

---

# Offline Inverse Constrained Reinforcement Learning for Safe-Critical Decision Making in Healthcare

---

Anonymous Author(s)  
Affiliation  
Address  
email

## Abstract

1 Reinforcement Learning (RL) applied in healthcare can lead to unsafe medical  
2 decisions and treatment, such as excessive dosages or abrupt changes, often due  
3 to agents overlooking common-sense constraints. Consequently, Constrained  
4 Reinforcement Learning (CRL) is a natural choice for safe decisions. However,  
5 specifying the exact cost function is inherently difficult in healthcare. Recent  
6 Inverse Constrained Reinforcement Learning (ICRL) is a promising approach that  
7 infers constraints from expert demonstrations. ICRL algorithms model Markovian  
8 decisions in an interactive environment. These settings do not align with the  
9 practical requirement of a decision-making system in healthcare, where decisions  
10 rely on historical treatment recorded in an offline dataset. To tackle these issues, we  
11 propose the Constraint Transformer (CT). Specifically, 1) utilize causal attention  
12 mechanism to incorporate historical decisions and observations into the constraint  
13 modeling and employ a non-Markovian layer for weighted constraints to capture  
14 critical states, 2) generative world model to perform exploratory data augmentation,  
15 thereby enabling offline RL methods to generate unsafe decision sequences. In  
16 multiple medical scenarios, empirical results demonstrate that CT can capture  
17 unsafe states and achieve strategies that approximate lower mortality rates, reducing  
18 the occurrence probability of unsafe behaviors.

## 19 1 Introduction

20 In recent years, the doctor-to-patient ratio imbalance has drawn attention, with the U.S. having  
21 only 223.1 physicians per 100,000 people [1]. AI-assisted therapy emerges as a promising solution,  
22 offering timely diagnosis, personalized care, and reducing dependence on experienced physicians.  
23 Therefore, the development of an effective AI healthcare assistant is crucial.

24 Reinforcement learning (RL) offers a promising approach  
25 to develop AI assistants by addressing sequential decision-  
26 making tasks. However, this method can still lead to  
27 unsafe behaviors, such as administering excessive drug  
28 dosages, inappropriate adjustments of medical parameters,  
29 or abrupt changes in medication dosages. These behaviors,  
30 such as “too high” or “sudden change” can significantly  
31 endanger patients, potentially resulting in acute hypoten-  
32 sion, hypertension, arrhythmias, and organ damage, with  
33 fatal consequences [4, 5, 6]. For example, in sepsis treat-  
34 ment, patients receiving vasopressors (vaso) at dosages  
35 exceeding  $1\mu\text{g}/(\text{kg}\cdot\text{min})$  have a mortality rate of 90%  
36 [7]. Moreover, the “sudden change” in vaso can rapidly  
37 affect blood vessels, causing acute fluctuations in blood  
38 pressure and posing life-threatening risks to patients [8]. Our experiments demonstrate that the work

Table 1: The proportion of unsafe behaviors occurrences in vaso suggested by physician and DDPG. The typical range for vaso is  $0.1 \sim 0.2\mu\text{g}/(\text{kg}\cdot\text{min})$ , with doses exceeding 0.5 considered high [2]. A cutoff value of 0.75 is identified as a critical threshold associated with increased mortality [3].

| Drug dosage ( $\mu\text{g}/(\text{kg}\cdot\text{min})$ ) | Physician | DDPG     |
|--|-----------|----------|
| vaso > 0.75  | 2.27%     | 7.44% ↑  |
| vaso > 0.9   | 1.71%     | 7.40% ↑  |
| $\Delta$ vaso > 0.75                                     | 2.45%     | 21.00% ↑ |
| $\Delta$ vaso > 0.9                                      | 1.88%     | 20.62% ↑ |

$\Delta$  vaso: The change in vaso between two-time points.

39 [9] applying the Deep Deterministic Policy Gradient (DDPG) algorithm in sepsis indeed exhibits  
40 “too high” and “sudden change”<sup>1</sup> unsafe behaviors in vaso recommendations, as shown in Table 1.

41 This paper aims to achieve safe healthcare policy learning to mitigate unsafe behaviors. The most  
42 common method for learning safe policies is Constrained Reinforcement Learning (CRL) [10, 11],  
43 with the key to its success lying in the constraints representation. However, in healthcare, we can  
44 only design the cost function based on prior knowledge, which limits its application due to a lack of  
45 personalization, universality, and reliance on prior knowledge. For more details about issues, please  
46 refer to Appendix A. Therefore, Inverse Constrained Reinforcement Learning (ICRL) [12] emerges as  
47 a promising approach, as it can infer the constraints adhered to by experts from their demonstrations.  
48 However, directly applying ICRL in healthcare presents several challenges:

49 **1) The Markov decision is not compatible with medi-**  
50 **cal decisions.** ICRL algorithms model Markov decisions,  
51 where the next state depends only on the current state and  
52 not on the history [13, 14]. However, in healthcare, the  
53 historical states of patients are crucial for medical decision-  
54 making [15], as demonstrated in the experiments shown  
55 in Figure 1. Therefore, ICRL algorithms based on Markov  
56 assumption can not capture patient history, and ignore in-  
57 dividual patient differences, thereby limiting effectiveness.

58 **2) Interactive environment is not available for health-**  
59 **care or medical decisions.** ICRL algorithms [12, 16]  
60 follow an online learning paradigm, allowing agents to  
61 explore and learn from interactive environments. How-  
62 ever, unrestricted exploration in healthcare often entails  
63 unsafe behaviors that could breach constraints and result  
64 in substantial losses. Therefore, it is necessary to infer constraints using only offline datasets.

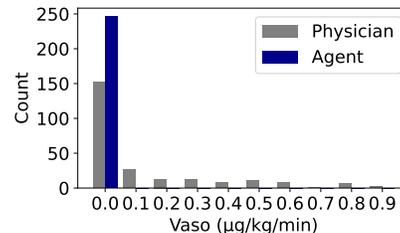


Figure 1: The distribution of vaso for patients with the same state. The physician makes different decisions due to referencing historical information, while the agent based on Markov decision-making can only make the same decision.

65 In this paper, we introduce offline Constraint Transformer (CT), a novel ICRL framework that  
66 incorporates patients’ historical information into constraint modeling and learns from offline data to  
67 infer constraints in healthcare. Specifically,

68 1) Inspired by the recent success of transformers in sequence modeling [17, 18, 19], we incorporate  
69 historical decisions and observations into constraint modeling using a causal attention mechanism. To  
70 capture key events in trajectories, we introduce a non-Markovian transformer to generate constraints  
71 and importance weights, and then define constraints using weighted sums. CT takes trajectories as  
72 input, allowing for the observation of patients’ historical information and evaluation of key states.

73 2) To learn from an offline dataset, we introduce a model-based offline RL method that simultaneously  
74 learns a policy model and a generative world model via auto-regressive imitation of the actions and  
75 observations in medical decisions. The policy model employs a stochastic policy with entropy  
76 regularization to prevent it from overfitting and improve its robustness. Utilizing expert datasets,  
77 the generative world model uses an auto-regressive exploration generation paradigm to effectively  
78 discover a set of violating trajectories. Then, CT can infer constraints in healthcare through these  
79 unsafe trajectories and expert trajectories.

80 In the medical scenarios of sepsis and mechanical ventilation, we conduct experimental evaluations of  
81 offline CT. Experimental evaluations demonstrate that offline CT can capture patients’ unsafe states  
82 and assign higher penalties, thereby providing more interpretable constraints compared to previous  
83 works [9, 20, 21]. Compared to unconstrained and custom constraints, CT achieves strategies that  
84 closely approximate lower mortality rates with a higher probability (improving by 8.85% compared to  
85 DDPG). To investigate the avoidance of unsafe behaviors with offline CT, we evaluate the probabilities  
86 of “too high” and “sudden changes” occurring in the sepsis. The experimental results show that CRL  
87 with CT can reduce the probability of unsafe behaviors to zero.

## 88 2 Related Works

89 **Reinforcement Learning in Healthcare.** RL has made great progress in the realm of healthcare, such  
90 as sepsis treatment [9, 20, 21, 22], mechanical ventilation [23, 24, 25], sedation [26] and anesthesia

<sup>1</sup>In sepsis, “too high” indicates that the dosage of the vaso medication exceeds the threshold. And “sudden change” indicates that the change in vaso medication dosage between two time points exceeds the threshold.

91 [27, 28]. However, these works mentioned above have not addressed potential safety issues such as  
 92 sudden changes or too high doses of medication. Therefore, the development of policies that are both  
 93 safe and applicable across various healthcare domains is crucial.

94 **Inverse Constrained Reinforcement Learning.** Previous works inferred constraint functions by  
 95 determining the feasibility of actions under current states. In discrete state-action spaces, Chou *et al.*  
 96 [29] and Park *et al.* [30] learned constraint sets to differentiate constrained state-action pairs. Scobee  
 97 & Sastry [31] proposed inferring constraint sets based on the principle of maximum entropy, while  
 98 some studies [32, 33] extended this approach to stochastic environments using maximum causal  
 99 entropy [34]. In continuous domains, Malik *et al.* [12], Gaurav *et al.* [16], and Qiao *et al.* [35] used  
 100 neural networks to approximate constraints. Some works [11, 29] applied Bayesian Monte Carlo and  
 101 variational inference to infer the posterior distribution of constraints in high-dimensional state spaces.  
 102 Xu *et al.* [36] modeled uncertainty perception constraints for arbitrary and epistemic uncertainties.  
 103 However, these methods can only be applied online and lack historical dependency.

104 **Transformers for Reinforcement Learning.** Transformer has produced exciting progress on RL  
 105 sequential decision problems [17, 18, 37, 38]. These works no longer explicitly learn Q-functions  
 106 or policy gradients, but focus on action sequence prediction models driven by target rewards. Chen  
 107 *et al.* [18] and Janner *et al.* [37] perform auto-regressive modeling of trajectories to achieve policy  
 108 learning in an offline environment. Furthermore, Zheng *et al.* [17] unify offline pretraining and  
 109 online fine-tuning within the Transformer framework. Liu *et al.* [38] and Kim *et al.* [19] integrate the  
 110 transformer architecture into constraint learning and preference learning. The transformer architecture,  
 111 with its sequence modeling capability and independence from the Markov assumption, can capture  
 112 temporal dependencies in medical decision-making. Thus, it is well-suited for trajectory learning and  
 113 personalized learning in medical settings.

### 114 3 Problem Formulation

115 We model the medical environment with a Constrained Markov Decision Process (CMDP)  $\mathcal{M}^c$  [39],  
 116 which can be defined by a tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \mathcal{C}, \gamma, \kappa, \rho_0)$ . Similar to studies [23, 40], we extract data  
 117 within 72 hours of patient admission, with each 4-hour interval constituting a window or time step.  
 118 The state indicators of the patient at each time step are denoted as  $s \in \mathcal{S}$ . The administered drug  
 119 doses or instrument parameters of interest are considered as actions  $a \in \mathcal{A}$ , while reward function  
 120  $\mathcal{R}$  is used to describe the quality of the patient’s condition and provided by experts based on prior  
 121 work [9, 23]. At each time step  $t$ , an agent performs an action  $a_t$  at a patient’s state  $s_t$ . This process  
 122 generates the reward  $r_t \sim \mathcal{R}(s_t, a_t)$ , the cost  $c_t \sim \mathcal{C}$  and the next state  $s_{t+1} \sim \mathcal{P}(\cdot | s_t, a_t)$ , where  
 123  $\mathcal{P}$  defines the transition probabilities.  $\gamma$  denotes the discount factor.  $\kappa \in \mathbb{R}_+$  denotes the bound of  
 124 cumulative costs.  $\rho_0$  defines the initial state distribution. The goal of the CRL policy  $\pi$  is to maximize  
 125 the reward return while limiting the cost in a threshold  $\kappa$ :

$$\arg \max_{\pi} \mathbb{E}_{\pi, \rho_0} [\sum_{t=1}^T \gamma^t r_t], \quad \text{s.t.} \quad \mathbb{E}_{\pi, \rho_0} [\sum_{t=1}^T \gamma^t c_t] \leq \kappa. \quad (1)$$

126 where  $T$  is the length of the trajectory  $\tau$ . CRL commonly assumes that constraint signals are directly  
 127 observable. However, in healthcare, such signals are not easily obtainable. Therefore, Our objective  
 128 is to infer reasonable constraints for CRL to achieve safe policy learning in healthcare.

