# Exploring a Graph-based Approach to Offline Reinforcement Learning for Sepsis Treatment

**Taisiya Khakharova**
Chair of Practical Computer Science/Software Systems Engineering
Brandenburg University of Technology Cottbus-Senftenberg, Germany
`Taisiya.Khakharova@b-tu.de`

**Lucas Sakizloglou**
Chair of Practical Computer Science/Software Systems Engineering
Brandenburg University of Technology Cottbus-Senftenberg, Germany
`Lucas.Sakizloglou@b-tu.de`

**Leen Lambers**
Chair of Practical Computer Science/Software Systems Engineering
Brandenburg University of Technology Cottbus-Senftenberg, Germany
`Leen.Lambers@b-tu.de`

## Abstract

Sepsis is a serious, life-threatening condition. When treating sepsis, it is challenging to determine the correct amount of intravenous fluids and vasopressors for a given patient. While automated reinforcement learning (RL)-based methods have been used to support these decisions with promising results, previous studies have relied on relational data. Given the complexity of modern healthcare data, representing data as a graph may provide a more natural and effective approach. This study models patient data from the well-known MIMIC-III dataset as a heterogeneous graph that evolves over time. Subsequently, we explore two Graph Neural Network architectures - GraphSAGE and GATv2 - for learning patient state representations, adopting the approach of decoupling representation learning from policy learning. The encoders are trained to produce latent state representations, jointly with decoders that predict the next patient state. These representations are then used for policy learning with the dBCQ algorithm. The results of our experimental evaluation confirm the potential of a graph-based approach, while highlighting the complexity of representation learning in this domain.

## 1 Introduction

Sepsis is a severe condition in which the body dysfunctionally responds to infection injuring the body's own tissues and organs, which may escalate to septic shock and lead, consequently, to death. Sepsis is not uncommon: an estimated 15 out of every 1000 hospitalized patients will develop sepsis while sepsis-related deaths account for 20% of all global deaths [27].

Vasopressors and intravenous (IV) fluids are essential components in sepsis treatment [34, 24]. IV fluids are used to restore circulating blood volume, while vasopressors address hypotension.

The Surviving Sepsis Campaign guidelines [24] proposed a universal approach, suggesting a fixed dosage of the drug per kilogram of patient body weight. However, this approach has been called into

question by more recent studies [19, 5], which indicate that personalized treatment strategies should be preferred.

Sepsis treatment has been addressed with reinforcement learning (RL) because it naturally fits a Markov decision process (MDP) framework. Many studies have demonstrated promising results in developing personalized treatment policies that could improve patient outcomes compared to standard protocols [15, 26, 39, 38]. RL agents are trained to recommend drug prescriptions based on publicly available medical datasets. Specifically, the agents learned to continually recommend a dosage of vasopressors and IV fluids that should be administered to a patient, taking into account multiple clinical variables such as vital signs, laboratory results, and demographic information.

To the best of our knowledge, all the previous studies on RL-based sepsis treatment optimization rely on relational data. As shown by related work [36], the complexity of medical data requires new architectures for encoding the information, and utilizing graph data structures can positively influence the performance of machine learning models. Graphs can capture semantic relationships for medical data, making further data analysis more efficient and accurate [1].

Graph neural network (GNNs) is a type of neural network that learns from graph-structured data. GNNs are designed to capture complex relationships by exploiting structure. Besides being already useful to healthcare research, e.g., for disease prediction [28], recent research from other application domains indicates that GNNs may be more suitable to RL over graph-structured data [20, 22].

In summary, RL for sepsis treatment is an area of ongoing research that has already shown impressive results. However, the studies are limited to the relational approach: they use relational data and traditional neural networks (NNs), while several studies have highlighted the potential for increased accuracy and efficiency when using medical graph data as input to the machine learning models. This paper investigates how a graph-based approach impacts accuracy and efficiency in RL-optimized sepsis treatment by leveraging GNNs to process the patient graph data.

In this work, we present a way of modeling timestamped medical data from the dataset MIMIC-III [12] as a graph. Then, we compare two state-of-the-art GNN-based encoder architectures. The GNN-based encoders are used to encode the graph data into latent representations. The RL policies are trained with dBCQ algorithm [8]. Then, we evaluate the learned RL policies obtained using these latent representations, and compare them with eath other and with the policy learned from representations that were obtained from traditional NN.

The work of Killian et al. [14] informs the design of our study. Similar to their design we decouple representation learning from RL policy learning to isolate the effect of the different encoders alone and therefore allow for a fair evaluation. We also compare our results with those of this study.

The rest of this paper is organized as follows: Section 2 explains RL terminology, formulates sepsis treatment as POMDP, describes our data, reviews graphs and GNNs; Section 3 presents our approach, with key contributions in Section 3.2 and Section 3.3. Section 4 provides details of the experimental setup, evaluation and threats to validity;Section 5 presents related work; and Section 6 concludes the paper and outlines future work directions.

## 2 Preliminaries

### 2.1 Reinforcement Learning

RL [29] is a subfield of machine learning that introduces concepts of agent and environment. The aim for an agent is to learn how to maximize cumulative reward gained through sequential interactions with the environment. Markov Decision Processes (MDP) provide a mathetical framework to describe the environment.

Offline or batch RL refers to a setting when an agent cannot interract with the environment while learning, but only has experiences documented from the past. The medical domain is a safety critical domain, therefore, we are limited to the offline RL application. In this setting, we must prevent the estimated value function from assigning values to actions that do not appear in the provided data. Batch Constrained Q-Learning (BCQ) [9] addresses this problem. For a setting where the action space is discretized there is an adapted version called dBCQ [8].

Weighted Importance Sampling (WIS) is one of the fundamental approaches to off-policy evaluation. WIS addresses the distribution mismatch between behaviour and target policies by reweighting trajectories based on their likelihood ratios. The method has been adopted in healthcare applications, including sepsis treatment [14, 16].

## 2.2 Sepsis Treatment as POMDP Problem

We build on the optimization of sepsis treatments based on a partially observable MDP (POMDP), as presented in [14]. In this work, an RL agent is trained to make personalized vasopressors and IV fluid dosage suggestions based on the patient's vital signs, lab results and demographic information. In further detail, we reuse the following POMDP elements:

**Action space.** Dosages of vasopressors and IV fluids were split into five quartiles each resulting into 25 possible actions. Therefore, administration of medication is represented by the discrete value from 0 to 25.

**Reward.** Each patient trajectory is associated with a value of +1 if the patient survived and with a value of -1 if not. The reward is given only once per trajectory - in the end.

**Observation space.** The general idea is that a patient's vital signs, lab values and demographic information are aggregated over 4 hours and encoded into the observation space. We opt for a continuous observation space where each observation is the patient information on a currect time step encoded into a latent vector of the dimesionality $\hat{d}_s$, as done in [14].

## 2.3 Medical Dataset and its Modeling as a Graph

Similarly to [15, 14], we evaluate our approach based on the MIMIC-III dataset (v.1.4) that contains data from patients in an Intensive Care Unit [12]. Specifically, we reuse the Python implementation of Killian et al. [14] of the MATLAB solution of Komorowski et al. [15] for all preprocessing and sepsis cohort extraction, resulting in a cleaned sepsis dataset. The sepsis dataset consists of 18,908 patient trajectories that are independent from each other. Killian et al. additionally remove patient trajectories consisting only of one time step and ended up with 18,573 patient trajectories. The mortality rate in the data remained at approximately 6%. Each trajectory consists of up to 20 time steps. Each time step comprises 38 patient features. The features encompass demographic information, lab values and vital signs, which are aggregated over a period of 4 hours. A time step is followed by an administration of a combination of vasopressors and IV fluids.

As mentioned in Section 1, however, contrary to [14, 15], our work explores a graph-based approach. Specifically, we model the dataset as a *dynamic heterogeneous* graph—see Section 3.2: A heterogeneous graph contains multiple types of nodes and/or edges, allowing richer relational structures—e.g., authors, papers, and venues connected by "writes" or "cites" relations [35], whereas a dynamic graph changes over time, with nodes or edges appearing, disappearing, or altering attributes over time [13]. On a formal level, such graphs can be represented via suitably defined *typed attributed graphs* [7]. Typed attributed graphs are graphs whose nodes and edges can have different types, and furthermore be equipped with attributes. Therefore, typed attributed graphs are capable of capturing the information present in the dataset in Section 3.2. Typed attributed graphs are created based on a *type graph*, that is, an ontology that defines all valid graph instances, i.e., typed attributed graphs.

## 2.4 Graph Neural Networks

Graph Neural Networks (GNNs) [40] are deep learning models designed to work with graph-structured data. Unlike traditional neural networks that process fixed-format data (like images or sequences), GNNs can handle irregular structures where each node might have a different number of connections. The key idea behind GNNs is that they learn representations for nodes by aggregating information from their neighbors in the graph.

