the scaled dot-product:

$$\text{Attn}_b^{k,i} = \text{Softmax}\left(\frac{Q_b^{k,i} (K_b^{k,i})^T}{\sqrt{d_{qkv}}}\right) V_b^{k,i} \quad (77)$$

The multi-head self-attention is then calculated by concatenating the attention heads:

$$\text{MHA}(Q_b^k, K_b^k, V_b^k) = \text{Concat}(\text{Attn}_b^{k,1}, \dots, \text{Attn}_b^{k,n_h}) \in \mathbb{R}^{t \cdot d_{\text{um}}} \quad (78)$$

Since  $\text{MHA}(Q_b^k, K_b^k, V_b^k)$  has the same dimension as  $H_b^k$ , then we could generate a new set of queries for the cross-attentions with residual connection:

$$\tilde{Q}_b^k = H_b^k + \text{MHA}(Q_b^k, K_b^k, V_b^k) \quad (79)$$

The cross-attentions are put on top of the self-attention layers and uses the queries and keys from the other two sub-transformers. For example, the multi-head cross-attention of sub-transformer  $k$  with sub-transformer  $m$  is:

$$\text{MHA}(\tilde{Q}_b^k, K_b^m, V_b^m) = \text{Concat}(\text{CrossAttn}_b^{k,m,1}, \dots, \text{CrossAttn}_b^{k,m,n_h}), \quad (m \neq k), \quad (80)$$

where the CrossAttn is defined as:

$$\text{CrossAttn}_b^{k,m,i} = \text{Softmax}\left(\frac{\tilde{Q}_b^k (K_b^{m,i})^T}{\sqrt{d_{qkv}}}\right) V_b^{m,i}. \quad (81)$$

1067 Finally, the output is achieved by adding self-attention’s output and cross-attention’s output together:

$$Z_b^k = \tilde{Q}_b^k + \sum_{m \neq k} \text{MHA}(\tilde{Q}_b^k, K_b^m, V_b^m) + S. \quad (82)$$

1068 The static covariates  $S$  are added when pooling different cross-attention outputs (see section 4.1  
1069 from [20]). The operations ensures that the information is shared across sub-transformers. Layer  
1070 normalizations are added for all the self- and cross attentions [38].

1071 (2) Feed-forward networks The output of the attention mechanism is processed through a position-  
1072 wise feed-forward network (FFN) applied to each time step:

$$H_b^{k,\text{ff}} = \text{Linear}(\text{ReLU}(\text{Linear}(Z_b^k))), \quad (83)$$

1073 where linear layers are followed by a dropout [20]. The output  $H_b^{k,\text{ff}}$  serves as the input ( $H_{b+1}^k$ )  
1074 to the next transformer block. After processing through all  $B$  blocks, the final representation for  
1075 sub-transformer  $k$  becomes  $H_B^k$ . These final representations are then aggregated via a element-wise  
1076 average pooling followed by a linear transformation and a exponential linear unit (ELU):

$$H_B = \text{ELU}(\text{Linear}(\text{Pool}(H_B^1, H_B^2, H_B^3))) \in \mathbb{R}^{t \cdot d_{\text{hr}}}. \quad (84)$$

1077 Hence, the transformer-based sequential encoder  $\mathcal{E}_\theta$  outputs the latent vector:  $\mathcal{E}_\theta(H_t) = \mathbf{hr}_t =$   
1078  $H_{B,t} \in \mathbb{R}^{d_{\text{um}}}$  at time step  $t$ . We omit the positional encoding for better clarity. For more relevant  
1079 details see section 4.2 from [20].

1080 Our DeepBlip has two stages: (1) Nuisance network, and (2) Blip prediction network. Each stage  
1081 uses a separate sequential encoder:  $\mathcal{E}_{\theta_N}^N$  for (1) and  $\mathcal{E}_{\theta_B}^B$  for (2). To accelerate the training, we only  
1082 instantiate  $\mathcal{E}_{\theta_B}^B$  by the transformer, and use LSTM to instantiate (1). The underlying reason is that  
1083 the demand for accuracy in Stage (1) is not as high as Stage (2), as the  $L^1$ -moment loss is Neyman  
1084 orthogonal over the nuisance functions.

