

Injection optimization at particle accelerators via reinforcement learning: From simulation to real-world application

Awal Awal^{1,2}, Jan Hetzel², Ralf Gebel^{2,3} and Jörg Pretz^{1,3}

¹*RWTH Aachen University, 52056 Aachen, Germany*

²*GSI Helmholtzzentrum für Schwerionenforschung GmbH, 64291 Darmstadt, Germany*

³*Forschungszentrum Jülich GmbH, 52425 Jülich, Germany*



(Received 21 June 2024; accepted 14 February 2025; published 17 March 2025)

Optimizing the injection process in particle accelerators is crucial for enhancing beam quality and operational efficiency. This paper presents a framework for utilizing reinforcement learning (RL) to optimize the injection process at accelerator facilities. By framing the optimization challenge as an RL problem, we developed an agent capable of dynamically aligning the beam's transverse space with desired targets. Our methodology leverages the soft actor-critic algorithm, enhanced with domain randomization and dense neural networks, to train the agent in simulated environments with varying dynamics promoting it to learn a generalized robust policy. The agent was evaluated in live runs at the cooler synchrotron COSY and it has successfully optimized the beam cross section reaching human operator level but in notably less time. An empirical study further validated the importance of each architecture component in achieving a robust and generalized optimization strategy. The results demonstrate the potential of RL in automating and improving optimization tasks at particle acceleration facilities.

DOI: [10.1103/PhysRevAccelBeams.28.034601](https://doi.org/10.1103/PhysRevAccelBeams.28.034601)

I. INTRODUCTION

The field of accelerator physics has seen a growing interest in leveraging machine learning techniques to enhance the operation and efficiency of particle accelerators. Among the machine learning techniques, reinforcement learning (RL) has emerged as a powerful tool for optimizing complex systems [1].

Particle accelerators are complex machines designed to accelerate charged particles to desired energies for a variety of applications ranging from basic research in physics to applied sciences and medical treatments [2,3]. The cooler synchrotron COSY [4,5], located at Forschungszentrum Jülich in Germany, is an accelerator facility primarily focused on hadron physics research [6–8]. Optimizing the injection process, which involves transferring particles into the accelerator's storage ring, is a nontrivial challenge. This process is critical for ensuring high-quality beam properties, such as intensity and stability, which directly impact the effectiveness and precision of experiments conducted at accelerator facilities [9,10].

Integrating machine learning methods into accelerator facilities has been a subject of increasing interest within the

scientific community. Several studies have demonstrated the potential of machine learning techniques in improving various aspects of accelerator operations. For instance, machine learning methods have been employed for beam diagnostics [11–13], optimization and control tasks [14–17], and anomaly detection [18,19]. These works, among others, highlight the versatility and effectiveness of machine learning approaches in addressing the challenges faced by accelerator facilities. Within the domain of utilizing RL in accelerator operations, several studies utilized RL methods to optimize the beam in simulation, showing their promising potentials [20–22]. Recent studies have managed to run RL agents successfully in live runs, including optimizing a beam with linear settings [23], employing an RL agent that is trained only in simulation to optimize the beam at a real machine [24,25], and in addition optimize a beam with nonlinear settings [25]. These recent methods, however, are either limited to linear optimization problems or low-dimensional optimization in comparison to the injection beam line (IBL) at COSY. Furthermore, achieving operator-level optimization remains a challenging task requiring complex RL methods that cannot be trained in live settings due to their need for a large amount of training. Achieving this level of optimization from simulation training only requires a more comprehensive approach in RL.

The optimization of the injection process in particle accelerators, such as the one at COSY, presents a machine learning challenge that is impractical for traditional supervised machine learning methods. Supervised learning depends on a labeled dataset, where the correct output

Published by the American Physical Society under the terms of the Creative Commons Attribution 4.0 International license. Further distribution of this work must maintain attribution to the author(s) and the published article's title, journal citation, and DOI.

for each input is known and used to train the model. This is not the case in the context of injection optimization as the optimal actions are not readily available or easily computed due to the highly dynamic and complex nature of particle accelerators. In addition, there are no expert samples of optimal actions and it is difficult to build one within the domain of injection optimization because the optimal settings can vary significantly based on the specific experimental setup and objectives. If sufficient expert samples of optimal actions are available then other machine learning methods like imitation learning [26] and inverse RL [27] might be considered.