GNNs update each node's representation through learnable message-passing. In each layer, a node first aggregates features from its neighbors and then applies a trainable transformation—typically a weighted linear layer followed by an activation function—to produce a new feature vector. Stacking L such layers means each node's features are updated L times, allowing the GNN to capture progressively a wider graph structure.
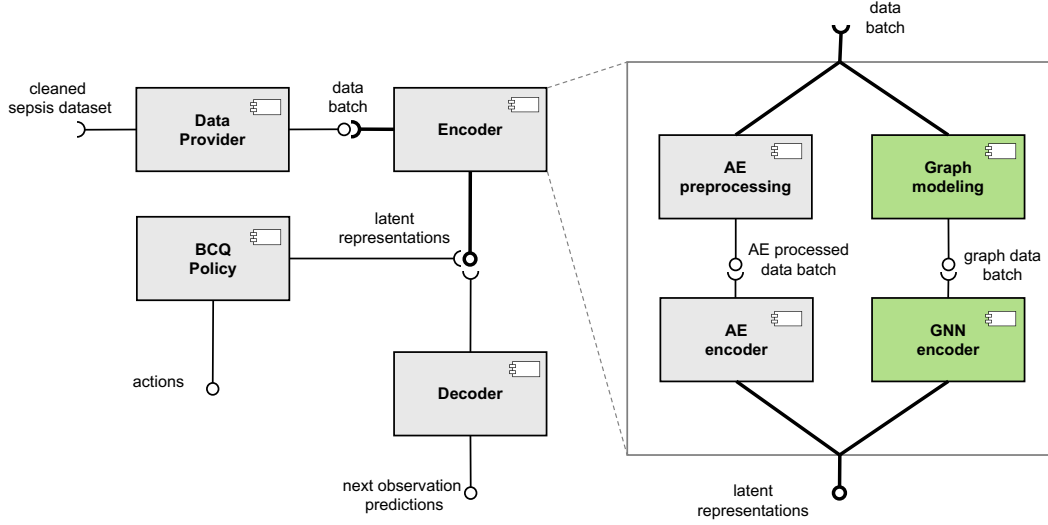
Figure 1: Component diagram of the proposed approach. The components representing our main contribution are shown in green.

# 3 Approach

In this section, we provide an overview of our approach and highlight the main contributions: we present a way of modeling medical data a graph and the design of GNN-based encoders.

## 3.1 Overview

The approach consists of the following building blocks: the Data provider, the Encoder, the Decoder, the BCQ policy. We use the Data provider, the Encoder, and the Decoder for representation learning, while the Data provider, the Encoder, and the BCQ policy are used for policy learning.

The Encoder component is subdivided into two flows: the GNN flow and the AE flow (see Fig. 1). The GNN flow encompasses the Graph modeling component and the GNN encoder. These two components are the primary contribution of this paper and we describe them in the following subsections (see Section 3.2, Section 3.3). The implementation of them the AE flow is reused from [14]. The aim of the AE flow is to present a reference point to compare the GNN flow with.

We also reuse the implementation of the Data provider, the Decoder and the BCQ policy to be able to compare with the results of [14]. The BCQ policy is an abstraction for the BCQ algorithm. The role of the Data provider is to perform training, validation, and testing split with 70/15/15 proportion and then organize the data into batches of size $b$, with each batch containing $b$ trajectories. After the batches are encoded into latent state representations, the Decoder concatenates them with the next performed action and predicts the next observation of the patient. The BCQ policy inputs latent state representations as observations from the environment and outputs the treatment recommendations for clinitians.

Following Killian et al. [14], we decouple representation learning from policy learning. Thus, training contains two phases. In the first phase, we train autoencoders to predict features of the next patient observation, effectively training the AE encoder and the GNN encoders to extract meaning in the form of a latent vector. The autoencoders' training process is detailed in Section 4.1, and the results of the representation learning phase are presented in Appendix A. In the second phase, the trained encoders are used to encode the data into latent representations. The representations serve as environment observations to train the BCQ policy. We evaluate the policy and present the results in Section 4. This setup enables a controlled, like-for-like comparison, meaning that any differences observed can be attributed to the representation (GNN vs. NN) rather that to changes in the BCQ.

## 3.2 Graph modeling

In this section, we explain the transformation process of tabular data into trajectory graphs. Obtaining trajectory graphs is an intermediate transformation. Each trajectory graph is processed further into sequences of graph snapshots to be used as an input to the GNN encoder (see Fig. 1). When modeling we take into account that exclusion of treatment history from the input data might lead to unsafe treatment recommendations [18] and therefore we opted for inclusion of this information into each graph snapshot.

In order to construct a trajectory graph, we divided the patient's features into two groups: time-invariant and time-variant features. We defined time-variant features as the features that change within one trajectory according to the data. We classified gender, age, readmission and ventilation statuses as time-invariant features. Unlike [14] we consider a patient's body weight to be a time-variant feature, as it may change within one trajectory.
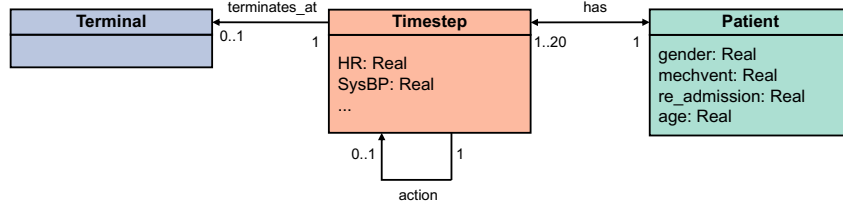


Figure 2: Type graph for patient trajectory graph.

We modeled each patient trajectory as a dynamic heterogeneous graph. We describe the language of such dynamic heterogeneous graph using an attributed type graph, which is depicted in Fig. 2. An exemplary instance graph is shown in Fig. 3. In Fig. 2 and Fig. 3, colors are used solely to facilitate the mapping of node types in instance graphs to the corresponding nodes types in the type graph. Time-invariant features $\mathcal{I}$ are encapsulated in a node type called Patient and time-variant features $\mathcal{T}$ – in Timestep nodes. Table 5 provides an overview of the features grouped into time-invariant and time-variant categories. Every patient trajectory has exactly one Patient node as a central element of the graph. Each Timestep node is connected to the next Timestep with a unidirectional edge that stores information about the action taken in-between the corresponding time steps. Each such edge has a feature vector represented as a one-hot encoding of length 25, containing all zeros except for a single one indicating the action taken. The Patient node is connected to each Timestep node via bidirectional edges, with all edge weights set to 1. The aim is to represent the Patient-Timestep association and to avoid introducing any other domain-specific information by using unequal weights. The bidirectional edges allow messages to be passed into both directions later at the GNN encoding stage. Terminal-type nodes indicate the end of the trajectory and are therefore connected with the last Timestep node in each trajectory. A Terminal node contains information on a reward for an RL agent: 1 if the patient survived and -1 if not (see Section 2.2). We introduced Terminal nodes for convenience of the subsequent transformation; the edge connecting the final Timestep to the Terminal node carries no attributes.

The subsequent stage of the transformation process involves graph snapshots generation from each trajectory graph (see Fig. 12). Each snapshot $g_t$ contains information about the patient's current observation, as well as their previous observations and treatment history. A snapshot $g_t$ comprises a Patient node from the trajectory graph and Timestep nodes from the initial time point up to the current time point. Each $g_t$ is encoded at the next stage to a latent representation $l_t$. A batch of graphs refers to a certain number of sequences of graph snapshots.

## 3.3 GNN encoder

The GNN encoder component produces a GNN designed to work with heterogeneous graphs containing different types of nodes and edges. The encoder consists of several heterogeneous graph convolution layers, followed by pooling operations and a linear layer. The final output is a fix-size latent representation for each graph. We use 2 alternative convolution layer types: SAGEConv operators from the GraphSAGE framework (henceforth SAGE) [10] and the GATv2Conv operator
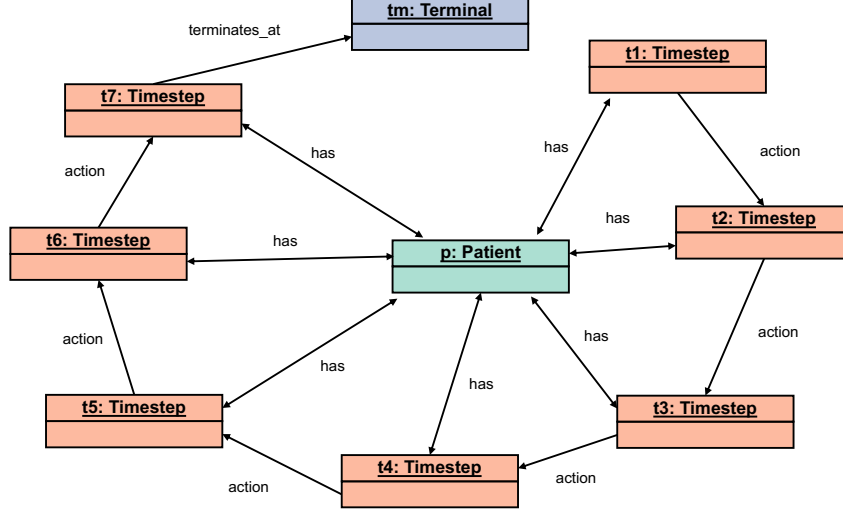
Figure 3: Patient trajectory graph with 7 time steps.

from the GATv2 architecture[4]. The variable parameters for the GNN are the number of output features $f_{out}$ in convolutional layers and the number of convolutional layers $n_{conv}$.