129 **Safe-Critical Decision Making with Constraint Inference in Healthcare.** Our general goal is for  
 130 our policy to approximate the optimal policy, which refers to the strategy under which the patient’s  
 131 mortality rate is minimized (achieving a zero mortality rate is often difficult since there are patients  
 132 who can not recover, regardless of all potential future treatment sequences [41]). Decision-making  
 133 with constraints can formulate safer strategies by discovering and avoiding unsafe states, thereby  
 134 approaching the optimal policy.

135 However, most offline RL algorithms rely on online evaluation, where the agent is evaluated in  
 136 an interactive environment, whereas in medical scenarios, only offline evaluation can be utilized.  
 137 In previous works [5, 9, 40, 42], they qualitatively analyzed by comparing the differences (DIFF)  
 138 between the drug dosage recommended by our policy  $\pi$  and the dosage administered by clinical  
 139 physicians  $\hat{\pi}$ , and its relationship with mortality rates, through graphical analysis. In the graph  
 140 depicting the relationship between the DIFF and mortality rate, at the point when DIFF is zero, the  
 141 lower the mortality rate of patients, the better the performance of the policy [40]. To provide a more  
 142 accurate quantitative evaluation, we introduce the concept of the probability of approaching the  
 143 optimal policy, defined as  $\omega$ :

$$\omega = \frac{\text{Number of survivors among the top } N \text{ patients}}{N} \quad (2)$$

144 We randomly collect  $2N$  patients (with an equal number of known survivors and non-survivors under  
 145 doctor’s policy  $\hat{\pi}$ ) from the offline dataset. We then calculate the DIFF and sort it in ascending order.  
 146 The optimality of the policy can be evaluated through the following two points: 1) The higher the  
 147 survival probability (i.e.,  $\omega$ ) of the top  $N$  patients, the lower the mortality rate can be achieved by  
 148 executing  $\pi$ ; 2) The smaller the DIFF among the surviving patients in the top  $N$ , the greater the  
 149 probability that  $\pi$  is optimal.

## 150 4 Method

151 To infer constraints and achieve safe decision-making in healthcare, we introduce the Offline Con-  
 152 straint Transformer (Figure 2), a novel ICRL framework.

153 **Inverse Constrained Reinforcement Learning.** ICRL aims to recover the cost function  $C^*$  by  
 154 leveraging a set of trajectories  $\mathcal{D}_e = \{\tau_e^{(i)}\}_i^N$  sampled from an expert policy  $\pi_e$ , where  $N$  denotes  
 155 the number of the trajectories. ICRL is commonly based on the Maximum Entropy framework [31],  
 156 and the likelihood function is articulated as [12]:

$$p(\mathcal{D}_e | \mathcal{C}) = \frac{1}{(Z_{\mathcal{M}^{\mathcal{C}}})^N} \prod_{i=1}^N \exp [R(\tau^{(i)})] \mathbb{I}^{\mathcal{M}^{\mathcal{C}}}(\tau^{(i)}) \quad (3)$$

157 Here,  $Z_{\mathcal{M}} = \int \exp(\beta r(\tau)) \mathbb{I}^{\mathcal{M}}(\tau) d\tau$  is the normalizing term. The indicator  $\mathbb{I}^{\mathcal{M}^{\mathcal{C}}}(\tau^{(i)})$  signifies the  
 158 extent to which the trajectory  $\tau^{(i)}$  satisfies the constraints. It can be approximated using a neural  
 159 network  $\zeta_{\theta}(\tau^{(i)})$  parameterized with  $\theta$ , defined as  $\zeta_{\theta}(\tau^{(i)}) = \prod_{t=0}^T \zeta_{\theta}(s_t^i, a_t^i)$ . Consequently, the  
 160 cost function can be formulated as  $C_{\theta} = 1 - \zeta_{\theta}$ . Substituting the neural network for the indicator, we  
 161 can update  $\theta$  through the gradient of the log-likelihood function:

$$\nabla_{\theta} \mathcal{L}(\theta) = \mathbb{E}_{\tau^{(i)} \sim \pi_e} [\nabla_{\theta} \log[\zeta_{\theta}(\tau^{(i)})]] - \mathbb{E}_{\hat{\tau} \sim \pi_{\mathcal{M}^{\hat{C}_{\theta}}}} [\nabla_{\theta} \log[\zeta_{\theta}(\hat{\tau}^{(i)})]] \quad (4)$$

162 where  $\mathcal{M}^{\hat{C}_{\theta}}$  denotes the MDP obtained after augmenting  $\mathcal{M}$  with the cost function  $C_{\theta}$ , using the  
 163 executing policy  $\pi_{\mathcal{M}^{\hat{C}_{\theta}}}$ . And  $\hat{\tau}$  are sampled from the policy. In practice, ICRL can be conceptualized  
 164 as a bi-level optimization task [11]. We can 1) update this policy based on Equation 1, and 2) employ  
 165 Equation 4 for constraint learning. Intuitively, the objective of Equation 4 is to distinguish between  
 166 trajectories generated by expert policies and imitation policies that may violate the constraints.

167 Specifically, task 1) involves updating the policy using advanced CRL methods. Significant progress  
 168 has been made in some works such as BCQ-Lagrangian (BCQ-Lag), COpiDICE [43], VOCE [44],  
 169 and CDT [38]. Meanwhile, task 2) focuses on learning the constraint function, as shown in Figure  
 170 2. Our research primarily improves the latter process due to two main challenges facing ICRL  
 171 in healthcare: **Challenge 1**) pertains to the limitations of the Markov property, and **Challenge 2**)  
 172 involves the issue of inferring constraints only from offline datasets. To address these challenges, we  
 173 propose the offline CT as our solution.

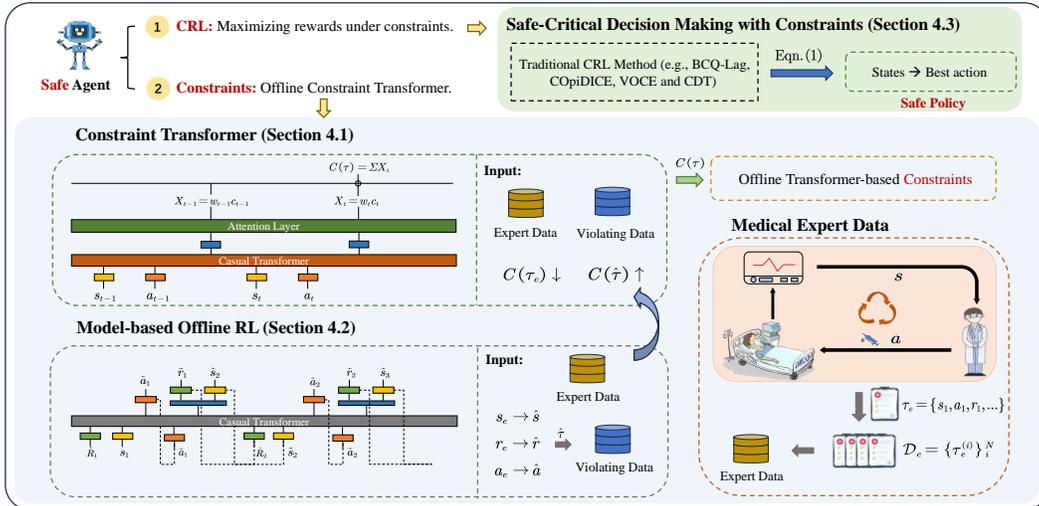


Figure 2: The overview of the safe healthcare policy learning with offline CT.

174 **Offline Constraint Transformer.** To address the first challenge, we delve into the inherent issues of  
 175 applying the Markov property to healthcare and draw inspiration from the successes of Transformer  
 176 in decision-making, redefining the representation of the constraints. To realize the offline training, we  
 177 consider the essence of ICRL updates, proposing a model-based RL to generate unsafe behaviors  
 178 used to train CT. We outline three parts: establishing the constraint representation model (Section  
 179 4.1), creating an offline RL for violating data (Section 4.2), and learning safe policies (Section 4.3).

## 180 4.1 Constraint Transformer

181 ICRL methods relying on the Markov property overlook patients’ historical informa-  
 182 tion, focusing only on the current state. However, both current and historical states,  
 183 along with vital sign changes are crucial for a human doctor’s decision-making pro-  
 184 cess [15]. To emulate the observational approach of humans, we draw inspiration  
 185 from the Decision Transformer (DT) [18] to incorporate historical information into  
 186 constraints for a more comprehensive observation and judgment. We propose a  
 187 constraint modeling approach based on a causal attention mechanism, as shown in Figure 3. The structure comprises a causal Transformer for  
 188 sequential modeling and a non-Markovian layer for weighted constraints learning.

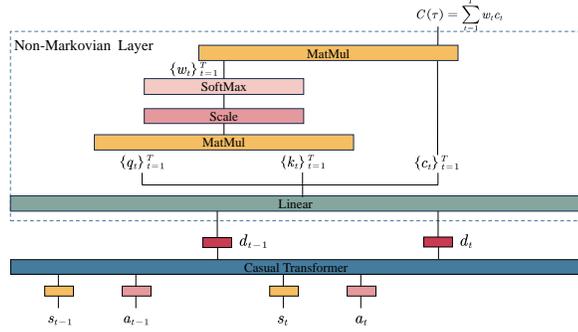


Figure 3: The structure of the Constraint Transformer. The structure comprises a causal Transformer for sequential modeling and a non-Markovian layer for weighted constraints learning.

196 **Sequential Modeling for Constraints Inference.** For a trajectory segment of length  $T$ ,  $2T$  input  
 197 embeddings are generated, with each position containing state  $s$  and action  $a$  embeddings. Addi-  
 198 tionally, these embeddings undergo linear and normalization layers before being fed into the causal  
 199 Transformer, which produces output embeddings  $\{d_t\}_{t=1}^T$  determined by preceding input embeddings  
 200 from  $(s_1, a_1, \dots, s_T, a_T)$ . Here,  $d_t$  depends only on the previous  $t$  states and actions.

201 **Modeling Non-Markovian for Weighted Constraints Learning.** Although  $d_t$  represents the cost  
 202 function  $c_t$  derived from observations over long trajectories, it doesn’t pinpoint which previous key  
 203 actions or states led to its increase. In healthcare, identifying key actions or states is vital for analyzing  
 204 risky behaviors and status, and enhancing model interpretability. To address this, we draw inspiration  
 205 from the design of the preference attention layer in [19] and introduce an additional attention layer.  
 206 This layer is employed to define the cost weight for non-Markovians. It takes the output embeddings  
 207 from the causality transformer as input and generates the corresponding cost and importance weights.  
 208 The output of the attention layer is computed by weighting the values through the normalized dot  
 209 product between the query and other keys:

$$\sum_{t=1}^T \text{softmax} \left( \left\{ \langle q_t, k_{t'} \rangle \right\}_{t'=1}^T \right)_t \cdot c_t = \sum_{t=1}^T w_t \cdot c_t \quad (5)$$

210 Here, the key  $k_t \in \mathbb{R}^m$ , query  $q_t \in \mathbb{R}^m$ , and value  $c_t \in \mathbb{R}^m$  are derived from the  $t$ -th input  $d_t$   
 211 through linear transformations, where  $m$  denotes the embedding dimension. Furthermore, for each  
 212 time step  $t$ , since  $d_t$  depends only on the previous state-action pairs  $\{(s_i, a_i)\}_{i=1}^t$  and serves as the  
 213 input embedding for the attention layer,  $c_t$  is also associated solely with the preceding  $t$  time steps.  
 214 The representation of the cost function as a weighted sum is defined as  $C(\tau) = \sum_{t=1}^T w_t \cdot c_t$ . Then,  
 215 we can also determine the constraint function values for each preceding subsequence. Introducing the  
 216 newly defined cost function, we redefine Equation 4 for CT as:

$$\nabla_{\phi} \mathcal{L}(\phi) = \mathbb{E}_{\hat{\tau} \sim \mathcal{D}_v} [\nabla_{\phi} \log[C_{\phi}(\hat{\tau})]] - \mathbb{E}_{\tau \sim \mathcal{D}_e} [\nabla_{\phi} \log[C_{\phi}(\tau)]] \quad (6)$$

217 where  $\phi$  is the parameter of CT,  $\mathcal{D}_e$  and  $\mathcal{D}_v$  represent the expert data and the violating data. This  
 218 formulation implies that the constraint should be minimized on the expert policy and maximized on  
 219 the violating policy. We construct an expert and a violating dataset to evaluate Equation 6 in offline.  
 220 The expert data can be acquired from existing medical datasets or hospitals. Regarding the violating  
 221 dataset, we introduce a generative model to establish it, as detailed in Section 4.2.