RL is well suited for this type of optimization tasks as it does not require labeled examples or expert samples. RL agents learn optimal policies through trial and error by interacting with the environment and are guided by a reward function that reflects the quality of the actions based on their outcomes. This approach allows RL agents to autonomously discover and refine strategies for optimizing the injection process [28].

Numerical optimization methods such as Bayesian optimization and evolution strategies are popular methods within the domain of optimizations in accelerator facilities due to their flexibility and adaptability [14,29–31]. These methods are easier to implement and are well suited for scenarios, where a simulation model is not available or when quick adjustments are needed for different optimization targets. However, they tend to require more time to converge with less consistency, especially in high-dimensional cases. On the other hand, RL, while requiring more initial investment in terms of computational resources and time, provides efficient, fast, and consistent optimization performance through targeted exploration by learning the system dynamics.

Building on our previous work and interest [29], this research introduces a framework for the application of RL in accelerator facilities and applies it to optimize the injection process at COSY. By framing the optimization challenge as an RL problem, we train an agent to make data-driven decisions that improve the injection efficiency. The RL agent is capable of adapting to varying conditions, learning from interactions with the environment to optimize the beam’s properties dynamically. The actions of RL agents are aligned and tailored to the specific domain they are trained on. This is because they build an internal concept of their environment and the potential consequences of their actions [32]. This paper is based on the thesis [33] which contains a more detailed description.

II. BEAM AND INJECTION OPTIMIZATION AT THE COOLER SYNCHROTRON COSY FACILITY

COSY accelerates protons or deuterons and is equipped to handle polarized and unpolarized ions. The facility, including the cyclotron and the IBL, is used in a wide range of experimental research. The IBL at COSY is specifically

designed for the transfer of negatively charged hydrogen and deuteron ions from the cyclotron JULIC [34] to the synchrotron. At the point where the ions are injected into COSY, the *injection point*, the electrons are stripped off via a stripping foil. During injection, the acceptance of COSY is filled with particles via multiturn phase space painting. Subsequently, the optimal region in phase space for the incoming beam during injection will be called *injection acceptance*. The IBL’s functionality is crucial, as it directly influences the intensity and quality of the beam delivered to COSY.

A. The injection beam line

The IBL at COSY, spanning a length of 94.15 m, is a transfer beam line designed for the efficient transfer of ions, see Fig. 1. It is logically partitioned into eight sections consisting of 15 quadrupole magnet families and 28 steerers, among other components, to guide and shape the beam into COSY. The IBL’s design allows for precise control over the beam, which is essential for achieving optimal performance in beam injection and experimentation. Figure 1 illustrates the layout of the COSY facility, highlighting the IBL and its various sections.

Section 8, the last section of the IBL, is composed of four quadrupoles and seven steerers and has high importance, see Fig. 2. The task of this section is to match the transferred beam from the preceding sections with the injection acceptance of COSY. The focus of the RL approach is to optimize this section due to its vital and direct role in controlling the final phase of the beam injection process and the subsequent impact on the beam intensity and quality inside COSY.

B. Injection optimization

The optimization of the IBL is a complex nonlinear problem, primarily due to the multidimensional nature of beam dynamics and the complex design of the IBL. The objective is to maximize the beam intensity inside COSY while ensuring minimal setup time. The optimization process involves adjusting the parameters of the IBL components, particularly the strengths of the quadrupole and steerer magnets, to fine-tune the beam’s trajectory and properties. The ultimate goal of this optimization is to align the phase space of the injected beam with the injection acceptance of the storage ring to ensure a stable and high intensity beam. A key challenge in the optimization is, that the beam parameters at the beginning of the IBL vary. This is in particular related to the cyclotron and its magnetic field, which drifts over time and thus influences the beam. This can even lead to beams where the cross section differs from a shape which can be estimated by a Gaussian and offers multiple clusters of particles instead. Given the complexity of this system which only can be set up by experts, a common mode of operation at COSY is to treat the beam from the cyclotron as a black box and make