**SAGEConv**    During the forward pass, the input passes through the $n_{conv}$ convolution SAGEConv layers, where each layer contains 3 SAGEConv operators corresponding to the number of edge types. Each layer is followed by a ReLU [21] activation function. The features are aggregated using 'mean' for each edge type, which means that if a node has more than one incoming edge of the same type, then these incoming messages are averaged. In our case, this parameter is relevant for the Patient node, because this node has many incoming edges of the same type from Timestep nodes, so, messages from each Timestep are averaged. After that, features from different edge types are summed up to accumulate information from different types of relationships. This step is important for Timestep nodes, as Timestep nodes have 2 different types of incoming edges.

After the convolution layers, graph-level mean pooling is applied in order to aggregate node features for each node type. As a result, we get a feature vector per node type. As we have two node types, Timestep and Patient, after the graph-level mean pooling we obtain two feature vectors. The next step is to aggregate them with the summation over different node types in order to get one vector. The size of this vector is defined by the hyperparameter $f_{out}$. As a next step, we add a linear layer of size $(f_{out}, l)$ to ensure that the output vector is of length $l$ even if $f_{out} \neq l$.

**GATv2Conv**    The second architecture uses the GATv2Conv operator in the layers. The difference with the SAGEConv architecture is that this operator (1) takes into account the weight of the edges in the input graph, (2) has additional training weights that indicate the importance of each neighbour node. This means that instead of using a simple average aggregation for each edge type, GATv2Conv calculates the weighted average.

## 4   Empirical Study

Our aim is to answer the question of how using representations derived from graph-structured data via trained GNN encoders influences the accuracy and efficiency of policy learning. To answer this question, we train GNN and AE encoders, along with a decoder, to predict the next patient's state. The downstream task for the encoders, as derived from the research question, is policy learning by a RL agent using the dBCQ. We evaluate the learned policies using WIS.

The hyperparameters that are common to both the autoencoders training of the encoders and to the RL policy learning are $l$, the size of the latent state representation vectors, which is set to 64; and $b$, the batch size, which is set to 128.

## 4.1 AE and GNN autoencoders training and hyperparameters

The objective of an autoencoder is to reconstruct the vector of time-variant patient features $c_{t+1}$, a subset of the observation $o_{t+1}$ available in the data (see Section 3.2).The reconstruction is based on the current observation $o_t$, the preceding action $a_{t-1}$ (or their history in the case of GNNs), and, at the decoder stage, the subsequent action $a_t$.

The goal of training is to obtain encoders that extract meaningful information from a patient observation in the form of a latent vector of size $l$. The training process work as follows: each encoder receives the current observation $o_t$ togther with the preceding action $a_{t-1}$. Based on this input, the encoder produces a latent representation vector of size $l$, which is concatenated with the subsequent action $a_t$ known from the data. The resulting vector serves as input to the decoder, which outputs a prediction of time-variant features $\hat{c}_{t+1}$ of the next patient observation. The loss is calculated as the differnce between the true time-variant patient features $c_{t+1}$ and the predicted features $\hat{c}_{t+1}$.

**Decoder**   While the encoder architectures differ, the decoder architecture remains nearly identical across all considered encoders. The decoder consists of three fully-connected layers $(l + 25, 64), (64, 128), (128, \text{obs\_dim})$ . For the GNNs case, $\text{obs\_dim} = 34$, and for the AE case, $\text{obs\_dim} = 33$. The hidden layers are followed by a ReLU activation function. The final layer of size obs\_dim produces a vector that serves as the mean of a unit-variance multivariate Gaussian distribution. This is then used to predict the time-variant features of the next patient observation $\hat{c}_{t+1}$.

**AE encoder**   For the AE encoder, we use the optimal hyperparameters from previous work [14]. It is a three-layered, fully connected NN with layer sizes of: $(\text{obs\_dim} + 25, 64), (64, 128)$ and $(128, l)$. The ReLU activation function is applied to the hidden layers. We train it for 600 epochs using the Adam optimizer with a learning rate of $5e{-}4$.

**GNN encoders**   We conduct an empirical study to determine the optimal hyperparameters for the GNN encoders. As for tuning the hyperparameters of the SAGEConv encoder, we set $(f_{out}, n_{conv}) \in \{2, 3\} \times \{64, 128\}$, i.e., $(f_{out}, n_{conv}) = (2, 64), (2, 128), (3, 64), (3, 128)$. The best-performing model is found at $(f^*_{out}, n^*_{conv}) = (2, 64)$. Setting $f_{out}$ to 128 results in clear overfitting. Therefore, for GATv2Conv, we fix $f^*_{out}$ to 64 and vary $n_{conv}$ within the range $\{1, 2, 3\}$. The best performing model has $n^*_{conv} = 1$. We train the GNN autoencoders for 200 training epohs using the Adam optimizer with a learning rate of $1e{-}3$. Appendix A contains the visualized training results. An overview of the selected hyperparameters is provided in Table 1.

Although all autoencoders share same training principle, they differ fundamentally in how they handle input data. The AE encoder has no access to the history of past observations or treatments: it processes each time step independently by concatenating the previous action $a_{t-1}$ with the current observation $o_t$. Any information about observations at a time step $\leq t - 1$ are not available to the network (see Section 4.4). By contrast, GNN encoders operate on snapshot graphs $g_t$ (see Section 3.2). Each snapshot is still processed independently, but its graph structure inherently contains all observations and actions from time 1 up to $t$. This allows GNNs to incorporate historical information without explicit recurrence.

## 4.2 dBCQ-Policy Learning and Evaluation

We trained policies using data encoded by the three best performing trained encoders: AE, GNN-SAGE and GNN-GATv2. The policy models were trained using the dBCQ algorithm with identical hyperparameters, and the evaluation is performed using WIS (see Section 2).

The dBCQ model was trained on the combined training and validation subsets of the data for one million iterations. The WIS evaluation was performed every 500 iterations on the testing subset. All subsets (training, validation, and testing) were first encoded into latent representations using the selected encoder. The Q-network consisted of three fully connected layers, with 64 nodes in each of the two hidden layers, followed by ReLU activations, and 25 nodes in the output layer. The BCQ action elimination threshold was set to 0.3, as in related work [14]. The learning rate was set to $1 \times 10^{-3}$ for all experiments. An overview of the selected hyperparameters is provided in Table 6.

To perform the WIS evaluation, we first trained clinician-like behavioral policies documented in the data. We used a neural network with three fully connected layers $(38, 128), (128, 128), (128, 25)$,

trained with cross-entropy loss. ReLU activations were applied to the hidden layers, and batch normalization was employed to regularize training. The supervised task consisted of imitating clinicians' actions based on patient observations. Since the actions were discretized (see Section 2.2), this was formulated as a classification problem. Training was carried out with the Adam optimizer at a learning rate of $1 \times 10^{-4}$. The behavioral policy was trained directly on the raw data, without encoding into latent representations. The set of hyperparameters used in our experiments is summarized in Table 7.

We run the experiments three times for each of the policy experiment curves, using the random seeds 1234, 2020, and 2025. We apply an exponential moving average with smoothing parameter $alpha = 0.1$. We then calculate the mean and standard deviation, plotting the results. Fig. 4 demonstrates the results of the WIS evalution of the policies learnt from observations obtained from different trained encoders – AE, GNN-SAGE and GNN-GATv2. We set the number of training iterations to 1M. Initially, we trained on 500K (see Fig. 10), but this number of iteration was not sufficient for the GNN cases to show any significant growth. The GNN-SAGE curve shows steady growth from 375K iterations, while GNN-GATv2 - starts to show growth only after 600K iterations. The AE learning curve demonstrates significatly faster growth than both GNN representation types: it demonstrates a rapid increase at 75K iterations. However, the WIS score for AE representations plateaux at 0.68 after around 350K iterations, while the score for GNN-SAGE representation grows to 0.75 by the end of the training. GNN-GATv2 representations show the worst performance with respect to both accuracy and efficiency, as steady growth only begins after 600K iterations, and by the end of the training, the WIS score is just above 0.5.
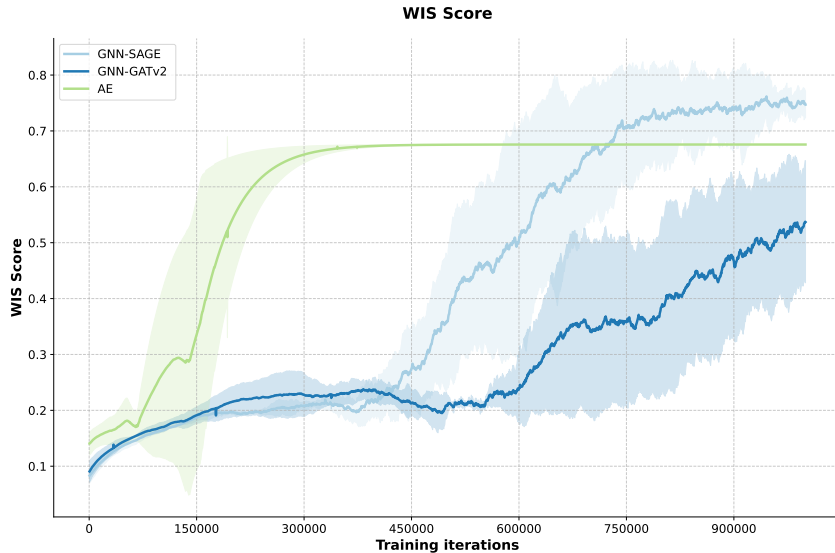


Figure 4: WIS on 1e6 training iteration on GNN-SAGEConv and GNN-GATv2Conv
.