1085 **Nuisance Network:** At each time step  $t$ , the nuisance network Estimates the nuisance functions:

$$q_{t,j,k}(h_{t+k}) := \mathbb{E}[Y_{t+\tau} \mid H_{t+k} = h_{t+k}], \quad 1 \leq t \leq T - \tau, 0 \leq k \leq \tau \quad (85)$$

$$q_{t,j,k}(h_{t+k}) := \mathbb{E}[Q_{t,j} \mid H_{t+k} = h_{t+k}], \quad 1 \leq t \leq T - \tau, 0 \leq k \leq j \leq \tau \quad (86)$$

1086 Therefore a **multi-head output layer** is added on top of the encoder  $\mathcal{E}_\theta^N$ . The heads receive the  
1087 representation  $\mathbf{hr}_t^N$  from the encoder and transform the input with multi-layer perceptron networks:

$$1. \hat{p}_{t,k}(h_{t+k}) = \text{MLP}_p^{(k)}(\mathbf{hr}_{t+k}^N) \approx p_{t,k}(h_{t+k}), \quad k = 0, 1, \dots, \tau$$

$$2. \hat{q}_{t,j,k}(h_{t+k}) = \text{MLP}_q^{(j,k)}(\mathbf{hr}_{t+k}^N) \approx q_{t,j,k}(h_{t+k}), \quad 0 \leq k \leq j \leq \tau$$

1090 where  $\text{MLP}_p^{(k)}$  and  $\text{MLP}_q^{(j,k)}$  are fully connected networks with ReLU activations. Each MLP has  
1091 input size of  $d_{\text{hr}}$ , a single hidden layer with  $d_{\text{hidden}}$  perceptrons and an single output transformation  
1092 that maps the hidden layer to a 1-dimensional output. For binary output in some  $q_{t,j,k}$ , an additional  
1093 sigmoid activation is added to contract the range into  $[0, 1]$ .

1094 **Blip prediction network’s output** Likewise a multi-head output layer is added on top of the  
1095 sequential encoder  $\mathcal{E}_{\theta_B}^B$  to predict the blip coefficients. For each horizon  $k \in \{0, 1, \dots, \tau\}$ :

$$\hat{\psi}_k(h_t) = \text{MLP}_{\theta_B}^{(k)}(\mathbf{hr}_t^B) \approx \psi_{t,k}(h_t)$$

1096 where  $\text{MLP}_{\theta_B}^{(k)}$  maps  $\mathbf{hr}_t^B$  to blip coefficients.  $\text{MLP}_{\theta_B}^{(k)}$  has the same structure as the MLPs in stage (1).

#### 1097 G.4 LSTM instantiation for DeepBlip

1098 DeepBlip-LSTM uses LSTM as the sequential encoders in two stages: 1. Stage (1) (Nuisance  
1099 Network): Estimates residuals for blip function estimation. 2. Stage (2) (Blip Prediction Network):  
1100 Predicts blip parameters  $\psi_{t,k}(h_t)$ .

1101 **Sequential Encoding via LSTM** Each stage shares an individual LSTM to encode patient history.  
1102 Let  $\mathbf{hz}_t$  be the hidden state at time  $t$ . At each time step  $t$ , the input vector  $\mathbf{v}_t \in \mathbb{R}^{d_{\text{input}}}$  concatenates:

- 1103 1. Static features  $X_s \in \mathbb{R}^{d_{\text{static}}}$  (remain constant across time)
- 1104 2. Current time-varying covariates  $X_d \in \mathbb{R}^{d_{\text{dynamic}}}$  ( $d_{\text{dynamic}} + d_{\text{static}} = d_x$ )
- 1105 3. Previous treatments  $a_{t-1} \in \mathbb{R}^{d_a}$
- 1106 4. Previous outcomes  $Y_{t-1} \in \mathbb{R}$

1107 to form the input vector:

$$V_t = \text{Concat}(X_s, X_d, A_{t-1}, Y_{t-1}) \in \mathbb{R}^{d_{\text{input}}}, \quad d_{\text{input}} = d_x + d_a + 1.$$

1108 Then the LSTM processes  $\{V_t\}_{t=1}^T$  into hidden state  $\text{hs}_t$  and cell state  $c_t$ :

$$(\text{hs}_t, c_t) = \text{LSTM}(v_t, (\text{hs}_{t-1}, c_{t-1})),$$

1109 where  $\mathbf{h}_t \in \mathbb{R}^{d_{\text{hidden}}}$ . Next we take the hidden state and derive the temporal representation of the  
 1110 patient state at time  $t$ :

$$\mathbf{hr}_t = \text{ELU}(\mathbf{W}_{\text{hr}} \cdot \text{dropout}(\text{hs}_t) + b_{\text{hr}}),$$

1111 where  $\mathbf{W}_{\text{hr}} \in \mathbb{R}^{d_{\text{hr}} \times d_{\text{hidden}}}$  is a learnable projection matrix.