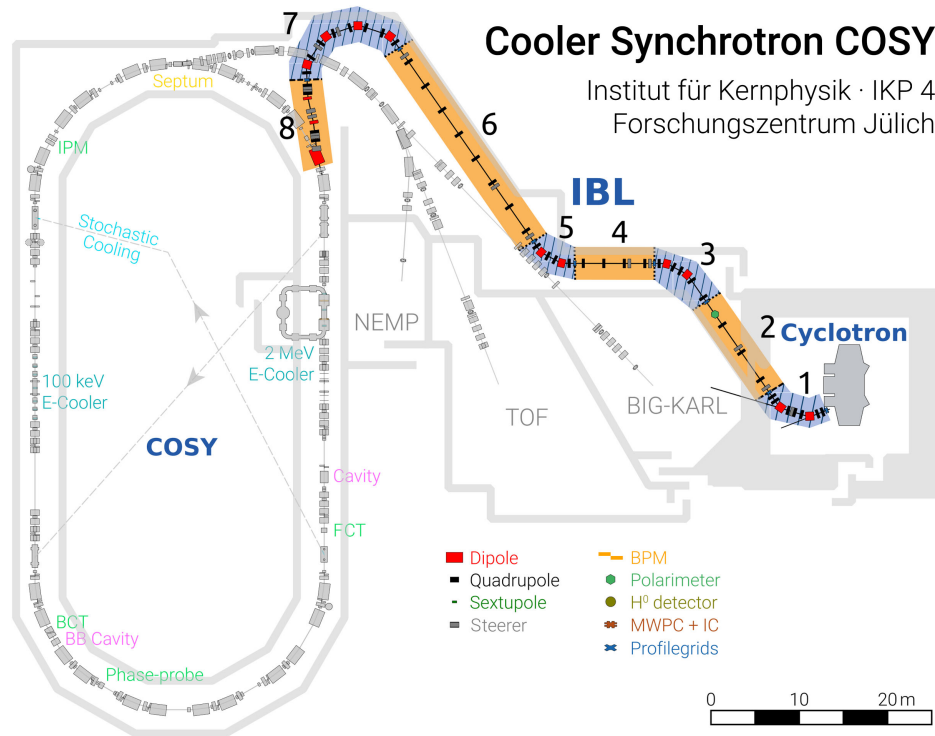


FIG. 1. The COSY facility at Forschungszentrum Jülich. Illustrated are the cyclotron (right), the cooler synchrotron COSY (left), and the interconnecting injection beam line (IBL). For the latter, its division into sections is indicated with colors and numerals.

adjustments to the IBL to compensate for the unknown beam parameters. Currently, the optimization process involves adjusting the magnets manually to maximize the beam current inside COSY, a process that is time consuming and lacks consistency. Automating this process is of high significance to ensure minimal setup time and higher availability for the experiments. An alternative approach is to measure the parameters of the beam at the beginning of the IBL, as described in [35]. As the measurement of the parameters to the necessary precision is lengthy as well, this approach is usually not followed in standard operation.

The automation of the injection optimization, as described in the following, is based on the assumption that the optimal beam occupation of the phase space at the

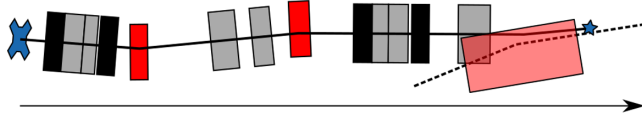


FIG. 2. A sketch of section 8 of the IBL showing the central path and the different components along the beam line. Dipoles are colored in red, quadrupoles are colored in black, and steerers are colored in gray. The injection point is marked with a blue star, where the charge exchange foil is placed. The trajectory of the p^- beam through the IBL is indicated by the solid black line. The dashed line indicates the central orbit through the synchrotron COSY.