## 222 4.2 Model-based Offline RL

223 To train CT offline, we introduce a model-  
 224 based offline RL method (Figure 4) to gener-  
 225 ate violating data that refers to unsafe  
 226 behavioral data and can be represented as  
 227  $\tau_v = (s_1, a_1, r_1, s_2, \dots) \in \mathcal{D}_v$ . The model  
 228 simultaneously learns a policy model and a  
 229 generative world model via auto-regressive imitation of the actions and observations in healthcare.  
 230 The model processes a trajectory,  $\tau_e \in \mathcal{D}_e$ , as a sequence of tokens encompassing the return-to-go,  
 231 states, and actions, defined as  $(\hat{R}_1, s_1, a_1, \dots, \hat{R}_T, s_T, a_T)$ . Notably, the return-to-go  $\hat{R}_t$  at timestep  
 232  $t$  is the sum of future rewards, calculated as  $\hat{R}_t = \sum_{t'=t}^T r_{t'}$ . At each timestep  $t$ , it employs  
 233 the tokens from the preceding  $K$  timesteps as its input, where  $K$  represents the context length.  
 234 Thus, the input tokens for it at timestep  $t$  are denoted as  $h_t = \{\hat{R}_{-K:t}, s_{-K:t}, a_{-K:t-1}\}$ , where  
 235  $\hat{R}_{-K:t} = \{\hat{R}_K, \dots, \hat{R}_t\}$ ,  $s_{-K:t} = \{s_K, \dots, s_t\}$  and  $a_{-K:t-1} = \{a_K, \dots, a_{t-1}\}$ .

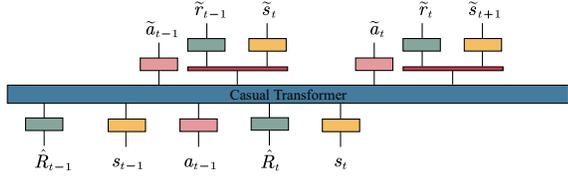


Figure 4: The structure of the model-based offline RL.

236 **Policy Model.** The input tokens are encoded through a linear layer for each modality. Subsequently,  
 237 the encoded tokens pass through a casual transformer to predict future action tokens. We use  
 238 a stochastic policy [38] to achieve policy learning. Additionally, we utilize a Shannon entropy  
 239 regularizer  $\mathcal{H}[\pi_{\vartheta}(\cdot | h)]$  to prevent policy overfitting and enhance robustness. The optimization  
 240 objective is to minimize the negative log-likelihood loss while maximizing the entropy with weight  $\lambda$ :

$$241 \min_{\vartheta} \mathbb{E}_{h_t \sim \mathcal{D}_e} [-\log \pi_{\vartheta}(\cdot | h_t) - \lambda \mathcal{H}[\pi_{\vartheta}(\cdot | h_t)]] \quad (7)$$

242 where the policy  $\pi_{\vartheta}(\cdot | h_t) = \mathcal{N}(\mu_{\vartheta}(h_t), \Sigma_{\vartheta}(h_t))$  adopts the stochastic Gaussian policy represen-  
 243 tation and  $\vartheta$  is the parameter.

244 **Generative World Model.** To predict states and rewards, we use  $x_t = \{h_t \cup a_t\}$  as input encoded  
 245 by linear layers. The encoded tokens pass through the casual transformer to predict hidden tokens.  
 246 Then we utilize two linear layers to fit the rewards and states. The optimization objective for the two  
 247 linear layers  $\ell$  with the parameters  $\varphi$  and  $\mu$  can be defined as:

$$\min_{\varphi, \mu} \mathbb{E}_{s_t, r_{t-1} \in x_t \sim \mathcal{D}_e} [(s_t - \ell_{\varphi}(x_t))^2 + (r_{t-1} - \ell_{\mu}(x_t))^2] \quad (8)$$

248 **Generating Violating Data.** In RL, excessively high rewards, surpassing those provided by domain  
 249 experts, may incentivize agents to violate the constraints in order to maximize the total reward [11].  
 250 Therefore, we set a high initial target reward  $\hat{R}_1$  to obtain violation data. We feed  $\hat{R}_1$  and initial state  
 251  $s_1^{(i)}$  into the model-based offline RL to generate  $\tau_v^{(i)}$  in an auto-regressive manner, as depicted in  
 252 model-based offline RL of Figure 2, where  $\tilde{a}$ ,  $\tilde{r}$  and  $\tilde{s}$  are predicted by the model. The target reward  
 253  $\hat{R}$  decreases incrementally and can be represented as  $\hat{R}_{t+1} = \hat{R}_t - \tilde{r}_t$ . Considering the average error  
 254 in trajectory prediction, we generate trajectories with the length  $K = 10$ , as detailed in Appendix  
 255 B.3. Repeating  $N$  initial states, we can get violating data  $\mathcal{D}_v = \{\tau_v^{(i)}\}_{i=1}^N$ .

256 Note that certain other generative models, such as Variational Auto-Encoder (VAE) [45], Generative  
 257 Adversarial Networks (GAN) [46, 47], and Denoising Diffusion Probabilistic Models (DDPM)  
 258 [48, 49], may be better at generating data. We introduce the model-based offline RL primarily  
 259 because it has been shown to generate violating data with exploration [38] and possess the ability to  
 260 process time-series features efficiently.

### 261 4.3 Safe-Critical Decision Making with Constraints.

262 To train offline CT, we gather the medical expert dataset  $\mathcal{D}_e$  from the environment. Then, we employ  
 263 gradient descent to train the model-based offline RL, guided by Equation 7 and Equation 8, continuing  
 264 until the model converges. Using this RL model, we automatically generate violating data denoted  
 265 as  $\mathcal{D}_v$ . Subsequently, CT is optimized based on Equation 6 to get the cost function  $C$ , leveraging  
 266 samples from both  $\mathcal{D}_e$  and  $\mathcal{D}_v$ . To learn a safe policy, we train the policy  $\pi$  using  $C$  until it converges  
 267 based on Equation 1. The detailed training procedure is presented in Algorithm 1.

## 268 5 Experiment

269 In this section, we first provide a brief overview of the task, as well as data extraction and prepro-  
 270 cessing. Subsequently, in Section 5.1, we demonstrate that CT can describe constraints in healthcare  
 271 and capture critical patient states. We emphasize its applicability to various CRL methods and its  
 272 ability to approach the optimal policy for reducing mortality rates in Section 5.2. Finally, Section 5.3  
 273 discusses the realization of the objective of safe medical policies.

---

**Algorithm 1** Safe Policy Learning with Offline CT
 

---

**Input:** Expert trajectories  $\mathcal{D}_e$ , context length  $K$ , target reward  $\hat{R}_1$ , samples  $N$ , episode length  $T$

- 1: Train model-based offline RL  $\mathcal{M}$ : Update  $\vartheta$ ,  $\varphi$  and  $\mu$  using the Equation (7) and Equation (8)
- 2: **for**  $t = 1, \dots, T$  **do**
- 3:   Sample initial states  $S_1$  from  $\mathcal{D}_e$
- 4:   Generate the violating dataset:  $\mathcal{D}_v \leftarrow \mathcal{M}.\text{generate\_data}(S_1, \hat{R}_1, K)$
- 5:   Sample set of trajectories  $\{\tau_e^{(i)}\}_{i=1}^N$  and  $\{\tau_v^{(i)}\}_{i=1}^N$  from  $\mathcal{D}_e$  and  $\mathcal{D}_v$
- 6:   Train offline CT: Use  $\{\tau_e^{(i)}\}_{i=1}^N$  and  $\{\tau_v^{(i)}\}_{i=1}^N$  to update  $\phi$  based on Equation (6)
- 7:   Safe policy learning: Update  $\pi$  using the cost function  $C_\phi(\tau)$  based on Equation (1)
- 8: **end for**

**Output:**  $\pi$  and  $C(\tau)$

---

274 **Tasks.** We primarily use the sepsis task that is commonly used in previous works [9, 20, 42, 22], and  
 275 supplement some experiments on the mechanical ventilator task [23, 50]. The detailed definition of  
 276 the two tasks mentioned above can be found in Appendix B.1 and B.2.

277 **Data Extraction and Pre-processing.** Our medical dataset is derived from the Medical Information  
 278 Mart for Intensive Care III (MIMIC-III) database [51]. For each patient, we gather relevant physio-  
 279 logical parameters, including demographics, lab values, vital signs, and intake/output events. Data is  
 280 grouped into 4-hour windows, with each window representing a time step. In cases of multiple data  
 281 points within a step, we record either the average or the sum. We eliminate variables with significant  
 282 missing values and use the  $k$ -nearest neighbors method to fill in the rest. Notably, the training dataset  
 283 consists of data from surviving patients, while the validation set includes survivors and non-survivors.

284 **Model-based Offline RL Evaluation.** To ensure the rigor of the experiments, we evaluate the validity  
 285 of the model-based offline RL, as detailed in Appendix B.3.

### 5.1 Can Offline CT Learn Effective Constraints?

287 In this section, we primarily assess the efficacy of the cost function  
 288 learned by offline CT in sepsis, focusing particularly on its capa-  
 289 bility to evaluate patient mortality rates and capture critical events.  
 290 First, we employ the cost function to compute cost values for the  
 291 validation dataset. Subsequently, we statistically analyze the rela-  
 292 tionship between these cost values and mortality rates. As shown in  
 293 Figure 5, there is an increase in patient mortality rates with rising  
 294 cost values. It’s noteworthy that such increases in mortality rates are  
 295 often attributed to suboptimal medical decisions. Therefore, these  
 296 experimental findings affirm that the cost values effectively reflect  
 297 the quality of medical decision-making. To observe the impact of the attention layer (non-Markovian  
 298 layer), we conduct experiments by removing the attention layer from CT. The results reveal that the  
 299 penalty values do not correlate proportionally with mortality rates. This indicates that the attention  
 300 layer plays a crucial role in assessing constraints.

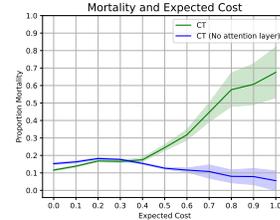


Figure 5: The relationship between cost and mortality.