1112 Above is the whole process of generating the temporal representation at time  $t$ . Since the DeepBlip  
 1113 has two stages, therefore we need to train **two separate** sequential encoders:  $\mathcal{E}_{\theta_N}^N(h_t) = \mathbf{hr}_t^N$  for  
 1114 the nuisance network and  $\mathcal{E}_{\theta_B}^B(h_t) = \mathbf{hr}_t^B$  for the blip prediction network. The output layer of both  
 1115 stages remain the same as the DeepBlip-TRM.

Table 4: Hyperparameter Tuning for Methods on Tumor Growth Data: We perform a random grid search with 20 iterations and choose the best hyperparameters for each task.  $C = d_{\text{input}}$  is the dimension of the input. For the causal transformer, the CDC coefficient is  $\alpha = 0.01$  [20]. The number of MC samples of G-Net is 50 [19]. Models are either instantiated by LSTMs or transformers. Here we display the search ranges of the *original* instantiations for each method. In the experiments, we adopt a universal instantiation type for all the baselines for fairness. Hence, the performance comparison in experiment section is fair.

Method	Component	Hyperparameter	Tuning Range
HA-TRM [11]	(end-to-end)	Transformer blocks ( $B$ )	[1, 2]
		Learning rate ( $\eta$ )	[0.01, 0.001, 0.0001]
CRN [3]	Encoder	Batch size	[64, 128]
		Attention heads ( $n_h$ )	2
		Transformer units ( $d_{\text{trm}}$ )	[0.5C, 1C]
		hidden representation ( $d_{\text{tr}}$ )	[0.5C, 1C]
		FC hidden units ( $d_{\text{tr}}$ )	[0.5d <sub>tr</sub> , 1d <sub>tr</sub> ]
		Number of epochs ( $n_e$ )	10
	Decoder	LSTM layers ( $l$ )	[1, 2, 3]
		Learning rate ( $\eta$ )	[0.1, 0.01, 0.001]
		LSTM hidden units ( $d_h$ )	[32, 64, 128, 256]
		LSTM dropout rate ( $p$ )	[0.1, 0.2, 0.3]
		Balanced representation size ( $d_r$ )	[0.5d <sub>h</sub> , 1d <sub>h</sub> , 2d <sub>h</sub> ]
CT [20]	(end-to-end)	FC hidden units ( $d_{\text{tr}}$ )	[0.5d <sub>tr</sub> , 1d <sub>tr</sub> ]
		Number of epochs ( $n_e$ )	20
R-MSNs [19]	Treatment/History Propensity Network	Transformer blocks ( $B$ )	[1, 2]
		Learning rate ( $\eta$ )	[0.01, 0.001, 0.0001]
		Batch size	[64, 128]
		Attention heads ( $n_h$ )	2
		Transformer units ( $d_{\text{trm}}$ )	[0.5C, 1C]
	Encoder / Decoder	hidden representation ( $d_{\text{tr}}$ )	[0.5C, 1C]
		FC hidden units ( $d_{\text{tr}}$ )	[0.5d <sub>tr</sub> , 1d <sub>tr</sub> ]
		LSTM dropout rate ( $p$ )	[0.1, 0.2, 0.3]
		Number of epochs ( $n_e$ )	20
		Max gradient norm	[0.5, 1.0, 2.0]
G-Net [32]	(end-to-end)	Number of epochs ( $n_e$ )	10
		Number of epochs ( $n_e$ )	20
GT [14]	(end-to-end)	LSTM layers ( $l$ )	[1, 2, 3]
		Learning rate ( $\eta$ )	[0.1, 0.01, 0.001]
DeepBlip (ours)	Nuisance Network	LSTM hidden units ( $d_h$ )	[64, 128, 256]
		LSTM output size ( $d_o$ )	[0.5d <sub>h</sub> , 1d <sub>h</sub> , 2d <sub>h</sub> ]
		Feed-forward hidden units ( $d_{\text{tr}}$ )	[0.5d <sub>h</sub> , 1d <sub>h</sub> , 2d <sub>h</sub> ]
		LSTM dropout rate ( $p$ )	[0.1, 0.2]
		Number of epochs ( $n_e$ )	20
	Blip Prediction Network	Transformer blocks ( $B$ )	[1, 2]
		Learning rate ( $\eta$ )	[0.01, 0.001, 0.0001]
		Batch size	[64, 128]
		Attention heads ( $n_h$ )	2
		Transformer units ( $d_{\text{trm}}$ )	[0.5C, 1C]
		hidden representation ( $d_{\text{tr}}$ )	[0.5C, 1C]
		FC hidden units ( $d_{\text{tr}}$ )	[0.5d <sub>tr</sub> , 1d <sub>tr</sub> ]
		Number of epochs ( $n_e$ )	20