charge exchange foil for injection is known and can be restored by an appropriate setting of the IBL. In the topical approach the task of finding this setting is divided in two subtasks: The first is to transport the beam from the beginning of the IBL to the beginning of section 8 without losing intensity. A possible approach to achieve this is Bayesian optimization as it is described in [29]. The second task is then to use the elements in section 8 to match the desired parameters at the injection point of COSY. Our approach is to use an RL agent to carry out the latter task. As in this approach the beam parameters at the beginning of the IBL are unknown and the first step only ensures full transmission of the beam, the beam parameters of the beginning of section 8 are to be regarded as unknown as well. Dividing the optimization in two tasks is chosen to reduce the number of free parameters for the RL agent to a reasonable amount while keeping the necessary degrees of freedom for a successful optimization. At COSY, we have the opportunity to record the position and the cross section of the beam at the end of the IBL with a fluorescent screen. As this gives access to the spacial components only, the goal of the RL agent is to match the transverse position and spread of the beam at this location to an operator given specification. For testing and demonstration of the capabilities of the RL agent the specification may differ from the optimal settings for injection. The RL agent is trained in simulation only prior to its application to the real machine. The experimental demonstration includes the described

matching to the user's specifications as well as an investigation of the influence of several architecture components on the result.

III. THEORETICAL BACKGROUND

In this section, we introduce the theoretical foundations of reinforcement learning, along with the concepts of partially observable Markov decision process (POMDP) and domain randomization. These concepts are essential for understanding the RL framework applied in this research. A comprehensive introduction to RL can be found in [36] and a detailed description is given in [28].

A. Reinforcement learning

Reinforcement Learning (RL) is the core principle of our approach to optimize the IBL at COSY. The general setup of RL involves an agent interacting with its environment at a sequence of time steps to maximize cumulative rewards. In this setup, the environment's state at each time step t is represented by $s_t \in \mathcal{S}$. Assuming full observability of the state, the agent's policy $\pi(a|s)$ dictates the probability distribution over actions $a \in \mathcal{A}$ given a state s . The agent then receives a reward signal $r_t = r(s_t, a_t)$ from the environment, providing feedback on the desirability of the action. The agent's objective is to maximize the expected return over a horizon, which is the sum of discounted rewards obtained during an episode, expressed as

$$R_t = \sum_{k=t}^T \gamma^{k-t} r_k, \quad (1)$$

where $\gamma \in [0, 1]$ is the discount factor which balances the immediate and future rewards and T is the horizon of each episode.

The state-value function $V^\pi(s)$ is the expected return over the horizon, starting from state s and following a policy π . The state-value function $V^\pi(s)$ is defined as

$$V^\pi(s) = \mathbb{E}_{\tau \sim p(\tau|\pi)} \left[\sum_{k=t}^T \gamma^{k-t} r_k | S_t = s \right], \quad (2)$$

where $\tau = (s_0, a_0, s_1, \dots, a_{T-1}, s_T)$ denotes a trajectory through the state-action space and $p(\tau|\pi)$ is the trajectory's probability under policy π , and the discount factor γ balances the immediate and future rewards. The action-value function $Q^\pi(s, a)$, also known as the Q -function, representing the expected return when starting in state s and taking action a then following policy π , is closely related to $V^\pi(s_t)$ and is defined as

$$Q^\pi(s, a) = \mathbb{E}_{\tau \sim p(\tau|\pi)} \left[\sum_{k=t}^T \gamma^{k-t} r_k | S_t = s, A_t = a \right]. \quad (3)$$

The state-value function $V^\pi(s_t)$ can be expressed in terms of the action-value function $Q^\pi(s, a)$ as

$$V^\pi(s) = \mathbb{E}_{a \sim \pi} [Q^\pi(s, a) | S_t = s]. \quad (4)$$

$J(\pi)$ represents the expected return by following a certain policy π . This measure indicates how good a policy is and maximizing it results in improving the policy to achieve the optimal policy π^* . The expected return of a policy $J(\pi)$ is defined as

$$J(\pi) = \mathbb{E}[R_0 | \pi] = \mathbb{E}_{\tau \sim p(\tau|\pi)} \left[\sum_{k=0}^T \gamma^k r(s_k, a_k) \right]. \quad (5)$$