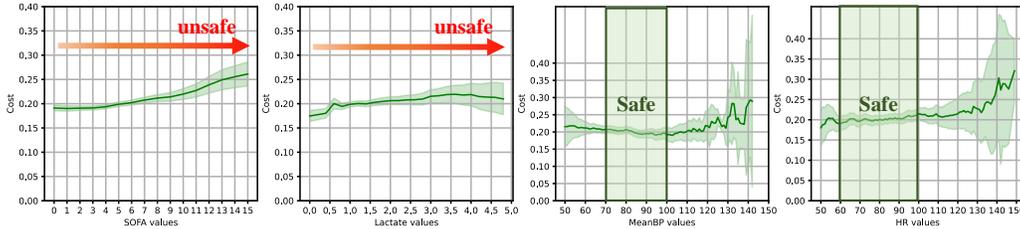


Figure 6: The relationship between physiological indicators and cost values. As SOFA and lactate levels become increasingly unsafe, the cost increases. Mean BP and HR at lower values within the safe range incur a lower cost, but as they move into unsafe ranges, the cost increases, penalizing previous state-action pairs. The cost can differentiate between relatively safe and unsafe regions.

301 To assess the capability of the cost function to capture key events, we analyze the relationship  
 302 between physiological indicators and cost values. We focus on four key indicators in sepsis treatment:  
 303 Sequential Organ Failure Assessment (SOFA) score [52], lactate levels [53], Mean Arterial Pressure

304 (MeanBP) [54], and Heart Rate (HR) [55]. The SOFA score and lactate levels are critical indicators  
 305 for assessing sepsis severity, with higher values indicating greater patient risk. MeanBP and HR  
 306 are essential physiological metrics, typically ranging from 70 to 100 mmHg and 60 to 100 beats,  
 307 respectively. Deviations from these ranges can signify patient risk. As depicted in Figure 6, the cost  
 308 values effectively distinguish between high-risk and safe conditions, reflecting changes in patient  
 309 status. Additional details on other parameters’ relationship with cost are in Appendix B.4.

## 310 5.2 Can Offline CT Improve the Performance of CRL?

311 **Baselines.** We adopt the DDPG method as the baseline in sepsis research [9], and the Double Deep  
 312 Q-Learning (DDQN) and Conservative Q-Learning (CQL) methods as baselines in ventilator research  
 313 [23]. Since there are no other offline inverse reinforcement learning works available for reference,  
 314 we have included two additional settings: no cost and custom cost. In the case of no cost, the cost is  
 315 set to zero, while the design of custom constraints is outlined in Appendix A. These settings help  
 316 evaluate whether CT can infer effective constraints.

317 **Metrics.** To assess effectiveness, we use  $\omega$  to indicate the probability that the policy is optimal and  
 318 analyze the relationship between DIFF and mortality rate through a graph. Recently, Kondrup *et*  
 319 *al.* [23] use the Fitted Q Evaluation (FQE) [56] to evaluate the policy in healthcare. However, the  
 320 value estimates of FQE depend solely on the dataset  $\mathcal{D}$  and the actions chosen by the policy  $\pi$  used to  
 321 train FQE. This reliance can lead to inaccurate estimates when evaluating unseen state-action pairs.  
 322 Therefore, we do not adopt this method as an evaluation metric.

323 **Results.** We combine our method CT with common CRL algorithms (e.g., VOCE, CopiDICE,  
 324 BCQ-Lag, and CDT), and compare them with both no-cost and custom cost settings. Each  
 325 CRL model is trained using no cost, custom cost,  
 326 and CT separately, with other parameters set the  
 327 same during training. For evaluation metrics,  
 328 we use IV difference (IV DIFF), vaso difference (VASO DIFF), and combined [IV, VASO]  
 329 difference (ACTION DIFF) as the metrics to  
 330 be ranked. We measure the mean and variance  
 331 of  $\omega\%$  in 10 sets of random seeds, and the re-  
 332 sults are shown in Table 2. From the results,  
 333 we can conclude: (1) In different CRL meth-  
 334 ods, CT consistently makes the strategy closer  
 335 to the one with lower mortality rates, with a  
 336 probability 8.85% higher than DDPG. (2) We  
 337 find that CDT+CT achieves better results on all  
 338 three metrics. CDT is also a transformer-based  
 339 method, which indicates that transformer-based  
 340 architecture indeed exhibits more outstanding  
 341 performance in healthcare.  
 342  
 343  
 344

345 Figure 7 illustrates the relationship between IV  
 346 and VASO DIFF with mortality rates under the  
 347 DDPG and CDT+CT methods in sepsis. In  
 348 VASO DIFF, when the gap is zero, the mor-  
 349 tality rate under CDT+CT is lower than that  
 350 under DDPG, indicating that following the for-  
 351 mer strategy could lead to a lower mortality  
 352 rate. Similarly, in IV DIFF, the same trend is ob-  
 353 served. Notably, for the IV strategy, the lowest  
 354 mortality rate for DDPG does not occur at the  
 355 point where the difference is zero, indicating a  
 356 significant estimation bias.

357 In addition, corresponding experiments are conducted on the mechanical ventilator, as shown in  
 358 Figure 8. Compared to previous methods DDQN and CQL, under the CDT+CT approach, a noticeable  
 359 trend is observed where the proportion of mortality rates increases with increasing differences. When

Table 2: Performance of sepsis strategies under various offline CRL models and different constraints.

| $\omega\%$       | COST               | IV DIFF $\uparrow$               | VASO DIFF $\uparrow$             | ACTION DIFF $\uparrow$           |
|------------------|--------------------|----------------------------------|----------------------------------|----------------------------------|
| DDPG             | -                  | 50.95 $\pm$ 1.34                 | 51.45 $\pm$ 0.75                 | 51.15 $\pm$ 1.15                 |
| VOCE             | No cost            | 47.45 $\pm$ 0.52                 | 46.35 $\pm$ 1.82                 | 51.00 $\pm$ 0.86                 |
|                  | Custom cost        | 46.45 $\pm$ 0.46                 | 52.00 $\pm$ 0.98                 | 49.40 $\pm$ 1.04                 |
|                  | CT                 | <b>53.33<math>\pm</math>0.94</b> | <b>59.04<math>\pm</math>1.13</b> | <b>56.15<math>\pm</math>1.08</b> |
| CopiDICE         | No cost            | 48.30 $\pm$ 0.91                 | 60.10 $\pm$ 0.6                  | 51.25 $\pm$ 0.70                 |
|                  | Custom cost        | <b>53.05<math>\pm</math>1.35</b> | 55.20 $\pm$ 0.24                 | 53.90 $\pm$ 1.04                 |
|                  | CT                 | 51.95 $\pm$ 0.41                 | <b>60.85<math>\pm</math>1.08</b> | <b>54.60<math>\pm</math>0.60</b> |
| BCQ-Lag          | No cost            | 47.50 $\pm$ 1.32                 | 51.05 $\pm$ 0.61                 | 49.35 $\pm$ 1.08                 |
|                  | Custom cost        | 51.54 $\pm$ 0.16                 | <b>56.23<math>\pm</math>1.43</b> | 53.69 $\pm$ 1.62                 |
|                  | CT                 | <b>52.45<math>\pm</math>1.01</b> | 55.34 $\pm$ 1.20                 | <b>54.39<math>\pm</math>0.86</b> |
| CDT              | No cost            | 56.50 $\pm$ 0.81                 | 62.45 $\pm$ 1.20                 | 58.90 $\pm$ 1.34                 |
|                  | Custom cost        | 54.70 $\pm$ 1.12                 | 59.85 $\pm$ 1.51                 | 57.80 $\pm$ 1.00                 |
|                  | CT                 | <b>57.15<math>\pm</math>1.67</b> | <b>65.20<math>\pm</math>1.22</b> | <b>60.00<math>\pm</math>1.49</b> |
| CDT              | Without CT         | 56.50 $\pm$ 0.81                 | 62.45 $\pm$ 1.20                 | 58.90 $\pm$ 1.34                 |
| CDT              | No attention layer | 55.25 $\pm$ 1.46                 | 64.00 $\pm$ 1.54                 | 57.90 $\pm$ 0.78                 |
| Generative Model | -                  | 55.49 $\pm$ 2.55                 | 56.60 $\pm$ 1.33                 | 57.00 $\pm$ 2.06                 |

**Blue:** Safe policy is closer to the optimal policy.  $\uparrow$ : higher is better.

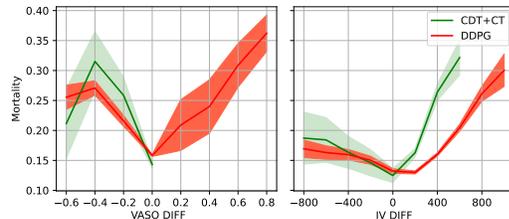
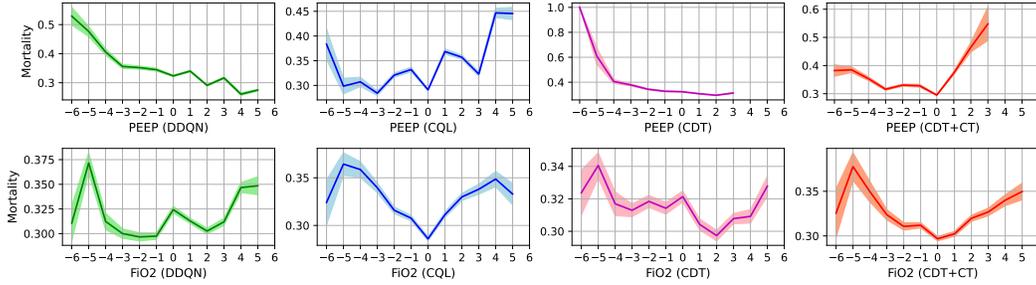


Figure 7: The relationship between DIFF and the mortality rate in sepsis. The x-axis represents the DIFF. The y-axis indicates the mortality rate of patients at a given DIFF. The solid line represents the mean, while the shaded area indicates the Standard Error of the Mean (SEM).

360 there is a significant difference in DIFF, the results may be unreliable, possibly due to the limited  
 361 data distribution in the tail.



362 Figure 8: The relationship between the DIFF of actions and mortality in mechanical ventilator. The  
 363 actions mainly consist of Positive End Expiratory Pressure (PEEP) and Fraction of Inspired Oxygen  
 364 (FiO2), which are crucial parameters in ventilator settings.

### 362 5.3 Can CRL with Offline CT Learn Safe Policies?

363 We have confirmed the existence  
 364 of two unsafe strategy issues,  
 365 namely “too high” and “sudden  
 366 change” in the treatment of sep-  
 367 sis, particularly in vaso in Sec-  
 368 tion 1. To validate whether  
 369 the CRL+CT approach could ad-  
 370 dress these concerns, we employ  
 371 the same statistical methods to

Table 3: The proportion of “too high” and “sudden change” oc-  
 currences in drug dosage recommended by RL methods.

| Drug dosage<br>( $\mu\text{g}/(\text{kg} \cdot \text{min})$ ) | Physician | DDPG   | No cost | CDT                           |                               |
|---|-----------|--------|---------|-------------------------------|-------------------------------|
|   |           |        |         | Custom cost                   | CT                            |
| vaso > 0.75   | 2.27%     | 7.44%  | 0.13%   | 0% ↓<br>(max = 0.00)          | 0% ↓<br>(max = 0.11)          |
| vaso > 0.9  | 1.71%     | 7.40%  | 0.09%   |                               |                               |
| $\Delta$ vaso > 0.75  | 2.45%     | 21.00% | 0.64%   | 0% ↓<br>(max $\Delta$ = 0.00) | 0% ↓<br>(max $\Delta$ = 0.10) |
| $\Delta$ vaso > 0.9   | 1.88%     | 20.62% | 0.48%   |                               |                               |

372 evaluate our methodology, shown in Table 3. To elucidate the efficacy of CT, we compare it with  
 373 CDT+No-cost and CDT+Custom-cost approaches. We find that only the custom cost and CT methods  
 374 successfully mitigated the risks associated with “too high” and “sudden change” behaviors. However,  
 375 the custom cost approach opts to avoid administering drugs to mitigate these risks. Without these  
 376 drugs, the patient’s condition may not be alleviated, potentially leading to patient mortality. The  
 377 CDT+CT approach can give a more appropriate drug dosage.