Table 5: Hyperparameter Tuning for Methods on MIMIC-III semi-synthetic Data: We perform a random grid search with 20 iterations and choose the best hyperparameters for each task.  $C = d_{\text{input}}$  is the dimension of the input. For the causal transformer, the CDC coefficient is  $\alpha = 0.01$  [20]. The number of MC samples of G-Net is 50 [19]. Models are either instantiated by LSTMs or transformers. Here we display the search ranges of the *original* instantiations for each method. In the experiments, we adopt a universal instantiation type for all the baselines for fairness. Hence, the performance comparison in experiment section is fair.

Method	Component	Hyperparameter	Tuning Range
HA-TRM [11]	(end-to-end)	Transformer blocks ( $B$ )	[1, 2]
		Learning rate ( $\eta$ )	[0.01, 0.001, 0.0001]
		Batch size	[64, 128]
		Attention heads ( $n_h$ )	2
		Transformer units ( $d_{\text{trm}}$ )	[0.5C, 1C]
		hidden representation ( $d_{\text{hr}}$ )	[0.5C, 1C]
CRN [3]	Encoder	FC hidden units ( $d_{\text{fc}}$ )	[0.5 $d_{\text{hr}}$ , 1 $d_{\text{hr}}$ ]
		Number of epochs ( $n_e$ )	10
		LSTM layers ( $l$ )	[1, 2, 3]
		Learning rate ( $\eta$ )	[0.1, 0.01, 0.001]
		LSTM hidden units ( $d_h$ )	[32, 64, 128, 256]
		LSTM dropout rate ( $p$ )	[0.1, 0.2, 0.3]
	Decoder	Balanced representation size ( $d_r$ )	[0.5 $d_h$ , 1 $d_h$ , 2 $d_h$ ]
		FC hidden units ( $d_{\text{fc}}$ )	[0.5 $d_h$ , 1 $d_h$ , 2 $d_h$ ]
		LSTM dropout rate ( $p$ )	[0.1, 0.2, 0.3]
		Number of epochs ( $n_e$ )	20
		LSTM layers ( $l$ )	[1, 2, 3]
		Learning rate ( $\eta$ )	[0.1, 0.01, 0.001]
CT [20]	(end-to-end)	LSTM hidden units ( $d_h$ )	[256, 512, 1024]
		Balanced representation size ( $d_r$ )	[0.5 $d_h$ , 1 $d_h$ , 2 $d_h$ ]
		FC hidden units ( $d_{\text{fc}}$ )	[0.5 $d_h$ , 1 $d_h$ , 2 $d_h$ ]
		LSTM dropout rate ( $p$ )	[0.1, 0.2, 0.3]
		Number of epochs ( $n_e$ )	20
		Transformer blocks ( $B$ )	[1, 2]
R-MSNs [19]	Treatment/History Propensity Network	Learning rate ( $\eta$ )	[0.01, 0.001, 0.0001]
		Batch size	[64, 128]
		Attention heads ( $n_h$ )	2
		Transformer units ( $d_{\text{trm}}$ )	[0.5C, 1C]
		hidden representation ( $d_{\text{hr}}$ )	[0.5C, 1C]
		FC hidden units ( $d_{\text{fc}}$ )	[0.5 $d_{\text{hr}}$ , 1 $d_{\text{hr}}$ ]
	Encoder / Decoder	Number of epochs ( $n_e$ )	10
		LSTM dropout rate ( $p$ )	[0.1, 0.2]
		Max gradient norm	[0.5, 1.0, 2.0]
		Number of epochs ( $n_e$ )	10
		LSTM layers ( $l$ )	1
		Learning rate ( $\eta$ )	[0.01, 0.001, 0.0001]
G-Net [32]	(end-to-end)	LSTM hidden units ( $d_h$ )	[256, 512, 1024]
		FC hidden units ( $d_{\text{fc}}$ )	[0.5 $d_{\text{hr}}$ , 1 $d_{\text{hr}}$ ]
		LSTM dropout rate ( $p$ )	[0.1, 0.2, 0.4]
		Max gradient norm	[0.5, 1.0, 2.0]
		Number of epochs ( $n_e$ )	20
		LSTM layers ( $l$ )	[1, 2, 3]
GT [14]	(end-to-end)	Learning rate ( $\eta$ )	[0.1, 0.01, 0.001]
		LSTM hidden units ( $d_h$ )	[64, 128, 256]
		LSTM output size ( $d_o$ )	[0.5 $d_h$ , 1 $d_h$ , 2 $d_h$ ]
		Feed-forward hidden units ( $d_{\text{ff}}$ )	[0.5 $d_h$ , 1 $d_h$ , 2 $d_h$ ]
		LSTM dropout rate ( $p$ )	[0.1, 0.2]
		Number of epochs ( $n_e$ )	20
DeepBlip (ours)	Nuisance Network	Transformer blocks ( $B$ )	[1, 2]
		Learning rate ( $\eta$ )	[0.01, 0.001, 0.0001]
		Batch size	[64, 128]
		Attention heads ( $n_h$ )	2
		Transformer units ( $d_{\text{trm}}$ )	[0.5C, 1C]
		hidden representation ( $d_{\text{hr}}$ )	[0.5C, 1C]
	Blip Prediction Network	FC hidden units ( $d_{\text{fc}}$ )	[0.5 $d_{\text{hr}}$ , 1 $d_{\text{hr}}$ ]
		Number of epochs ( $n_e$ )	10
		LSTM dropout rate ( $p$ )	[0.1, 0.2]
		Balanced representation size ( $d_r$ )	[0.5 $d_h$ , 1 $d_h$ , 2 $d_h$ ]
		FC hidden units ( $d_{\text{fc}}$ )	[0.5 $d_h$ , 1 $d_h$ , 2 $d_h$ ]
		Number of epochs ( $n_e$ )	10