The next state is dictated by the state-transition function $P(s_{t+1} | s_t, a_t)$ which is determined by the dynamics of the environment [28]. Contemporary methods in RL utilize a key concept in RL which is the Bellman equation [37]. It provides a recursive decomposition for the value function. The Bellman equation for $Q^\pi(s)$ is given as

$$Q^\pi(s, a) = \mathbb{E}[r(s, a) + \gamma V^\pi(s_{t+1}) | S_t = s, A_t = a]. \quad (6)$$

The objective in RL is to find the optimal policy π^* that maximizes the expected return from any initial state. This is often framed as maximizing the state-value function V^π or the action-value function Q^π for all states $s \in \mathcal{S}$. Mathematically, the objective is to find π^* such that

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{s \sim \mathcal{S}} [V^\pi(s)]. \quad (7)$$

Table I shows an overview of the variables used in the equations and their definitions.

TABLE I. Definitions of variables used throughout the paper.

| Variable | Meaning |
|-----------|--|
| t | Time step |
| τ | Trajectory |
| a | Action |
| s | State |
| o | Observation |
| h | History |
| π | Policy |
| γ | Discount factor |
| r | Reward |
| R_t | Future expected rewards from time step t |
| $J(\pi)$ | Expected return of a policy π |
| $Q(s, a)$ | Action-value function |
| $V(s)$ | State-values function |
| ρ | Environment dynamics vector |
| g | The goal of the optimization |

B. Partially observable Markov decision processes

RL methods are primarily designed to solve decision-making problems formulated as Markov decision processes (MDPs), where the environment setup is inherently Markovian [38]. An MDP is characterized by a set of states \mathcal{S} , a set of actions \mathcal{A} , a transition function $P[s_{t+1}|s_t, a_t]$ that determines the probability of transitioning from state s_t to state s_{t+1} after taking action a_t , and a reward function $r(s_t, a_t)$ which assigns rewards to state-action pairs. The fundamental assumption in an MDP is that the current state encapsulates all necessary information for decision making, implying that the future state s_{t+1} is conditionally independent of past states given the current state s_t and action a_t . Formally, this is expressed as

$$P[s_{t+1}|s_t, a_t] = P[s_{t+1}|s_t, a_t, s_{t-1}, a_{t-1}, \dots, s_0, a_0]. \quad (8)$$

Partially observable Markov decision processes (POMDPs) extend the MDP framework to scenarios with partial observability [39]. In POMDP, at each time step t the agent receives an observation $o_t \in \Omega$, which provides partial information about the actual state s_t . The challenge in POMDPs lies in the agent's need to infer the hidden state of the environment from the history of observations and actions. This is a significant complication over the full observability assumption in MDPs, where the current state encapsulates all necessary information for decision making.

To address the challenges posed by partial observability, strategies often involve utilizing the history of observations and actions or an encoding of this history to approximate the unobservable states, effectively transforming the problem back into a solvable MDP. The agent observes at each time step an encoding of the history h_t , from which it needs to infer the hidden state of the environment. The history encoding h_t can be expressed as

$$h_t = f(s_{t-1}, a_{t-1}, \dots, s_0, a_0), \quad (9)$$

where f represents a function that encodes the history of states and actions into a format that the agent can use for its decision-making process. This encoding can be the stacking of a predefined number of the most recent observations and actions [40,41], or it can be achieved using a recurrent neural network [42,43], which is capable of maintaining and updating internal state representations based on the sequence of observations and actions.

C. Domain randomization

The challenge of transferring the trained machine learning models, particularly in RL, from simulation to the real world is a major challenge in robotics and control systems [44]. This challenge, often referred to as the sim-to-real transfer problem, arises because the transition probabilities in a simulated environment, $P_{\text{sim}}[s_{t+1}|s_t, a_t]$, do not perfectly match those in the real world, $P_{\text{real-world}}[s_{t+1}|s_t, a_t]$.

The discrepancy between these transition probabilities usually leads to a model that performs well in simulation but fails to generalize in real-world conditions.