378 **Ablation Study.** To investigate the impact of each component on the model’s performance, we  
 379 conducted experiments by sequentially removing each component from the CDT+CT model. The  
 380 results are presented in the lower half of Table 2. Both CT and its non-Markovian layer (attention  
 381 layer) are indispensable and crucial components; removing either one results in a decrease in perfor-  
 382 mance. Additionally, we observed that even a pure generative model outperforms DDPG in terms  
 383 of performance. This is primarily because it inherently operates as a sequence-based reinforcement  
 384 learning model, possessing exploration and consideration for long-term history. Therefore, this  
 385 further underscores the effectiveness of sequence-based approaches in healthcare applications.

## 386 6 Conclusion

387 In this paper, we propose offline CT, a novel ICRL algorithm designed to address safety issues  
 388 in healthcare. This method utilizes a causal attention mechanism to observe patients’ historical  
 389 information, similar to the approach taken by actual doctors and employs non-Markovian importance  
 390 weights to effectively capture critical states. To achieve offline learning, we introduce a model-based  
 391 offline RL for exploratory data augmentation to discover unsafe decisions and train CT. Experiments  
 392 in sepsis and mechanical ventilation demonstrate that our method avoids risky behaviors while  
 393 achieving strategies that closely approximate the lowest mortality rates.

394 **Limitations.** There are also several limitations of offline CT: (1) Lack of rigorous theoretical analysis:  
 395 We did not precisely define the types of constraint sets, thereby conducting rigorous theoretical  
 396 analysis on constraint sets remains challenging; (2) Need for more computational resources: Due  
 397 to the Transformer architecture, more computational resources are required; (3) Fewer evaluation  
 398 metrics: There is a lack of more medical-specific evaluation metrics in the experimental evaluation  
 399 section; (4) Unrealistic assumptions of expert demonstrations: we assume that expert demonstrations  
 400 are optimal in both constraint satisfaction and reward maximization. However, in reality, this  
 401 assumption may not always hold. Therefore, researching a more effective approach to address the  
 402 aforementioned issues holds promise for the field of secure medical reinforcement learning.

403 **References**

- 404 [1] Stephen Petterson, Robert McNellis, Kathleen Klink, David Meyers, and Andrew Bazemore.  
405 The state of primary care in the united states: A chartbook of facts and statistics. *Washington,*  
406 *DC: Robert Graham Center*, 2018.
- 407 [2] Estevão Bassi, Marcelo Park, Luciano Cesar Pontes Azevedo, et al. Therapeutic strategies for  
408 high-dose vasopressor-dependent shock. *Critical care research and practice*, 2013, 2013.
- 409 [3] Thomas Auchet, Marie-Alix Regnier, Nicolas Girerd, and Bruno Levy. Outcome of patients  
410 with septic shock and high-dose vasopressor therapy. *Annals of Intensive Care*, 7:1–9, 2017.
- 411 [4] Davide Tommaso Andreis and Mervyn Singer. Catecholamines for inflammatory shock: a  
412 jekyll-and-hyde conundrum. *Intensive care medicine*, 42:1387–1397, 2016.
- 413 [5] Yan Jia, John Burden, Tom Lawton, and Ibrahim Habli. Safe reinforcement learning for sepsis  
414 treatment. In *2020 IEEE International conference on healthcare informatics (ICHI)*, pages 1–7.  
415 IEEE, 2020.
- 416 [6] Rui Shi, Olfa Hamzaoui, Nello De Vita, Xavier Monnet, and Jean-Louis Teboul. Vasopressors  
417 in septic shock: which, when, and how much? *Annals of Translational Medicine*, 8(12), 2020.
- 418 [7] Claude Martin, Sophie Medam, Francois Antonini, Julie Alingrin, Malik Haddam, Emmanuelle  
419 Hammad, Bertrand Meyssignac, Coralie Vigne, Laurent Zieleskiewicz, and Marc Leone. Nore-  
420 pinephrine: not too much, too long. *Shock*, 44(4):305–309, 2015.
- 421 [8] Kristin Lavigne Fadale, Denise Tucker, Jennifer Dungan, and Valerie Sabol. Improving nurses’  
422 vasopressor titration skills and self-efficacy via simulation-based learning. *Clinical Simulation*  
423 *in Nursing*, 10(6):e291–e299, 2014.
- 424 [9] Yong Huang, Rui Cao, and Amir Rahmani. Reinforcement learning for sepsis treatment: A  
425 continuous action space solution. In *Machine Learning for Healthcare Conference*, pages  
426 631–647. PMLR, 2022.
- 427 [10] Yongshuai Liu, Avishai Halev, and Xin Liu. Policy learning with constraints in model-free  
428 reinforcement learning: A survey. In *The 30th International Joint Conference on Artificial*  
429 *Intelligence (IJCAI)*, 2021.
- 430 [11] Guiliang Liu, Yudong Luo, Ashish Gaurav, Kasra Rezaee, and Pascal Poupart. Benchmarking  
431 constraint inference in inverse reinforcement learning. *arXiv preprint arXiv:2206.09670*, 2022.
- 432 [12] Shehryar Malik, Usman Anwar, Alireza Aghasi, and Ali Ahmed. Inverse constrained reinforce-  
433 ment learning. In *International conference on machine learning*, pages 7390–7399. PMLR,  
434 2021.
- 435 [13] Masaaki Kijima. *Markov processes for stochastic modeling*. Springer, 2013.
- 436 [14] Zhiyue Zhang, Hongyuan Mei, and Yanxun Xu. Continuous-time decision transformer for  
437 healthcare applications. In *International Conference on Artificial Intelligence and Statistics*,  
438 pages 6245–6262. PMLR, 2023.
- 439 [15] Catherine Plaisant, Brett Milash, Anne Rose, Seth Widoff, and Ben Shneiderman. Lifelines:  
440 visualizing personal histories. In *Proceedings of the SIGCHI conference on Human factors in*  
441 *computing systems*, pages 221–227, 1996.
- 442 [16] Ashish Gaurav, Kasra Rezaee, Guiliang Liu, and Pascal Poupart. Learning soft constraints from  
443 constrained expert demonstrations. *arXiv preprint arXiv:2206.01311*, 2022.
- 444 [17] Qinqing Zheng, Amy Zhang, and Aditya Grover. Online decision transformer. In *international*  
445 *conference on machine learning*, pages 27042–27059. PMLR, 2022.
- 446 [18] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter  
447 Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning  
448 via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097,  
449 2021.

- 450 [19] Changyeon Kim, Jongjin Park, Jinwoo Shin, Honglak Lee, Pieter Abbeel, and Kimin Lee.  
451 Preference transformer: Modeling human preferences using transformers for rl. *arXiv preprint*  
452 *arXiv:2303.00957*, 2023.
- 453 [20] Aniruddh Raghu, Matthieu Komorowski, Imran Ahmed, Leo Celi, Peter Szolovits, and Marzyeh  
454 Ghassemi. Deep reinforcement learning for sepsis treatment. *arXiv preprint arXiv:1711.09602*,  
455 2017.
- 456 [21] Xuefeng Peng, Yi Ding, David Wihl, Omer Gottesman, Matthieu Komorowski, H Lehman Li-  
457 wei, Andrew Ross, Aldo Faisal, and Finale Doshi-Velez. Improving sepsis treatment strategies  
458 by combining deep and kernel-based reinforcement learning. In *AMIA Annual Symposium*  
459 *Proceedings*, volume 2018, page 887. American Medical Informatics Association, 2018.
- 460 [22] Thanh Cong Do, Hyung Jeong Yang, Seok Bong Yoo, and In-Jae Oh. Combining reinforcement  
461 learning with supervised learning for sepsis treatment. In *The 9th International Conference on*  
462 *Smart Media and Applications*, pages 219–223, 2020.
- 463 [23] Flemming Kondrup, Thomas Jiralerspong, Elaine Lau, Nathan de Lara, Jacob Shkrob, My Duc  
464 Tran, Doina Precup, and Sumana Basu. Towards safe mechanical ventilation treatment using  
465 deep offline reinforcement learning. In *Proceedings of the AAAI Conference on Artificial*  
466 *Intelligence*, volume 37, pages 15696–15702, 2023.
- 467 [24] Wei Gong, Linxiao Cao, Yifei Zhu, Fang Zuo, Xin He, and Haoquan Zhou. Federated inverse  
468 reinforcement learning for smart icus with differential privacy. *IEEE Internet of Things Journal*,  
469 2023.
- 470 [25] Chao Yu, Guoqi Ren, and Yinzhaodong. Supervised-actor-critic reinforcement learning for  
471 intelligent mechanical ventilation and sedative dosing in intensive care units. *BMC medical*  
472 *informatics and decision making*, 20:1–8, 2020.
- 473 [26] Niloufar Eghbali, Tuka Alhanai, and Mohammad M Ghassemi. Patient-specific sedation  
474 management via deep reinforcement learning. *Frontiers in Digital Health*, 3:608893, 2021.
- 475 [27] Giulia Calvi, Eleonora Manzoni, and Mirco Rampazzo. Reinforcement q-learning for closed-  
476 loop hypnosis depth control in anesthesia. In *2022 30th Mediterranean Conference on Control*  
477 *and Automation (MED)*, pages 164–169. IEEE, 2022.
- 478 [28] Gabriel Schamberg, Marcus Badgeley, Benyamin Meschede-Krasa, Ohyeon Kwon, and  
479 Emery N Brown. Continuous action deep reinforcement learning for propofol dosing dur-  
480 ing general anesthesia. *Artificial Intelligence in Medicine*, 123:102227, 2022.
- 481 [29] Glen Chou, Dmitry Berenson, and Necmiye Ozay. Learning constraints from demonstrations. In  
482 *Algorithmic Foundations of Robotics XIII: Proceedings of the 13th Workshop on the Algorithmic*  
483 *Foundations of Robotics 13*, pages 228–245. Springer, 2020.
- 484 [30] Daehyung Park, Michael Noseworthy, Rohan Paul, Subhro Roy, and Nicholas Roy. Inferring  
485 task goals and constraints using bayesian nonparametric inverse reinforcement learning. In  
486 *Conference on robot learning*, pages 1005–1014. PMLR, 2020.
- 487 [31] Dexter RR Scobee and S Shankar Sastry. Maximum likelihood constraint inference for inverse  
488 reinforcement learning. *arXiv preprint arXiv:1909.05477*, 2019.
- 489 [32] David L McPherson, Kaylene C Stocking, and S Shankar Sastry. Maximum likelihood constraint  
490 inference from stochastic demonstrations. In *2021 IEEE Conference on Control Technology*  
491 *and Applications (CCTA)*, pages 1208–1213. IEEE, 2021.
- 492 [33] Mattijs Baert, Pietro Mazzaglia, Sam Leroux, and Pieter Simoens. Maximum causal entropy  
493 inverse constrained reinforcement learning. *arXiv preprint arXiv:2305.02857*, 2023.
- 494 [34] Brian D Ziebart, J Andrew Bagnell, and Anind K Dey. Modeling interaction via the principle of  
495 maximum causal entropy. 2010.
- 496 [35] Guanren Qiao, Guiliang Liu, Pascal Poupart, and Zhiqiang Xu. Multi-modal inverse constrained  
497 reinforcement learning from a mixture of demonstrations. *Advances in Neural Information*  
498 *Processing Systems*, 36, 2024.