## 4.3 Discussion

Compared to Killian et al. [14], we successfully reproduced the AE learning curve, although their plot covers 200K iterations while ours extends to 1M. As Killian et al. do not present results beyond 200K, it remains unclear how performance evolves, and our work provides an inside into what happens to the AE curve on the later iterations. However, their work clearly demonstrates that policies learned from recurrent encodings – especially CDE – converge much faster than those based on GNNs. In their results, the CDE curve plateaux at a WIS score of approximately 0.78 after 75K iterations. In contrast, the GNN-SAGE curve requires significantly more steps to begin showing steady improvement, but it ultimately demonstrates competitive performance, reaching a WIS score of approximately 0.75. This suggests that, when sufficiently trained, a GNN-SAGE based policy can match – or nearly match – the acccucy of a CDE based policy. Therefore , our future work will investigate encoders based on a combination of RNNs and GNNs.

The accuracy results for autoencoders training did not correspond to those for the downstream policy learning task. AE and GNN-SAGE performed worse than GNN-GATv2 in the autoencoders training task, but better in policy learning. While the difference in autoencoders training results between AE and GNN-SAGE was marginal, the shapes of their respective policy training learning curves were very different.

It appears that the additional learned weights in GNN-GATv2, which represent the importance of each node neighbour, might mislead the policy learning agent, thereby slowing down the learning curve compared to GNN-SAGE, where edge weights are neither taken into account nor learned. This is consistent with the findings of Killian et al., who discovered that more complex models do not necessarily result in a better learning curve.

The experiments were conducted on a system comprising an AMD EPYC 9554 processor with 128 cores running at 3.76 GHz, three NVIDIA L40 GPUs each with 46 GB of memory, and 1.5 TB of system RAM. The system runs CUDA version 12.7. Training the BCQ policy on the GPU took approximately three hours for each experiment with 500K iterations, and around six hours with 1M iterations. Traning the autoencoders on the GPU for each GNN took around 30 hours, whereas AE took between three and four hours, around ten times faster. We trained the models using both the CPU and the GPU. For GNN autoencoders training, using the GPU does not speed things up. One possible reason for this is that the graphs we use are relatively small; this issue will be addressed in future work.

## 4.4 Threats to Validity

This section discusses threats to validity. Specifically, we consider 4 aspects of validity–construct, internal, external, reliability [37]–which we outline here and aim to address in future work.

**Construct Validity**  We identify several construct-validity threats. First, because our goal was to study GNNs, we converted the data into a graph, but this was only one possible modeling approach, so it remains unclear how GNNs would perform with a different data representation. Second, we evaluated learned policies using a quantitative metric–WIS score, which, on the one hand, was used in related work to evaluate sepsis treatment policies leaned in the offline RL settings [14], but, on the other hand, we did not perform any qualitative assessment, so it is uncertain whether high scores truly reflect clinically meaningful decisions. Finally, to estimate the RL policy's value via WIS, we must first learn a NN approximation of the clinitians' behavior policy; because this surrogate is imperfect, the importance weights are biased and the estimated policy values may not reflect true performance of the RL policies.

**Internal Validity**  We address an internal-validity threat by decoupling representation learning from policy learning. This isolation ensures that any performance differences can be attributed to the encoder itself rather than to entangled training dynamics. Another internal-validity concern arises from our experimental setup: the GNN models utilize historical data directly incorporated into the graph structure, whereas the AE does not. This difference makes it challenging to isolate the architectural effects from simply having more input information. Nonetheless, the results show, that even when both SAGE and GATv2 encoder architectures received identical inputs (including history), their performances diverged–SAGEConv outperformed AE, while GATv2 underperformed relative to AE. These mixed results suggest that architectural differences may matter more than mere access to historical data. Finally, to reduce random error in the WIS analysis, we ran each experiment multiple times, averaged the results, and reported the standard deviations.

**External Validity**  We addressed external-validity threats by using the standard evaluation metric–WIS, which is common in RL-based sepsis treatment and allows comparison with the previous study of Killian et al. [14]. However, since we only used the MIMIC-III dataset, it remains unclear whether our results would generalize to other datasets—this limitation necessitates further investigation.

**Reliability**  In terms of reliability, we consistently used the same stratified training, validation, and test splits throughout our experiments; however, because we did not fix a random seed when creating those splits, a future researcher attempting to reproduce our work could obtain different partitions. We mitigated imbalance by stratifying on trajectory outcomes, ensuring each subset

contained proportional positive and negative cases. Regarding random NN initialisation for the RL policy learning, we set and report random seeds. To further reduce variability, we report mean results of multiple runs and standard deviations, providing transparency about stability. Finally, Fig. 4 shows that the results of policy training from the AE representations match those of Killian et al. [14] for the first 200K training iterations, despite a different data split, supporting the reliability of our findings.

## 5 Related work

GNNs are used in the context of representation learning in the heathcare domain [36, 25]. The survey [23] shows that there has been an increase in the number of studies using GNN for clinical risk prediction since 2020. Unlike us, most of the papers focus on diagnosis prediction rather than treatment policy. The most popular architecture used is GAT [33].

One of the first mature studies on the RL applied to sepsis treatment was published by Komorowski et al. [15]. They used a discrete action space for treatment recommendations, which was adopted by many subsequent studies, as summarised by Roggeveen et al. [26]. Huang et al., on the other hand utilize continuous action spaces [11]. Komorowski et al.[15] used a discrete state space by applyng clustering to patient states and utilizing tabular Q-learning. However, most later works use continuous state spaces and apply deep RL [6, 39, 38, 30, 17]. Killian et al. and Huang et al. encoded patient states into latent state observations [14, 11]. The reward in related work is either terminal, based on the survival of the patient, or includes intermediate signals based on short-term changes in the patient's condition, with no clear preference for either variant. The field of RL for sepsis treatment is continually being explored, with new studies appearing every year [32].

GNNs have also been successfully combined with RL in domains such as Google Research Football [22] and routing optimization [2]. In contrast to our work, these studies optimized the RL policy and the representation jointly through end-to-end training.

To our knowledge, there are no papers that use GNNs to find an optimal sepsis treatment strategy. In the domain of sepsis, there was a study that focuses on predicting the onset of sepsis [31]. They also model medical data as a graph and use GNNs to learn the representations. The dataset they use is not openly available.

## 6 Conclusion and Future Work

We applied a graph-based approach to optimize sepsis treatment with reinforcement learning. Patient data from MIMIC-III were modeled as dynamic heterogeneous graphs, and we trained GNN based autoencoders (based on the GraphSAGE and GAT architectures) to generate latent state representations. Following prior work, representation learning was decoupled from policy learning and the learned representations were used to train an offline RL agent. Our results show that the GraphSAGE encoder outperforms traditional relational methods in accuracy and is comparable to the (recurrent-encoding) prior approach, although it requires more training iterations to converge.

In future work we will address further threats to validity from Section 4.4. Firstly, we will perform a qualitative analysis of the learned policy for comparison with other approaches, and will evaluate its clinical interpretability. Secondly, we will explore alternative graph modeling approaches. Specifically, we will distribute the patient features to multiple nodes instead of modeling them as attributes of one node. This should address the issue of long training time. Finally, we plan to explore dynamic GNN encoders that incorporate recurrent neural networks in order to better capture temporal dependencies.

## References

[1] Bader Aldughayfiq, Farzeen Ashfaq, N. Z. Jhanjhi, and Mamoona Humayun. Capturing semantic relationships in electronic health records using knowledge graphs: An implementation using MIMIC III dataset and GraphDB. 11(12):1762.

[2] Paul Almasan, José Suárez-Varela, Krzysztof Rusek, Pere Barlet-Ros, and Albert Cabellos-Aparicio. Deep reinforcement learning meets graph neural networks: exploring a routing optimization use case. 196:184–194.

[3] Lukas Biewald. Experiment tracking with weights and biases, 2020. Software available from wandb.com.

[4] Shaked Brody, Uri Alon, and Eran Yahav. How attentive are graph attention networks?

[5] E. Carlos Sanchez, Michael R. Pinsky, Sharmili Sinha, Rajesh Chandra Mishra, Ahsina Jahan Lopa, and Ranajit Chatterjee. Fluids and early vasopressors in the management of septic shock: Do we have the right answers yet? 9(3):138–147.