Domain randomization is a technique designed to address this discrepancy by training the agent in a variety of simulated environments with randomly altered physics parameters, sensor noise, and other environmental conditions. The core idea is to expose the agent to a wide range of possible conditions during the training phase, which helps in improving its ability to generalize from the simulated environment to the real world [45,46]. In standard training, an agent is trained under a fixed dynamics vector $\rho \in \mathcal{P}$, where \mathcal{P} represent the space of all possible dynamics parameters. In domain randomization, for each training episode i , a dynamics vector ρ_i is sampled from a predefined distribution over \mathcal{P} , i.e., $\rho_i \sim P(\mathcal{P})$. The environment dynamics for that episode are then defined by ρ_i , and they remain constant for the duration of the episode. This process can be formalized as follows:

$$\rho_i \sim P(\mathcal{P}), \quad \forall i \in \{1, 2, \dots, N\}, \quad (10)$$

where N is the number of episodes and $P(\mathcal{P})$ is the probability distribution over the dynamics space \mathcal{P} . The goal of domain randomization is to train an agent such that the real-world environment appears as another sample from the distribution $P(\mathcal{P})$. This approach effectively broadens the distribution of the simulation to incorporate a wider range of real-world variations. The agent, therefore, learns a generalized policy π that is robust across a variety of dynamics.

IV. METHODOLOGY

The optimization of the injection process at COSY currently involves manual adjustments of all sections (1 to 8). This manual method aims to enhance the beam current within COSY but is time consuming and often results in inconsistent beam properties within the storage ring.

A better alternative strategy is to focus on optimizing the beam's phase space at the point of injection, rather than merely aiming to increase the beam current inside the storage ring. The objective is to align the phase space of the incoming beam with the storage ring's injection acceptance to ensure a consistently high-quality beam. In this research, the agent is assigned to manipulate the 11 magnets (4 quadrupoles and 7 steerers) of section 8 to optimize the transverse space of the beam at the injection point. The settings of sections 1–7 can be chosen independently, as long as 100% transmission to the end of section 7 is reached. A possible method to optimize these sections is discussed in [29]. While acknowledging the importance of the angle of the injected beam, it was excluded from the optimization process due to limitations with the IBL

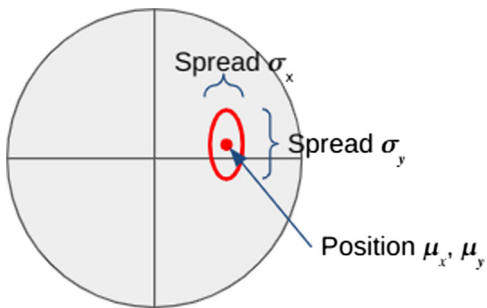


FIG. 3. The optimization goal of the agent is the beam’s cross section at the injection point and is set by operators. It is defined by the center of the beam (μ_x, μ_y) and its spread (σ_x, σ_y) .

sensors at COSY. The direct feedback on the beam’s characteristics is provided by a camera positioned at the end of the IBL.

A. Explicit goal

The goal $g \in \mathbb{R}^d$ of the optimization is explicitly dictated by the operators. We set the goal as the cross section (position and width) of the beam at the injection point. Therefore, the goal $g \in \mathbb{R}^4$ is defined as the beam parameters corresponding to the beam center (μ_x, μ_y) and its spread (σ_x, σ_y) at the injection point as observed through the camera, see Fig. 3. This requires the value function $V(s)$ to generalize not only over states, but goals too as a universal value function approximator $V(s, g)$ [47]. This adjustment demands the introduction of a goal-oriented reward function, $r(s, a, g)$, and modifies the agent’s policy correspondingly to $\pi(a|s, g)$. This modification enables the learning of a universal policy adaptable to various target configurations at the injection point by presenting the agent with a randomly sampled goal g at the beginning of each episode. The agent’s task is to manipulate the magnets to match the observed beam characteristics with the desired goal, thereby optimizing the transverse position and width of the beam at the point of injection.