- 499 [36] Sheng Xu and Guiliang Liu. Uncertainty-aware constraint inference in inverse constrained  
500 reinforcement learning. In *The Twelfth International Conference on Learning Representations*,  
501 2023.
- 502 [37] Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big  
503 sequence modeling problem. *Advances in neural information processing systems*, 34:1273–  
504 1286, 2021.
- 505 [38] Zuxin Liu, Zijian Guo, Yihang Yao, Zhepeng Cen, Wenhao Yu, Tingnan Zhang, and Ding  
506 Zhao. Constrained decision transformer for offline safe reinforcement learning. *arXiv preprint*  
507 *arXiv:2302.07351*, 2023.
- 508 [39] Eitan Altman. Constrained markov decision processes with total cost criteria: Lagrangian  
509 approach and dual linear program. *Mathematical methods of operations research*, 48:387–417,  
510 1998.
- 511 [40] Aniruddh Raghu, Matthieu Komorowski, Leo Anthony Celi, Peter Szolovits, and Marzyeh  
512 Ghassemi. Continuous state-space models for optimal sepsis treatment: a deep reinforcement  
513 learning approach. In *Machine Learning for Healthcare Conference*, pages 147–163. PMLR,  
514 2017.
- 515 [41] Mehdi Fatemi, Taylor W Killian, Jayakumar Subramanian, and Marzyeh Ghassemi. Medi-  
516 cal dead-ends and learning to identify high-risk states and treatments. *Advances in Neural*  
517 *Information Processing Systems*, 34:4856–4870, 2021.
- 518 [42] Matthieu Komorowski, Leo A Celi, Omar Badawi, Anthony C Gordon, and A Aldo Faisal. The  
519 artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care.  
520 *Nature medicine*, 24(11):1716–1720, 2018.
- 521 [43] Jongmin Lee, Cosmin Paduraru, Daniel J Mankowitz, Nicolas Heess, Doina Precup, Kee-Eung  
522 Kim, and Arthur Guez. Coptidice: Offline constrained reinforcement learning via stationary  
523 distribution correction estimation. *arXiv preprint arXiv:2204.08957*, 2022.
- 524 [44] Jiayi Guan, Guang Chen, Jiaming Ji, Long Yang, Zhijun Li, et al. Voce: Variational optimization  
525 with conservative estimation for offline safe reinforcement learning. *Advances in Neural*  
526 *Information Processing Systems*, 36, 2024.
- 527 [45] Yongju Kim, Hyung Keun Park, Jaimyun Jung, Peyman Asghari-Rad, Seungchul Lee, Jin You  
528 Kim, Hwan Gyo Jung, and Hyoung Seop Kim. Exploration of optimal microstructure and  
529 mechanical properties in continuous microstructure space using a variational autoencoder.  
530 *Materials & Design*, 202:109544, 2021.
- 531 [46] Tim Hsu, William K Epting, Hokon Kim, Harry W Abernathy, Gregory A Hackett, Anthony D  
532 Rollett, Paul A Salvador, and Elizabeth A Holm. Microstructure generation via generative  
533 adversarial network for heterogeneous, topologically complex 3d materials. *Jom*, 73:90–102,  
534 2021.
- 535 [47] Akshay Iyer, Biswadip Dey, Arindam Dasgupta, Wei Chen, and Amit Chakraborty. A conditional  
536 generative model for predicting material microstructures from processing methods. *arXiv*  
537 *preprint arXiv:1910.02133*, 2019.
- 538 [48] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion  
539 models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,  
540 2023.
- 541 [49] Christian Dürer, Paul Seibert, Dennis Rucker, Stephanie Handford, Markus Kästner, and  
542 Maik Gude. Conditional diffusion-based microstructure reconstruction. *Materials Today*  
543 *Communications*, 35:105608, 2023.
- 544 [50] Arne Peine, Ahmed Hallawa, Johannes Bickenbach, Guido Dartmann, Lejla Begic Fazlic, Anke  
545 Schmeink, Gerd Ascheid, Christoph Thiemermann, Andreas Schuppert, Ryan Kindle, et al.  
546 Development and validation of a reinforcement learning algorithm to dynamically optimize  
547 mechanical ventilation in critical care. *NPJ digital medicine*, 4(1):32, 2021.

- 548 [51] Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad  
549 Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. MIMIC-III,  
550 a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.
- 551 [52] Yonglin Li, Chunjiang Yan, Ziyang Gan, Xiaotu Xi, Zhanpeng Tan, Jun Li, and Guowei Li.  
552 Prognostic values of sofa score, qsofa score, and lods score for patients with sepsis. *Annals of*  
553 *palliative medicine*, 9(3):1037044–1031044, 2020.
- 554 [53] Seung Mok Ryoo, JungBok Lee, Yoon-Seon Lee, Jae Ho Lee, Kyoung Soo Lim, Jin Won  
555 Huh, Sang-Bum Hong, Chae-Man Lim, Younsuck Koh, and Won Young Kim. Lactate level  
556 versus lactate clearance for predicting mortality in patients with septic shock defined by sepsis-3.  
557 *Critical care medicine*, 46(6):e489–e495, 2018.
- 558 [54] Nishant Raj Pandey, Yu-yao Bian, and Song-tao Shou. Significance of blood pressure variability  
559 in patients with sepsis. *World journal of emergency medicine*, 5(1):42, 2014.
- 560 [55] Marta Carrara, Bernardo Bollen Pinto, Giuseppe Baselli, Karim Bendjelid, and Manuela Ferrario.  
561 Baroreflex sensitivity and blood pressure variability can help in understanding the different  
562 response to therapy during acute phase of septic shock. *Shock*, 50(1):78–86, 2018.
- 563 [56] Hoang Le, Cameron Voloshin, and Yisong Yue. Batch policy learning under constraints. In  
564 *International Conference on Machine Learning*, pages 3703–3712. PMLR, 2019.
- 565 [57] Mervyn Singer, Clifford S Deutschman, Christopher Warren Seymour, Manu Shankar-Hari,  
566 Djillali Annane, Michael Bauer, Rinaldo Bellomo, Gordon R Bernard, Jean-Daniel Chiche,  
567 Craig M Coopersmith, et al. The third international consensus definitions for sepsis and septic  
568 shock (sepsis-3). *Jama*, 315(8):801–810, 2016.
- 569 [58] Flavio Lopes Ferreira, Daliana Peres Bota, Annette Bross, Christian Mélot, and Jean-Louis  
570 Vincent. Serial evaluation of the sofa score to predict outcome in critically ill patients. *Jama*,  
571 286(14):1754–1758, 2001.

## 572 A Design and Analysis of the Custom Constraint Function

573 We base our design on prior knowledge that intravenous (IV) intake exceeding  $2000\text{mL}/4\text{h}$  or  
 574 vasopressor (Vaso) dosage surpassing  $1\text{g}/(\text{kg}\cdot\text{min})$  is generally considered unsafe in sepsis treatment  
 575 [6]. To design a reasonable constraint function, we refer to the constraint function designed by Liu *et*  
 576 *al.* in the Bullet safety gym environments[38]. We define the cost function as shown in Equation 9.  
 577 Thus, during the treatment of sepsis, if the agent exceeds the maximum dosage thresholds of the two  
 578 medications, it incurs a cost due to constraint violation.

$$c(s, a) = \mathbf{1}(a_{IV} > a_{IV \max}) + \mathbf{1}(a_{Vaso} > a_{Vaso \max}) \quad (9)$$

579 where,  $s$  and  $a$  represent the patient’s state and action, respectively.  $a_{IV \max} = 2000$  indicates that  
 580 the maximum fluid intake through IV is  $2000\text{mL}$ , and  $a_{Vaso \max} = 1$  signifies that the maximum  
 581 Vaso dosage is  $1\mu\text{g}/(\text{kg}\cdot\text{min})$ .

582 We applied our custom constraint function in the CDT [38] method, and the results are shown in  
 583 Figure 9. Compared to the Vaso dosage recommended by doctors, our strategy exhibits excessive  
 584 suppression of the Vaso. The maximum dosage of Vaso is  $0.0011\mu\text{g}/(\text{kg}\cdot\text{min})$ , which is minimal  
 585 and insufficient to provide the patient with effective therapeutic effects.

586 Therefore, Equation 9 is not suitable. The primary issues may include uniform constraint strength  
 587 for excessive drug dosages, for instance, the cost for IV exceeding  $2000\text{ mL}$  and IV exceeding  
 588  $3000\text{ mL}$  is the same at 1; lack of generalization, where the constraint cost does not vary with the  
 589 patient’s tolerance. If a patient has an intolerance to VASO, the maximum value for VASO maybe 0,  
 590 which cannot be captured by the self-imposed constraint function. Moreover, it lacks generalization,  
 591 requiring redesign of the constraint function when addressing other unsafe medical issues; and it’s  
 essential to ensure the correctness of the underlying medical knowledge premises.

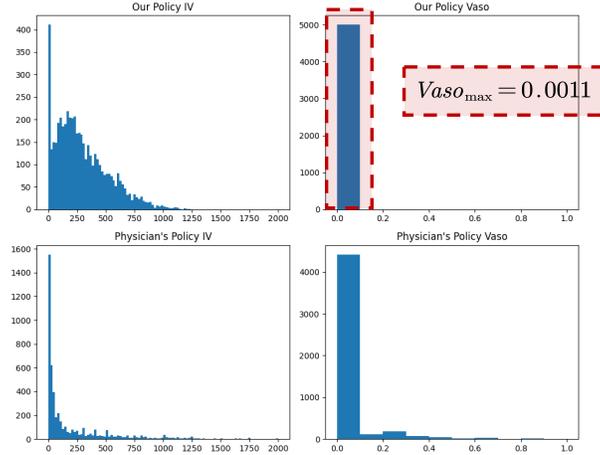


Figure 9: Drug dosage distribution under custom constraint functions in sepsis.

592

## 593 B Experiment Supplement

### 594 B.1 Sepsis Problem Define

595 Our definition is similar to [40]. We extract data from adult patients meeting the criteria for sepsis-3  
 596 criteria [57] and collect their data within the first 72 hours of admission.

597 **State Space.** We use a 4-hour window and select 48 patient indicators as the state for a one-time unit  
 598 of the patient. The state indicators include Demographics/Static, Lab Values, Vital Signs, and Intake  
 599 and Output Events, detailed as follows [40]:

- 600 • Demographics/Static: Shock Index, Elixhauser, SIRS, Gender, Re-admission, GCS - Glas-  
 601 gow Coma Scale, SOFA - Sequential Organ Failure Assessment, Age

- 602 • Lab Values Albumin: Arterial pH, Calcium, Glucose, Hemoglobin, Magnesium, PTT -  
603 Partial Thromboplastin Time, Potassium, SGPT - Serum Glutamic-Pyruvic Transaminase,  
604 Arterial Blood Gas, BUN Blood Urea Nitrogen, Chloride, Bicarbonate, INR - International  
605 Normalized Ratio, Sodium, Arterial Lactate, CO2, Creatinine, Ionised Calcium, PT - Pro-  
606 thrombin Time, Platelets Count, SGOT Serum Glutamic-Oxaloacetic Transaminase, Total  
607 bilirubin, White Blood Cell Count
- 608 • Vital Signs: Diastolic Blood Pressure, Systolic Blood Pressure, Mean Blood Pressure,  
609 PaCO2, PaO2, FiO2, PaO/FiO2 ratio, Respiratory Rate, Temperature (Celsius), Weight (kg),  
610 Heart Rate, SpO2
- 611 • Intake and Output Events: Fluid Output - 4 hourly period, Total Fluid Output, Mechanical  
612 Ventilation

613 **Action Space.** Regarding the treatment of sepsis, there are two main types of medications: in-  
614 travenous fluids and vasopressors. We select the total amount of intravenous fluids for each time  
615 unit and the maximum dose of vasopressors as the two dimensions of the action space, defined as  
616 (sum(IV), max(Vaso)). Each dimension is a continuous value greater than 0.

617 **Reward Function.** We refer to the reward function used in [9], as shown in the following equation:

$$r(s_t, s_{t+1}) = \lambda_1 \tanh(s_t^{\text{SOFA}} - 6) + \lambda_2 (s_{t+1}^{\text{SOFA}} - s_t^{\text{SOFA}}) \quad (10)$$

618 Where  $\lambda_0$  and  $\lambda_1$  are hyperparameters set to  $-0.25$  and  $-0.2$ , respectively. This reward function is  
619 designed based on the SOFA score, as it is a key indicator of the health status for sepsis patients and  
620 widely used in clinical settings. The formula describes a penalty when the SOFA score increases and  
621 a reward when the SOFA score decreases. We set 6 as the cutoff value because the mortality rate  
622 sharply increases when the SOFA score exceeds 6 [58].