[6] Yunho Choi, Songmi Oh, Jin Won Huh, Ho-Taek Joo, Hosu Lee, Wonsang You, Cheng-mok Bae, Jae-Hun Choi, and Kyung-Joong Kim. Deep reinforcement learning extracts the optimal sepsis treatment policy from treatment records. 4(1):245.

[7] Hartmut Ehrig, Ulrike Prange, and Gabriele Taentzer. Fundamental theory for typed attributed graph transformation. In Hartmut Ehrig, Gregor Engels, Francesco Parisi-Presicce, and Grzegorz Rozenberg, editors, *Graph Transformations*, volume 3256, pages 161–177. Springer Berlin Heidelberg. Series Title: Lecture Notes in Computer Science.

[8] Scott Fujimoto, Edoardo Conti, Mohammad Ghavamzadeh, and Joelle Pineau. Benchmarking batch deep reinforcement learning algorithms.

[9] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. Version Number: 3.

[10] William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. Version Number: 4.

[11] Yong Huang, Rui Cao, and Amir Rahmani. Reinforcement learning for sepsis treatment: A continuous action space solution. In Zachary Lipton, Rajesh Ranganath, Mark Sendak, Michael Sjoding, and Serena Yeung, editors, *Proceedings of the 7th Machine Learning for Healthcare Conference*, volume 182 of *Proceedings of Machine Learning Research*, pages 631–647. PMLR, 05–06 Aug 2022.

[12] Alistair E.W. Johnson, Tom J. Pollard, Lu Shen, Li-wei H. Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. MIMIC-III, a freely accessible critical care database. 3(1):160035.

[13] Seyed Mehran Kazemi. *Dynamic Graph Neural Networks*, pages 323–349. Springer Nature Singapore, Singapore, 2022.

[14] Taylor W. Killian, Haoran Zhang, Jayakumar Subramanian, Mehdi Fatemi, and Marzyeh Ghassemi. An empirical study of representation learning for reinforcement learning in healthcare.

[15] Matthieu Komorowski, Leo A. Celi, Omar Badawi, Anthony C. Gordon, and A. Aldo Faisal. The artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care. 24(11):1716–1720.

[16] Luchen Li, Matthieu Komorowski, and Aldo A. Faisal. Optimizing sequential medical treatments with auto-encoding heuristic search in POMDPs.

[17] Dayang Liang, Huiyi Deng, and Yunlong Liu. The treatment of sepsis: an episodic memory-assisted deep reinforcement learning approach. 53(9):11034–11044.

[18] MingYu Lu, Zachary Shahn, Daby Sow, Finale Doshi-Velez, and Li-Wei H. Lehman. Is deep reinforcement learning ready for practical applications in healthcare? a sensitivity analysis of duel-DDQN for hemodynamic management in sepsis patients. 2020:773–782.

[19] Charalampos D. Moschopoulos, Dimitra Dimopoulou, Anastasia Dimopoulou, Konstantina Dimopoulou, Konstantinos Protopapas, Nikolaos Zavras, Sotirios Tsiodras, Anastasia Kotanidou, and Paraskevi C. Fragkou. New insights into the fluid management in patients with septic shock. 59(6):1047.

[20] Sai Munikoti, Deepesh Agarwal, Laya Das, Mahantesh Halappanavar, and Balasubramaniam Natarajan. Challenges and opportunities in deep reinforcement learning with graph neural networks: A comprehensive review of algorithms and applications. *IEEE Transactions on Neural Networks and Learning Systems*, 35(11):15051–15071, 2024.

[21] Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, ICML'10, page 807–814, Madison, WI, USA, 2010. Omnipress.

[22] Yizhan Niu, Jinglong Liu, Yuhao Shi, and Jiren Zhu. Graph neural network based agent in google research football.

[23] Heloísa Oss Boll, Ali Amirahmadi, Mirfarid Musavian Ghazani, Wagner Ourique de Morais, Edison Pignaton de Freitas, Amira Soliman, Farzaneh Etminani, Stefan Byttner, and Mariana Recamonde-Mendoza. Graph neural networks for clinical risk prediction based on electronic health records: A survey. *Journal of Biomedical Informatics*, 151:104616, 2024.

[24] Andrew Rhodes, Laura E. Evans, Waleed Alhazzani, Mitchell M. Levy, Massimo Antonelli, Ricard Ferrer, Anand Kumar, Jonathan E. Sevransky, Charles L. Sprung, Mark E. Nunnally, Bram Rochwerg, Gordon D. Rubenfeld, Derek C. Angus, Djillali Annane, Richard J. Beale, Geoffrey J. Bellinghan, Gordon R. Bernard, Jean-Daniel Chiche, Craig Coopersmith, Daniel P. De Backer, Craig J. French, Seitaro Fujishima, Herwig Gerlach, Jorge Luis Hidalgo, Steven M. Hollenberg, Alan E. Jones, Dilip R. Karnad, Ruth M. Kleinpell, Younsuk Koh, Thiago Costa Lisboa, Flavia R. Machado, John J. Marini, John C. Marshall, John E. Mazuski, Lauralyn A. McIntyre, Anthony S. McLean, Sangeeta Mehta, Rui P. Moreno, John Myburgh, Paolo Navalesi, Osamu Nishida, Tiffany M. Osborn, Anders Perner, Colleen M. Plunkett, Marco Ranieri, Christa A. Schorr, Maureen A. Seckel, Christopher W. Seymour, Lisa Shieh, Khalid A. Shukri, Steven Q. Simpson, Mervyn Singer, B. Taylor Thompson, Sean R. Townsend, Thomas Van Der Poll, Jean-Louis Vincent, W. Joost Wiersinga, Janice L. Zimmerman, and R. Phillip Dellinger. Surviving sepsis campaign: International guidelines for management of sepsis and septic shock: 2016. 43(3):304–377.

[25] Emma Rocheteau, Catherine Tong, Petar Veličković, Nicholas Lane, and Pietro Liò. Predicting patient outcomes with graph representation learning. Version Number: 1.

[26] Luca Roggeveen, Ali El Hassouni, Jonas Ahrendt, Tingjie Guo, Lucas Fleuren, Patrick Thoral, Armand Rj Girbes, Mark Hoogendoorn, and Paul Wg Elbers. Transatlantic transferability of a new reinforcement learning model for optimizing haemodynamic treatment for critically ill patients with sepsis. 112:102003.

[27] Kristina E Rudd, Sarah Charlotte Johnson, Kareha M Agesa, Katya Anne Shackelford, Derrick Tsoi, Daniel Rhodes Kievlan, Danny V Colombara, Kevin S Ikuta, Niranjan Kissoon, Simon Finfer, Carolin Fleischmann-Struzek, Flavia R Machado, Konrad K Reinhart, Kathryn Rowan, Christopher W Seymour, R Scott Watson, T Eoin West, Fatima Marinho, Simon I Hay, Rafael Lozano, Alan D Lopez, Derek C Angus, Christopher J L Murray, and Mohsen Naghavi. Global, regional, and national sepsis incidence and mortality, 1990–2017: analysis for the global burden of disease study. 395(10219):200–211.

[28] Amit Sharma, Pradeep Kumar Singh, Polina Nikashina, Vadim Gavrilenko, Alexey Tselykh, and Alexander Bozhenyuk. *AI and GNN Model for Predictive Analytics on Patient Data and Its Usefulness in Digital Healthcare Technologies*, pages 331–345. Springer International Publishing, Cham, 2023.

[29] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

[30] Dipesh Tamboli, Jiayu Chen, Kiran Pranesh Jotheeswaran, Denny Yu, and Vaneet Aggarwal. Reinforced sequential decision-making for sepsis treatment: The PosNegDM framework with mortality classifier and transformer. pages 1–9.

[31] Ankur Teredesai, Sijin Huang, Tucker Stewart, Juhua Hu, Armaan Thakker, Katherine Stern, and Grant E O'Keefe. Sub-sequence graph representation learning on high variability data for dynamic risk prediction in critical care. In *2022 IEEE International Conference on Big Data (Big Data)*, pages 2082–2092, 2022.

[32] Rui Tu, Zhipeng Luo, Chuanliang Pan, Zhong Wang, Jie Su, Yu Zhang, and Yifan Wang. Offline safe reinforcement learning for sepsis treatment: Tackling variable-length episodes with sparse rewards.

[33] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks.

[34] Jason Waechter, Anand Kumar, Stephen E. Lapinsky, John Marshall, Peter Dodek, Yaseen Arabi, Joseph E. Parrillo, R. Phillip Dellinger, and Allan Garland. Interaction between fluids and vasoactive agents on mortality in septic shock: A multicenter, observational study*. 42(10):2158–2168.

[35] Xiao Wang, Deyu Bo, Chuan Shi, Shaohua Fan, Yanfang Ye, and Philip S. Yu. A survey on heterogeneous graph embedding: Methods, techniques, applications and sources.

[36] Tingyi Wanyan, Hossein Honarvar, Ariful Azad, Ying Ding, and Benjamin S. Glicksberg. Deep learning with heterogeneous graph embeddings for mortality prediction from electronic health records. 3(3):329–339.

[37] Claes Wohlin, Per Runeson, Martin Höst, Magnus C. Ohlsson, Björn Regnell, and Anders Wesslén. *Experimentation in Software Engineering*. Springer Berlin Heidelberg.