## 1118 **NeurIPS Paper Checklist**

### 1119 **1. Claims**

1120 Question: Do the main claims made in the abstract and introduction accurately reflect the  
1121 paper's contributions and scope?

1122 Answer: [\[Yes\]](#)

1123 Justification: The abstract and introduction claims that our proposed DeepBlip framework  
1124 addresses two key limitations as well as several practical strengths. These claims are  
1125 supported by SNMM (section 3) the framework (section 4) and experiment (section 5), with  
1126 both theoretical and empirical content.

1127 Guidelines:

- 1128 • The answer NA means that the abstract and introduction do not include the claims  
1129 made in the paper.
- 1130 • The abstract and/or introduction should clearly state the claims made, including the  
1131 contributions made in the paper and important assumptions and limitations. A No or  
1132 NA answer to this question will not be perceived well by the reviewers.
- 1133 • The claims made should match theoretical and experimental results, and reflect how  
1134 much the results can be expected to generalize to other settings.
- 1135 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
1136 are not attained by the paper.

### 1137 **2. Limitations**

1138 Question: Does the paper discuss the limitations of the work performed by the authors?

1139 Answer: [\[Yes\]](#)

1140 Justification: We discussed the limitations in Sec. 6.

1141 Guidelines:

- 1142 • The answer NA means that the paper has no limitation while the answer No means that  
1143 the paper has limitations, but those are not discussed in the paper.
- 1144 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1145 • The paper should point out any strong assumptions and how robust the results are to  
1146 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
1147 model well-specification, asymptotic approximations only holding locally). The authors  
1148 should reflect on how these assumptions might be violated in practice and what the  
1149 implications would be.
- 1150 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
1151 only tested on a few datasets or with a few runs. In general, empirical results often  
1152 depend on implicit assumptions, which should be articulated.
- 1153 • The authors should reflect on the factors that influence the performance of the approach.  
1154 For example, a facial recognition algorithm may perform poorly when image resolution  
1155 is low or images are taken in low lighting. Or a speech-to-text system might not be  
1156 used reliably to provide closed captions for online lectures because it fails to handle  
1157 technical jargon.
- 1158 • The authors should discuss the computational efficiency of the proposed algorithms  
1159 and how they scale with dataset size.
- 1160 • If applicable, the authors should discuss possible limitations of their approach to  
1161 address problems of privacy and fairness.
- 1162 • While the authors might fear that complete honesty about limitations might be used by  
1163 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
1164 limitations that aren't acknowledged in the paper. The authors should use their best  
1165 judgment and recognize that individual actions in favor of transparency play an impor-  
1166 tant role in developing norms that preserve the integrity of the community. Reviewers  
1167 will be specifically instructed to not penalize honesty concerning limitations.