## 623 B.2 Mechanical Ventilation Treatment Problem Define

624 The RL problem definition for Mechanical Ventilation Treatment is referenced from [23].

### 625 State Space.

- 626 • Demographics/Static: Elixhauser, SIRS, Gender, Re-admission, GCS, SOFA, Age
- 627 • Lab Values Albumin: Arterial pH, Glucose, Hemoglobin, Magnesium, PTT, BUN Blood  
628 Urea Nitrogen, Chloride, Bicarbonate, INR, Sodium, Arterial Lactate, CO2, Creatinine,  
629 Ionised Calcium, PT, Platelets Count, White Blood Cell Count, Hb
- 630 • Vital Signs: Diastolic Blood Pressure, Systolic Blood Pressure, Mean Blood Pressure,  
631 Temperature, Weight (kg), Heart Rate, SpO2
- 632 • Intake and Output Events: Urine output, vasopressors, intravenous fluids, cumulative fluid  
633 balance

634 **Action Space.** The action space mainly consists of Positive End Expiratory Pressure (PEEP) and  
635 Fraction of Inspired Oxygen (FiO2), which are crucial parameters in ventilator settings. Here, we  
636 consider a discrete space configuration, with each parameter divided into 7 intervals. Therefore, our  
action space is  $7 \times 7$ , depicted as 4.

Table 4: The action space of the mechanical ventilator.

| Action              | 0     | 1     | 2     | 3     | 4     | 5     | 6   |
|---------------------|-------|-------|-------|-------|-------|-------|-----|
| PEEP(cmH20)         | 0-5   | 5-7   | 7-9   | 9-11  | 11-13 | 13-15 | >15 |
| FiO2(Percentage(%)) | 25-30 | 30-35 | 35-40 | 40-45 | 45-50 | 50-55 | >55 |

637

638 **Reward Function.** The primary objective of setting respiratory parameters is to ensure the patient’s  
639 survival. We adopt the same reward function design as the work [23], defined as Equation 11. This  
640 reward function first considers the terminal reward: if the patient dies, the reward  $r$  is set to  $-1$ ;  
641 otherwise, it is  $+1$  in the terminal state. Additionally, to provide more frequent rewards, intermediate  
642 rewards are considered. Intermediate rewards mainly focus on the Apache II score, which evaluates

643 various parameters to describe the patient’s health status. This reward function utilizes the increase or  
 644 decrease in this score to reward the agent.

$$r(s_t, a_t, s_{t+1}) = \begin{cases} +1 & \text{if } t = T \text{ and } m_t = 1 \\ -1 & \text{if } t = T \text{ and } m_t = 0 \\ \frac{(A_{t+1} - A_t)}{\max_A - \min_A} & \text{otherwise} \end{cases} \quad (11)$$

645 In Equation 11,  $T$  represents the length of the patient’s trajectory,  $m$  indicates whether the patient  
 646 ultimately dies,  $A$  denotes the Apache II score, and  $\max_A$  and  $\min_A$  respectively denote the maximum  
 647 and minimum values.

### 648 B.3 The Evaluation of Model-based Offline RL

649 **Generating data within a reasonable range.** To validate model-based offline RL, we first check  
 650 whether the values it produces fall within the legal range. The results are depicted in Figure 10. After  
 651 analyzing the generated data, we find that the majority of state values have a probability of over 99%  
 652 of being within the legal range. A few values related to gender and re-admission range between 60%  
 653 and 70%. This could be due to these two indicators having limited correlation with other metrics,  
 making them more challenging for the model to assess.

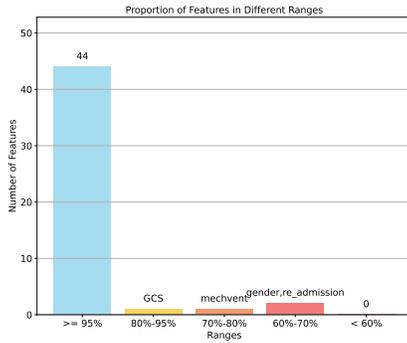


Figure 10: The accuracy of predicting different state values within the legal range.

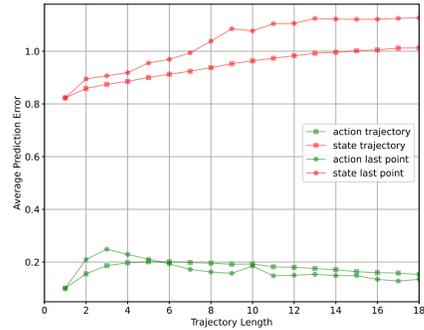


Figure 11: The relationship between average prediction error and trajectory length.

654

655 **Generating violating data.** In addition, we evaluate the violating actions generated by the model, as  
 656 shown in Figure 12. When compared with expert strategies and penalty distributions, we find that the  
 657 actions generated by the model mostly fall within the legal range. However, it occasionally produces  
 658 behaviors that are inappropriate for the current state, constituting violating data. This indicates that  
 our generative model can produce legally violating data.

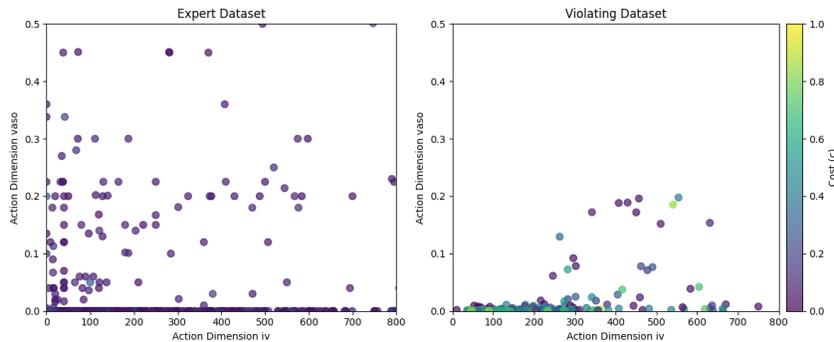


Figure 12: The distribution and penalty values of violating data and expert data.

659

660 **The length of a trajectory.** Regarding the selection of trajectory length, we consider the relationship  
 661 between the average prediction error, the error of the last point in the trajectory, and the trajectory

662 length. We use the model-based offline RL to generate trajectories and compare them with expert  
 663 data using the Euclidean distance to measure their differences. We evaluate the average error and  
 664 the error of the last point in the trajectory, as shown in Figure 11. We observe that with an increase  
 665 in trajectory length, the average prediction error at each time step decreases, while the state error  
 666 stabilizes. Taking into account the observation length and prediction accuracy, we ultimately choose  
 667 to generate trajectories with lengths ranging from 10 to 15.

#### 668 B.4 The Evaluation of Cost function in Sepsis

669 To validate that the CT method captures key states, we conduct statistical analysis on the relationship  
 670 between state values and penalty values. We collect penalty values under different state values  
 671 for all patients, and the complete information is shown in Figure 13. We find that the CT method  
 672 successfully captures unsafe states and imposes higher penalties accordingly. The safe range of state  
 673 values is shown in Table 5.

674 To validate the role of the attention layer in capturing states in CT, we conducted tests, and the  
 675 experimental results are presented in Figure 14 and 13. We found that the attention layer plays a  
 676 crucial role in state capture. For instance, in the case of an increase in the SOFA score, without the  
 677 attention layer, this increase cannot be captured, while with the attention layer, it clearly captures the  
 678 change. Thus, this indicates that SOFA, as a key diagnostic indicator of sepsis, with the help of the  
 attention layer, CT can accurately capture its changes.

Table 5: State indicators and their normal ranges.

| Indicator        | Safe Range | Indicator       | Safe Range | Indicator   | Safe Range |
|------------------|------------|-----------------|------------|-------------|------------|
| Albumin          | 3.5~5.1    | HCO3            | 25~40      | SGOT        | 0~40       |
| Arterial_BE      | -3~+3      | Glucose         | 70~140     | SGPT        | 0~40       |
| Arterial_lactate | 0.5~1.7    | HR              | 60~100     | SIRS        | ↓          |
| Arterial_PH      | 7.35~7.45  | Hb              | 12~16      | SOFA        | ↓          |
| BUN              | 7~22       | INR             | 0.8~1.5    | Shock_Index | ↓          |
| CO2_mEqL         | 20~34      | MeanBP          | 70~100     | Sodium      | 135~145    |
| Calcium          | 8.6~10.6   | PT              | 11~13      | SpO2        | 95~99      |
| Chloride         | 96~106     | PTT             | 23~37      | SysBP       | 90~139     |
| Creatinine       | 0.5~1.5    | PaO2_FiO2       | 400~500    | Temp_C      | 36.0~37.0  |
| DiaBP            | 60~89      | Platelets_count | 125~350    | WBC_count   | 4~10       |
| FiO2             | 0.5~0.6    | Potassium       | 4.1~5.6    | PaCO2       | 35~45      |
| GCS              | ↑          | RR              | 12~20      | PaO2        | 80~100     |

↑ indicates higher values are more normal, while ↓ indicates lower values are more normal.  
 The maximum value for GCS is 15. The minimum value for SIRS, SOFA, and Shock\_Index is 0.

679

#### 680 B.5 Experimental Settings

681 To train the CRL+CT model, we use a total of 3 NVIDIA GeForce RTX 3090 GPUs, each with  
 682 24GB of memory. Training a CRL+CT model typically takes 5-6 hours. We employ 5 random seeds  
 683 for validation. We use the Adam optimization algorithm to optimize all our networks, updating the  
 684 learning rate using a decay factor parameterization at each iteration. The main hyperparameters are  
 685 summarized in Table 6 and 7.