[38] XiaoDan Wu, RuiChang Li, Zhen He, TianZhi Yu, and ChangQing Cheng. A value-based deep reinforcement learning model with human expertise in optimal treatment of sepsis. 6(1):15.

[39] Tianyi Zhang, Yimeng Qu, Deyong Wang, Ming Zhong, Yunzhang Cheng, and Mingwei Zhang. Optimizing sepsis treatment strategies via a reinforcement learning model. 14(2):279–289.

[40] Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications.

# A    Representation Learning Phase

Figs. 5 to 9 were produced with the visualization tools provided by Weights & Biases [3]. The name of each run in the figures is automatically generated and provides no additional information other than being a unique, human-readable run identifier.

Because obs_dim varies depending on the encoder type (see Section 4.1), we normalize the loss when comparing the encoders performance . For both training and validation sets, we compute the loss per batch and average over all batches. Finally, each average batch loss is divided by the batch size to obtain an average trajectory loss.

## A.1    AE Encoder Training

To replicate Killian et al.'s AE training [14], we adopted their hyperparameters: 600 training epochs, a learning rate of $5 \times 10^{-4}$, and validation every 10 epochs. The AE was run with four random seeds (1234, 25, 32, and 2020); the resulting loss curves are shown in Fig. 5.

The loss reported (see Fig. 5) is the average trajectory loss. No smoothing was applied to the training loss; however, the validation loss was smoothed using a 10-step running average. Of the four runs, one began to overfit around epoch 30, but its validation loss resumed declining after peaking at approximately epoch 350. By epoch 600, the average validation loss across all runs was 15.448765.

## A.2    GNN Encoder Training

Both GNN-SAGEConv and GNN-GATv2Conv were trained for 200 epochs. For visualization, we smoothed the validation curves using a 10-step running average. Only encoder hyperparameters varied between runs; decoder hyperparameters remained fixed. Table 1 provides the overview of the GNN encoders hyperparameters.
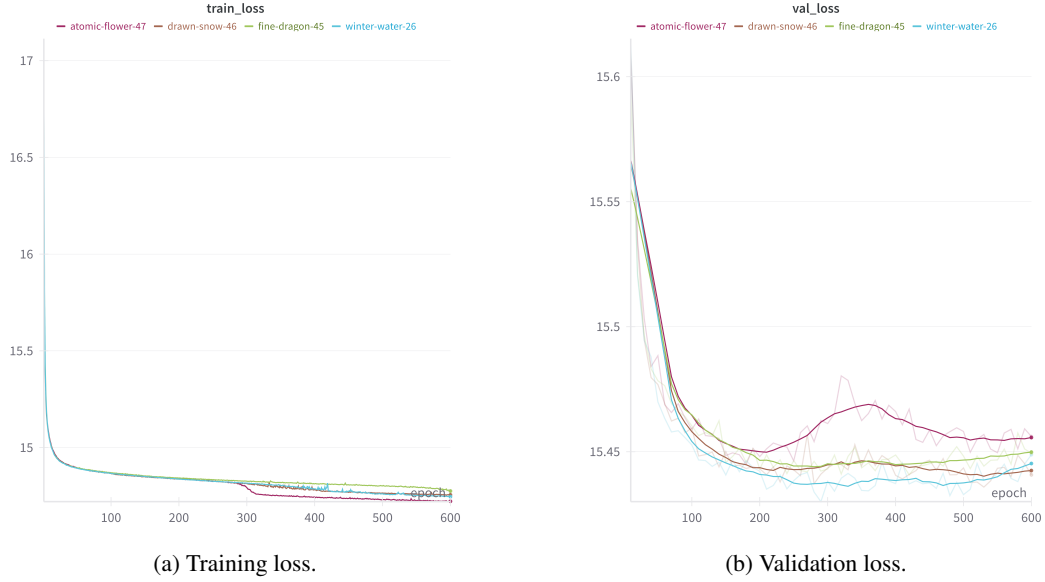
(a) Training loss.
(b) Validation loss.

Figure 5: Training AE autoencoder with four different seeds.

Table 1: Selected hyperparameters for the GNN encoders.

| Encoder | Epochs | lr | $f^*_{out}$ | $n^*_{conv}$ | Intra aggr. | Inter aggr. |
|---|---|---|---|---|---|---|
| GNN-SAGEConv | 200 | $1e{-}3$ | 64 | 2 | mean | sum |
| GNN-SAGEConv | 200 | $1e{-}3$ | 64 | 1 | attn. | sum |

**GNN-SAGEConv** GNN-SAGEConv was trained under four hyperparameter combinations: $f_{out} = \{64, 128\}$ and $n_{conv} = \{2, 3\}$ (Table 2). The best model—lowest validation loss—used $(f_{out}, n_{conv}) = (64, 2)$ (Fig. 6). Increasing the number of features in the graph convolutional layers to 128 led to overfitting for the *cerulean-river-39* and *golden-mountain-41* runs. Choosing the lower number of features and higher number of layers slightly increased the validation loss.

**GATv2Conv** We evaluated GNN-GATv2Conv by fixing $f_{out} = 64$ (since $f_{out} = 128$ caused overfitting) and varying the number of convolutional layers: $n_{conv} \in \{1, 2, 3\}$. Table 3 maps each run name to its $n_{conv}$ value. As shown in Fig. 7, all three configurations converged to nearly the same validation loss. We selected *balmy-bird-49* as the best model because it reached convergence slightly faster and exhibited fewer fluctuations after epoch 50 compared to *clear-lake-44*. We attribute this to oversmoothing: our graph is relatively small, so adding more GATv2 layers tends to blur node features, and a lower $n_{conv}$ yields optimal performance.

**GNN-SAGEConv and GNN-GATv2Conv Comparison** Putting the results of both the GNN-SAGEConv and the GNN-GATv2Conv in Fig. 8, we see that the validation loss curves for the GNN-SAGEConv never overlap with those of the GNN-GATv2Conv; all the GNN-GATv2Conv curves lie below the GNN-SAGEConv curves. Although the training losses are similar, the GNN-GATv2Conv

Table 2: Considered hyperparameters combinations for the GNN-SAGEConv.

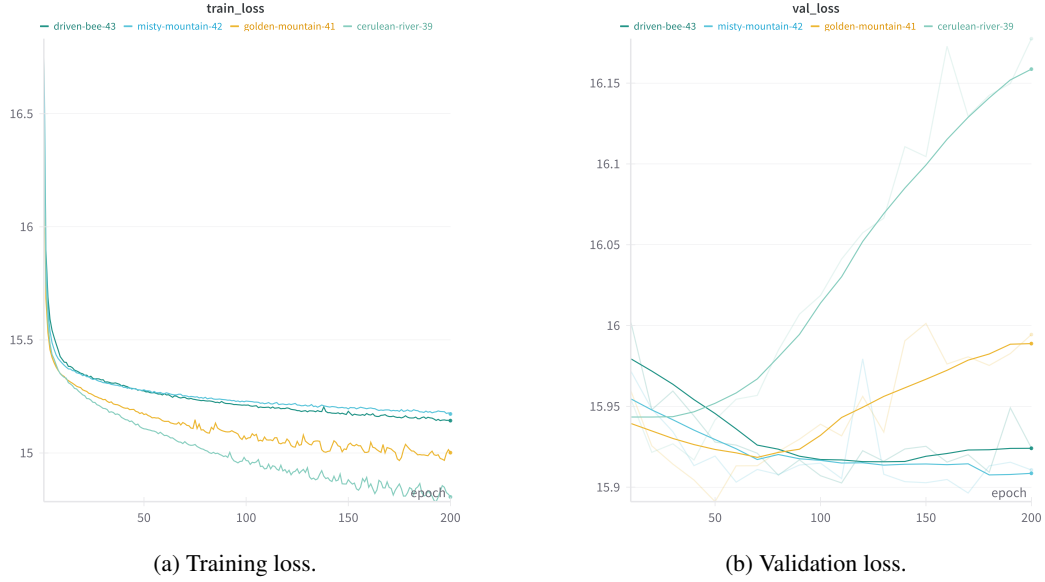| Run name | $f_{out}$ | $n_{conv}$ |
|---|---|---|
| golden-mountain-41 | 128 | 2 |
| cerulean-river-39 | 128 | 3 |
| driven-bee-43 | 64 | 3 |
| **misty-mountain-42** | **64** | **2** |

14

(a) Training loss.　　　　　　　(b) Validation loss.

Figure 6: Training GNN-SAGEConv encoder.

Table 3: Considered hyperparameters combinations for the GNN-GATv2Conv.

| Run name | $n_{conv}$ |
|---|---|
| **balmy-bird-49** | **1** |
| clear-lake-44 | 2 |
| atomic-snowball-48 | 3 |

clearly outperforms the GNN-SAGEConv on validation loss. This suggests that leveraging edge information in the GNN-GATv2Conv is critical for the encoder–decoder task.