### 1168 **3. Theory assumptions and proofs**

1169 Question: For each theoretical result, does the paper provide the full set of assumptions and  
1170 a complete (and correct) proof?

Answer: [Yes]

Justification: We provide full and detailed assumptions in Sec. 3 and Appendix B.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We disclose all the details needed to reproduce the experimntal results in Appendix D, G, E, G.6 and G.5.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

1224 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
1225 tions to faithfully reproduce the main experimental results, as described in supplemental  
1226 material?

1227 Answer: [Yes]

1228 Justification: We use open datasets from MIMIC in our experiment and provide the link to  
1229 an anonymous github repo: <https://anonymous.4open.science/r/DeepBlip-A39B>.

1230 Guidelines:

- 1231 • The answer NA means that paper does not include experiments requiring code.
- 1232 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 1233 • While we encourage the release of code and data, we understand that this might not be  
1234 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
1235 including code, unless this is central to the contribution (e.g., for a new open-source  
1236 benchmark).
- 1237 • The instructions should contain the exact command and environment needed to run to  
1238 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 1239 • The authors should provide instructions on data access and preparation, including how  
1240 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1241 • The authors should provide scripts to reproduce all experimental results for the new  
1242 proposed method and baselines. If only a subset of experiments are reproducible, they  
1243 should state which ones are omitted from the script and why.
- 1244 • At submission time, to preserve anonymity, the authors should release anonymized  
1245 versions (if applicable).
- 1246 • Providing as much information as possible in supplemental material (appended to the  
1247 paper) is recommended, but including URLs to data and code is permitted.

## 1250 6. Experimental setting/details

1251 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
1252 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
1253 results?

1254 Answer: [Yes]

1255 Justification: We include all the information needed to reproduce the results, includ-  
1256 ing dataset in Appendix D, implementation in Appendix G and hyperparameters in Ap-  
1257 pendix G.6,G.5.

1258 Guidelines:

- 1259 • The answer NA means that the paper does not include experiments.
- 1260 • The experimental setting should be presented in the core of the paper to a level of detail  
1261 that is necessary to appreciate the results and make sense of them.
- 1262 • The full details can be provided either with the code, in appendix, or as supplemental  
1263 material.

## 1264 7. Experiment statistical significance

1265 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
1266 information about the statistical significance of the experiments?

1267 Answer: [Yes]

1268 Justification: Yes we report standard deviation in the results for MIMIC. For the tumor  
1269 dataset we also take the average RMSE over five runs although std. dev is not directly  
1270 visualized in the plot.

1271 Guidelines:

- 1272 • The answer NA means that the paper does not include experiments.
- 1273 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
1274 dence intervals, or statistical significance tests, at least for the experiments that support  
1275 the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We report the computing resources in Appendix G.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: [NA]

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We discuss broader impact in the Sec. 6.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [\[Yes\]](#)

Justification: Our work does not pose such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [\[Yes\]](#)

Justification: We cite all the code and datasets used in our work.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: We do not provide any assets in our paper.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: Our work does not involve crowdsourcing or research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: Our work does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

1433 **16. Declaration of LLM usage**  
1434 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
1435 non-standard component of the core methods in this research? Note that if the LLM is used  
1436 only for writing, editing, or formatting purposes and does not impact the core methodology,  
1437 scientific rigorousness, or originality of the research, declaration is not required.  
1438 Answer: [NA]  
1439 Justification: The core method development for our work does not involve LLMs.  
1440 Guidelines:  
1441 • The answer NA means that the core method development in this research does not  
1442 involve LLMs as any important, original, or non-standard components.  
1443 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)  
1444 for what should or should not be described.