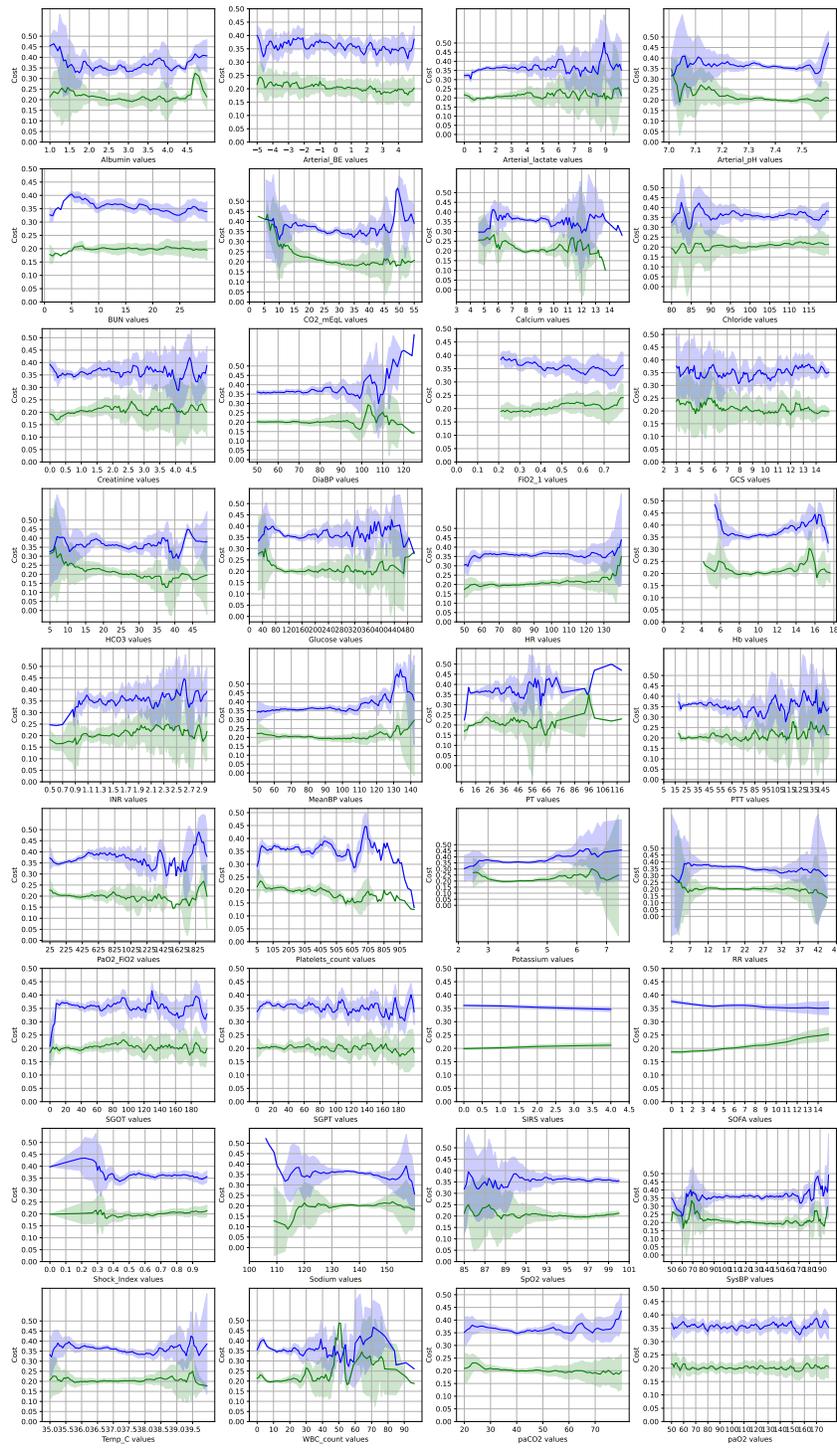


Figure 13: The relationship between all states and cost values

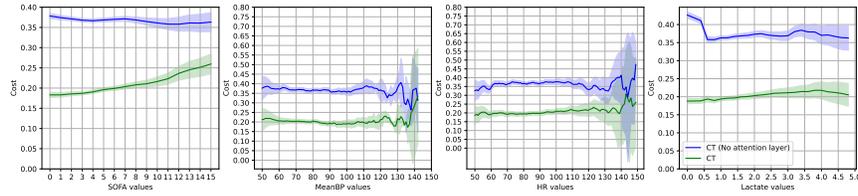


Figure 14: The performance contrast between CT with and without an attention layer. The blue line represents the absence of an attention layer, while the green line indicates the presence of an attention layer.

Table 6: List of the utilized hyperparameters in CT.

| Offline CT Parameters | values |
|-----------------------|--------|
| Genetivate Model      |        |
| Embedding_dim         | 128    |
| Layer                 | 3      |
| Head                  | 8      |
| Learning rate         | 1e-4   |
| Pre-train steps       | 5000   |
| Batch size            | 256    |
| CT                    |        |
| Embedding_dim         | 64     |
| Layer                 | 3      |
| Head                  | 1      |
| Learning rate         | 1e-6   |
| Update steps          | 30000  |
| Batch size            | 512    |
| CDT                   |        |
| Learning rate         | 1e-4   |
| Embedding_dim         | 128    |
| Layers                | 3      |
| Heads                 | 8      |
| Update steps          | 60000  |

Table 7: List of the utilized hyperparameters in CRL.

| Parameters                     | Sepsis  | Parameters                     | Mechanical Ventilation |
|--------------------------------|---------|--------------------------------|------------------------|
| General                        |         | General                        |                        |
| Expert data patient number     | 14313   | Expert data patient number     | 13846                  |
| Validation data patient number | 6275    | Validation data patient number | 5954                   |
| Max Length                     | 10      | Max Length                     | 10                     |
| Action_dim                     | 2       | Action_dim                     | 2                      |
| State_dim                      | 48      | State_dim                      | 36                     |
| Gamma                          | 0.99    | Gamma                          | 0.99                   |
| DDPG                           |         | DDQN                           |                        |
| Learning rate                  | 1e-3    | Learning rate                  | 1e-4                   |
| Policy Network                 | 256,256 | Policy Network                 | 64,64                  |
| Replay memory size             | 20000   | Update steps                   | 500000                 |
| Update steps                   | 20000   |                                |                        |
| VOCE                           |         | CQL                            |                        |
| Learning rate                  | 1e-3    | Learning rate                  | 1e-4                   |
| Policy Network                 | 256,256 | Policy Network                 | 64,64                  |
| Alpha scale                    | 10      | Update steps                   | 500000                 |
| KL constraint                  | 0.01    | Alphas                         | 0.05,0.1,0.5,1,2       |
| Dual constraint                | 0.1     |                                |                        |
| Update steps                   | 4000    |                                |                        |
| CoplDICE                       |         |                                |                        |
| Learning rate                  | 1e-4    |                                |                        |
| Policy Network                 | 256,256 |                                |                        |
| Alpha                          | 0.5     |                                |                        |
| Cost limit                     | 10      |                                |                        |
| Update steps                   | 100000  |                                |                        |
| BCQ-Lag                        |         |                                |                        |
| Learning rate                  | 1e-3    |                                |                        |
| Policy Network                 | 256,256 |                                |                        |
| Cost limit                     | 10      |                                |                        |
| Lambda                         | 0.75    |                                |                        |
| Beta                           | 0.5     |                                |                        |
| Update steps                   | 100000  |                                |                        |

## 686 **NeurIPS Paper Checklist**

### 687 **1. Claims**

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

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

691 Justification: In the abstract and introduction, we delineate the main motivations and  
692 contributions of this paper and its application in the field of safe reinforcement learning in  
693 healthcare.

694 Guidelines:

- 695 • The answer NA means that the abstract and introduction do not include the claims  
696 made in the paper.
- 697 • The abstract and/or introduction should clearly state the claims made, including the  
698 contributions made in the paper and important assumptions and limitations. A No or  
699 NA answer to this question will not be perceived well by the reviewers.
- 700 • The claims made should match theoretical and experimental results, and reflect how  
701 much the results can be expected to generalize to other settings.
- 702 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
703 are not attained by the paper.

### 704 **2. Limitations**

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

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

707 Justification: In the final section, this paper discusses the limitations of the method.

708 Guidelines:

- 709 • The answer NA means that the paper has no limitation while the answer No means that  
710 the paper has limitations, but those are not discussed in the paper.
- 711 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 712 • The paper should point out any strong assumptions and how robust the results are to  
713 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
714 model well-specification, asymptotic approximations only holding locally). The authors  
715 should reflect on how these assumptions might be violated in practice and what the  
716 implications would be.
- 717 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
718 only tested on a few datasets or with a few runs. In general, empirical results often  
719 depend on implicit assumptions, which should be articulated.
- 720 • The authors should reflect on the factors that influence the performance of the approach.  
721 For example, a facial recognition algorithm may perform poorly when image resolution  
722 is low or images are taken in low lighting. Or a speech-to-text system might not be  
723 used reliably to provide closed captions for online lectures because it fails to handle  
724 technical jargon.
- 725 • The authors should discuss the computational efficiency of the proposed algorithms  
726 and how they scale with dataset size.
- 727 • If applicable, the authors should discuss possible limitations of their approach to  
728 address problems of privacy and fairness.
- 729 • While the authors might fear that complete honesty about limitations might be used by  
730 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
731 limitations that aren't acknowledged in the paper. The authors should use their best  
732 judgment and recognize that individual actions in favor of transparency play an impor-  
733 tant role in developing norms that preserve the integrity of the community. Reviewers  
734 will be specifically instructed to not penalize honesty concerning limitations.

### 735 **3. Theory Assumptions and Proofs**

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

738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791

Answer: [Yes]

Justification: We have documented the relevant theories and assumptions in the paper or supplementary materials.

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: Our approach is reproducible, and our code can be made publicly available after the paper is published, including the relevant data processing procedures.

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

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

795 Answer: [Yes]

796 Justification: Our code can be made publicly available after the paper is published.

797 Guidelines:

- 798 • The answer NA means that paper does not include experiments requiring code.
- 799 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
800 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 801 • While we encourage the release of code and data, we understand that this might not be  
802 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
803 including code, unless this is central to the contribution (e.g., for a new open-source  
804 benchmark).
- 805 • The instructions should contain the exact command and environment needed to run to  
806 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
807 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 808 • The authors should provide instructions on data access and preparation, including how  
809 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 810 • The authors should provide scripts to reproduce all experimental results for the new  
811 proposed method and baselines. If only a subset of experiments are reproducible, they  
812 should state which ones are omitted from the script and why.
- 813 • At submission time, to preserve anonymity, the authors should release anonymized  
814 versions (if applicable).
- 815 • Providing as much information as possible in supplemental material (appended to the  
816 paper) is recommended, but including URLs to data and code is permitted.

## 817 6. Experimental Setting/Details

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

821 Answer: [Yes]

822 Justification: We provided a detailed description of the experimental setup and metrics.

823 Guidelines:

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

## 829 7. Experiment Statistical Significance

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

832 Answer: [Yes]

833 Justification: We tested our method with multiple random seeds and calculated the standard  
834 error.

835 Guidelines:

- 836 • The answer NA means that the paper does not include experiments.
- 837 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
838 dence intervals, or statistical significance tests, at least for the experiments that support  
839 the main claims of the paper.
- 840 • The factors of variability that the error bars are capturing should be clearly stated (for  
841 example, train/test split, initialization, random drawing of some parameter, or overall  
842 run with given experimental conditions).

- 843 • The method for calculating the error bars should be explained (closed form formula,  
844 call to a library function, bootstrap, etc.)
- 845 • The assumptions made should be given (e.g., Normally distributed errors).
- 846 • It should be clear whether the error bar is the standard deviation or the standard error  
847 of the mean.
- 848 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
849 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
850 of Normality of errors is not verified.
- 851 • For asymmetric distributions, the authors should be careful not to show in tables or  
852 figures symmetric error bars that would yield results that are out of range (e.g. negative  
853 error rates).
- 854 • If error bars are reported in tables or plots, The authors should explain in the text how  
855 they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

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

860 Answer: [Yes]

861 Justification: We explain the required computational resources and related information in  
862 the appendix.

863 Guidelines:

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

## 9. Code Of Ethics

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

875 Answer: [Yes]

876 Justification: Although our work is related to healthcare, we train and test our models on  
877 offline data, adhering to ethical standards.

878 Guidelines:

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

## 10. Broader Impacts

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

887 Answer: [Yes]

888 Justification: Our work has a positive impact on safe healthcare, promoting the expansion of  
889 artificial intelligence technology into the medical field.

890 Guidelines:

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

- 894
- 895
- 896
- 897
- 898
- 899
- 900
- 901
- 902
- 903
- 904
- 905
- 906
- 907
- 908
- 909
- 910
- 911
- 912
- 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).

## 913 11. Safeguards

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

917 Answer: [NA]

918 Justification: Our work does not pose security risks because it is based on publicly available  
919 datasets and models.

920 Guidelines:

- 921
- 922
- 923
- 924
- 925
- 926
- 927
- 928
- 929
- 930
- 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.

## 931 12. Licenses for existing assets

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

935 Answer: [Yes]

936 Justification: The code, data, and models we referenced are all cited, and we followed the  
937 licenses and terms of use throughout the process.

938 Guidelines:

- 939
- 940
- 941
- 942
- 943
- 944
- 945
- 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.

- 946
- 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.
- 947
- 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.
- 948
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.
- 949
- 950
- 951
- 952
- 953

### 954 13. **New Assets**

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

957 Answer: [Yes]

958 Justification: We will provide detailed data extraction code and model code as part of the  
959 submission files.

960 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.
- 961
- 962
- 963
- 964
- 965
- 966
- 967
- 968

### 969 14. **Crowdsourcing and Research with Human Subjects**

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

973 Answer: [NA]

974 Justification: This paper does not involve crowdsourcing nor research with human subjects.

975 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.
- 976
- 977
- 978
- 979
- 980
- 981
- 982
- 983

### 984 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 985 Subjects**

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

990 Answer: [NA]

991 Justification: This paper does not involve crowdsourcing nor research with human subjects.

992 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.
- 993
- 994
- 995
- 996
- 997

998  
999  
1000  
1001  
1002

- 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.