## A.3 AE and GNN Encoders Comparison

We normalised the loss to compare the performance of the encoders. As the loss documented in the plots is the average trajectory loss, we have derived from the data that the average trajectory length in each data set (training, validation and test) is 13.3 and the median is 13. Therefore, dividing the observed loss by 13.3 gives the average loss per single step prediction. Each step consists of either 34 or 33 features, so dividing by this number we get the average loss per feature. Therefore the final normalization coefficient for the GNN encoder is $\frac{1}{13,3\cdot34}$, while for AE it is $\frac{1}{13,3\cdot33}$.

We normalized the loss to compare the performance of the encoders. Since the loss shown in the plots is the average trajectory loss, we found that the mean trajectory length in each dataset (training, validation, and test) is 13.3 (median 13). Dividing the observed loss by 13.3 yields the average loss per single-step prediction. Each step consists of either 34 features (for GNN) or 33 features (for AE), so dividing by these numbers gives the average loss per feature. Thus, the normalization coefficient is $\frac{1}{13.3\times34}$ for the GNN encoder and $\frac{1}{13.3\times33}$ for the AE.

Table 4: Final training loss and normalized loss for each autoencoder model.

| Encoder–decoder | Loss | Normalized loss | Absolute difference with AE |
|---|---|---|---|
| AE | 15.44 | 0.035179 | 0 |
| GNN–SAGEConv | 15.91 | 0.035184 | 0.000005 |
| GNN–GATv2Conv | 15.75 | 0.034830 | 0.00034 |

15

(a) Training loss.

(b) Validation loss.

Figure 7: Training GNN-GATv2Conv encoder.



(a) Training loss.

(b) Validation loss.
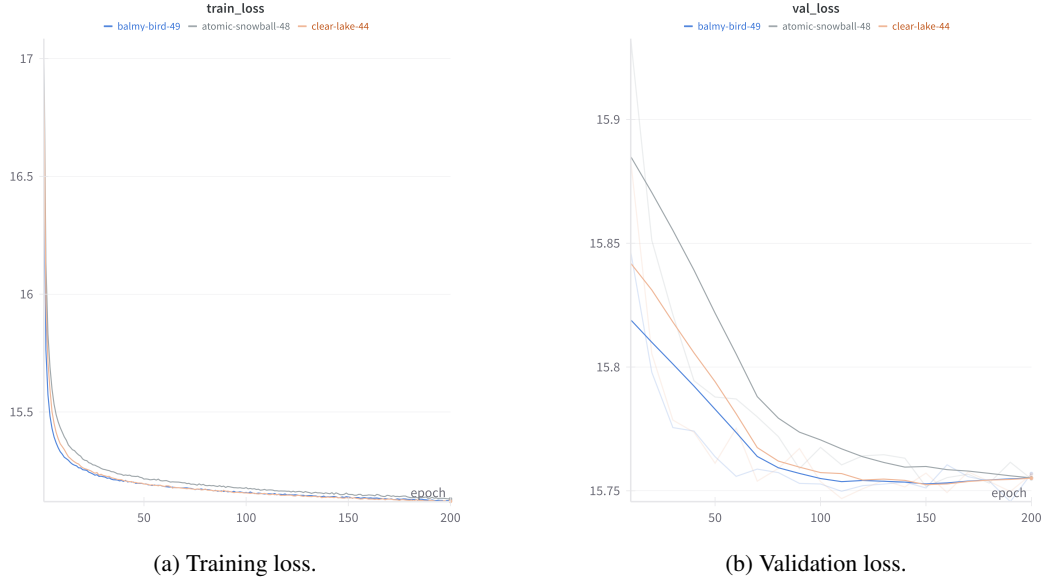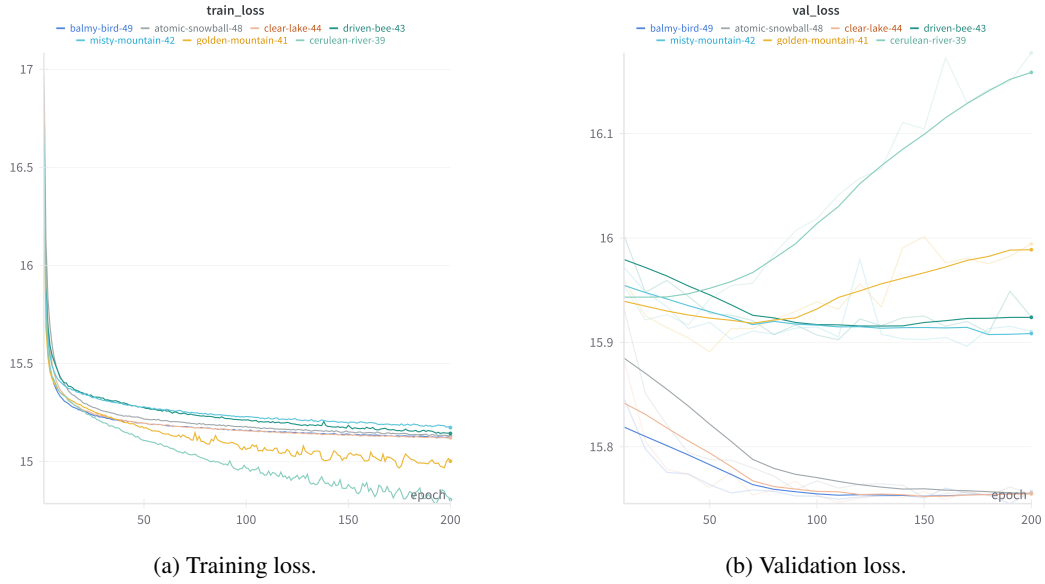
Figure 8: Training GNN-SAGEConv and GNN-GATv2Conv based autoencoders.

The normalization indicates that the best performing model is GNN-GATv2Conv, while GNN-SAGEConv and AE perform nearly the same and therefore share second place. The absolute difference in loss per feature between AE and GNN-SAGEConv is marginal (5e-6), whereas the difference between AE and GNN-GATv2Conv is about 3e-4 —substantially larger.

## A.4   Action Injection Experiment

As an additional experiment, we considered an alternative autoencoders training approach without $a_t$ injection into the decoder input. We trained both GNN encoders with the same hyperparameters, $f_{out} = 64$ and $n_{conv} = 2$, for 200 epochs, performing validation every 10 epochs. The *jumping-water-19* run corresponds to the architecture with action injection, while "upbeat-star-20" corresponds to the version without action injection.

In Fig. 9, the plotted losses are average batch losses. Each batch loss is the sum of all trajectory losses in that batch, and each trajectory loss is the sum of its per-step losses. Although the training losses are similar, the validation loss for the architecture without action injection is noticeably higher. Therefore, in the final encoder–decoder design, we include action injection.
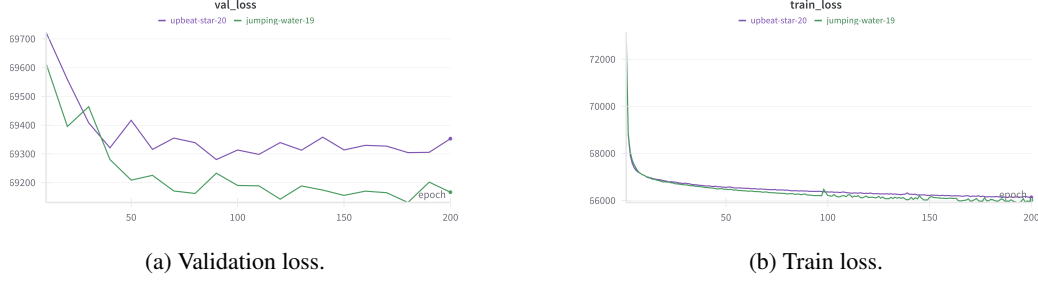


(a) Validation loss.

(b) Train loss.

Figure 9: Experiments with and without action $a_t$ injection.
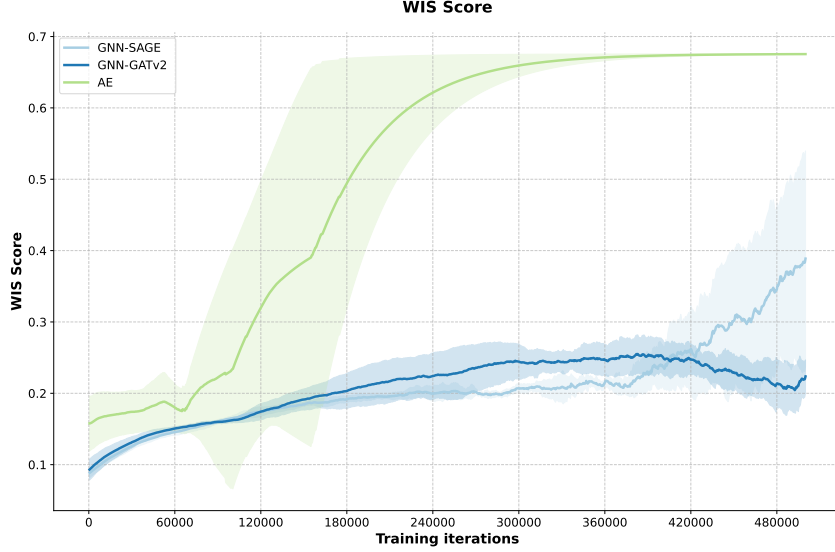


Figure 10: WIS on 500K training iteration on AE, GNN-SAGEConv and GNN-GATv2Conv.

## B   Policy Training with Untrained Encoders

We ran policy training four times for each encoder type (randomly initialized weights) using seeds 1234, 25, 32, and 2020 (see Fig. 11). Each run consisted of 1 million iterations, with evaluations every 500 iterations. The AE encoder failed to learn, while both GNN encoders—GATv2 and SAGEConv—achieved better policies than the AE. As expected, both GNNs exhibited high variance due to their untrained weights, with GATv2 outperforming SAGE.

17

Table 5: Patient features grouped into time-invariant ($\mathcal{I}$) and time-variant ($\mathcal{T}$) categories.

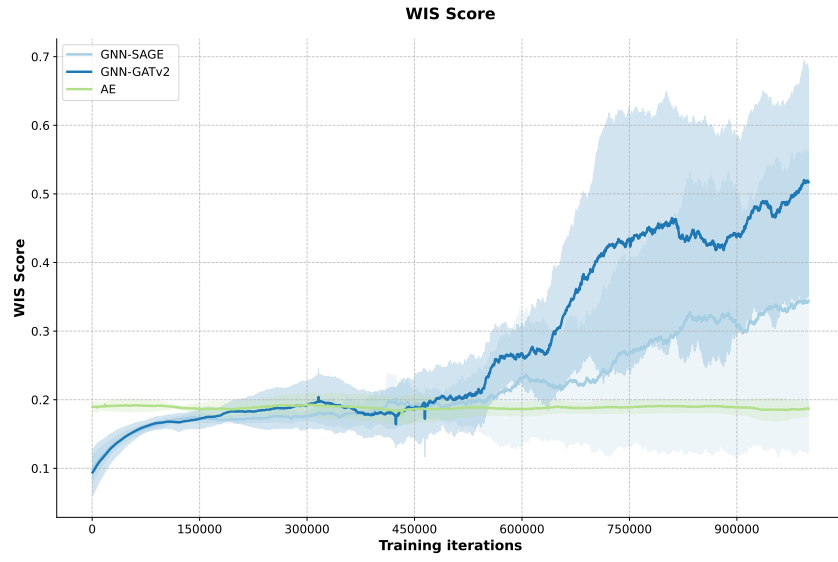| Category | Variable | Human-readable name |
|---|---|---|
| Time-invariant ($\mathcal{I}$) | o:gender | Gender |
| | o:age | Age |
| | o:re_admission | ICU Re-admission |
| | o:mechvent | Mechanical Ventilation (yes/no at admission) |
| Time-variant ($\mathcal{T}$) | o:max_dose_vaso | Maximum Dose of Vasopressor (over 4h) |
| | o:Weight_kg | Weight (kg) |
| | o:GCS | Glasgow Coma Scale |
| | o:HR | Heart Rate |
| | o:SysBP | Systolic Blood Pressure |
| | o:MeanBP | Mean Blood Pressure |
| | o:DiaBP | Diastolic Blood Pressure |
| | o:RR | Respiratory Rate |
| | o:Temp_C | Temperature (°C) |
| | o:FiO2_1 | Fraction of Inspired Oxygen ($FiO_2$) |
| | o:Potassium | Potassium |
| | o:Sodium | Sodium |
| | o:Chloride | Chloride |
| | o:Glucose | Glucose |
| | o:Magnesium | Magnesium |
| | o:Calcium | Calcium |
| | o:Hb | Hemoglobin |
| | o:WBC_count | White Blood Cell Count |
| | o:Platelets_count | Platelet Count |
| | o:PTT | Partial Thromboplastin Time |
| | o:PT | Prothrombin Time |
| | o:Arterial_pH | Arterial Blood pH |
| | o:paO2 | Partial Pressure of Oxygen ($PaO_2$) |
| | o:paCO2 | Partial Pressure of Carbon Dioxide ($PaCO_2$) |
| | o:Arterial_BE | Arterial Base Excess |
| | o:HCO3 | Bicarbonate ($HCO_3^-$) |
| | o:Arterial_lactate | Arterial Lactate |
| | o:SOFA | Sequential Organ Failure Assessment (SOFA) Score |
| | o:SIRS | Systemic Inflammatory Response Syndrome (SIRS) Score |
| | o:Shock_Index | Shock Index (HR / SBP) |
| | o:PaO2_FiO2 | $PaO_2$/$FiO_2$ Ratio |
| | o:cumulated_balance | Fluid Balance (cumulative) |
| | o:SpO2 | Oxygen Saturation ($SpO_2$) |
| | o:BUN | Blood Urea Nitrogen |
| | o:Creatinine | Creatinine |
| | o:SGOT | Aspartate Aminotransferase (AST/SGOT) |
| | o:SGPT | Alanine Aminotransferase (ALT/SGPT) |
| | o:Total_bili | Total Bilirubin |
| | o:INR | International Normalized Ratio (INR) |
| | o:input_total | Fluid Input (total) |
| | o:input_4hourly | Fluid Input (last 4h) |
| | o:output_total | Fluid Output (total) |
| | o:output_4hourly | Fluid Output (last 4h) |

Figure 11: WIS on 1M training iterations with not trained encoders.

Table 6: Selected hyperparameters for the BCQ training.

| BCQ threshold | Discount | Polyak target update | Target update frequency | Optimizer | Leaning rate |
|---|---|---|---|---|---|
| 0.3 | 0.99 | True ($\tau = 0,01$) | 1 | Adam | 1e−3 |

Table 7: Selected hyperparameters for the behavioral cloning training.

| Epochs | Optimizer | Weight decay | Learning rate |
|---|---|---|---|
| 5000 | Adam | 0.1 | 1e−4 |

(a) Snapshot at the first time step.

(b) Snapshot at the second time step.
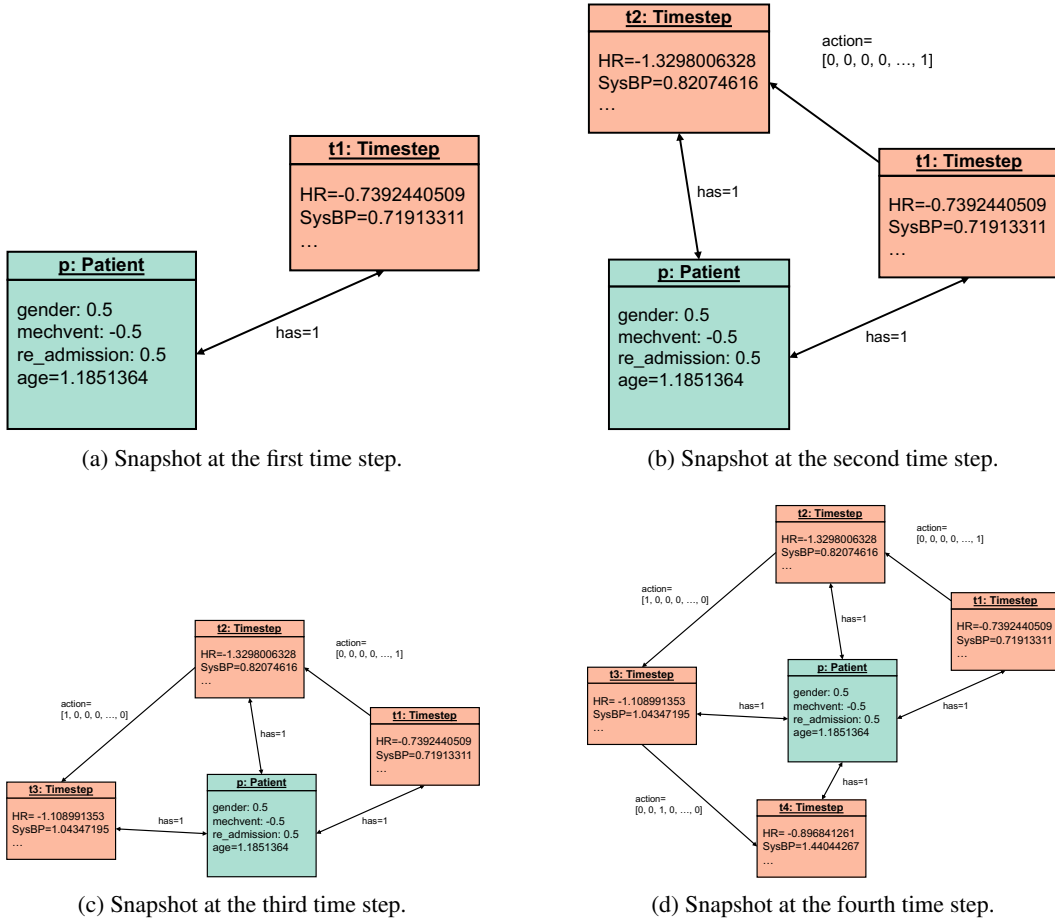
(c) Snapshot at the third time step.

(d) Snapshot at the fourth time step.

Figure 12: Graph evolution over the first four time steps.