B. Randomized simulation dynamics

To train the policy to perform under varying real-world dynamics, the concept of domain randomization is employed in simulation during the training phase. The goal to be achieved here is to train the agent across a multitude of simulated environments, each characterized by distinct physics parameters, sensor noise levels, and environmental factors. This is realized by adjusting the training environment’s dynamics through a set of parameters ρ which are varied at the start of each training episode but remain static throughout the same episode. In our experiment, the dynamics vector ρ consists of (i) initial bunch phase space parameters, (ii) initial beam angle, offset, and number of clusters, (iii) magnets’ strength, (iv) initial magnets values, and (v) offset to magnets reading value.

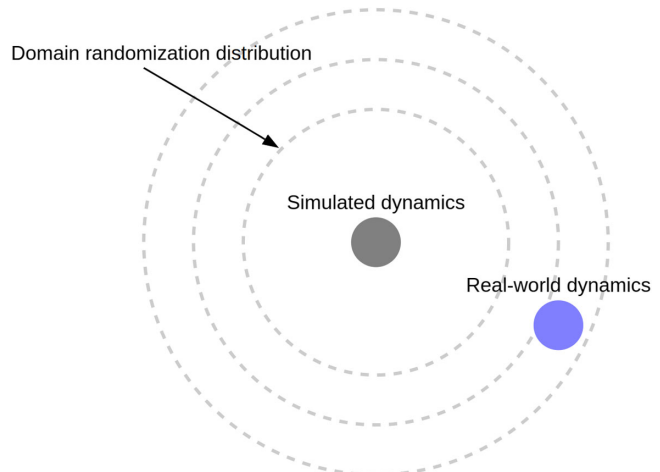


FIG. 4. Illustration of the domain randomization concept. Usually there is discrepancy between the simulated dynamics and the real-world dynamics shown here as the gray point and the blue point, respectively, in the space of all possible dynamics. By expanding the distribution of dynamics from which the simulated dynamics are sampled, the trained agent becomes more robust. If the domain randomization is sufficient, the real-world dynamics appear to the agent as another sample of the simulated environments [33].

When the agent is trained under sufficient domain randomization and can perform the optimization in simulation successfully, the agent is expected to perform well in the live run. The reason is, for the agent, the live IBL looks as another sample of the environment dynamics it encountered during the training. If the domain randomization is not sufficient, then the defined distribution over the dynamics space should be increased as shown in Fig. 4. Formally, this is correct when the defined distribution over the dynamics space $P(\mathcal{P})$ includes the dynamics of the live IBL $\rho_{\text{live-IBL}}$

$$\rho_{\text{live-IBL}} \in P(\mathcal{P}). \quad (11)$$

In addition, a Gaussian noise $\mathcal{N}(0, \sigma^2)$ is added to the observations which are sampled at each time step of the training. This noise enhances the agent’s ability to operate under realistic conditions by simulating the measurement inaccuracies and environmental noise. By exposing the agent policy to a broad range of dynamic conditions, it learns to generalize its performance by learning an internal concept of the environment, improving its applicability to the actual conditions of the IBL at COSY.

C. Reinforcement learning agent

The RL agent used in this study is based on the soft actor-critic (SAC) algorithm [48], a model-free and off-policy actor-critic method [49,50] that is effective in environments requiring continuous action decisions. SAC is distinguished by its incorporation of entropy into the

reward structure, promoting exploration by the agent and preventing premature convergence to suboptimal policies. Based on the earlier discussions, explicit goals g , encoded history h , and simulation dynamics ρ are incorporated into the policy $\pi_\phi(a|s, h, g)$ and the Q -function $Q_\theta(s, h, g, a, \rho)$ which are parameterized using neural networks. Here ϕ and θ denote the weights of the policy network and the Q -function network, respectively.

The SAC algorithm optimizes the stochastic policy during off-policy learning, which means it can learn from experiences generated by any policy, not just the current one. The policy network π_ϕ aims to maximize the expected return as well as the entropy of the policy to encourage exploration. The objective function for the policy network, parameterized by ϕ , is given by

$$J(\pi_\phi) = \mathbb{E}_{(s,h,g) \sim \mathcal{D}, a \sim \pi_\phi} [Q_\theta(s, h, g, a, \rho) - \alpha \log \pi_\phi(a|s, h, g)], \quad (12)$$

where α is the temperature parameter that determines the relative importance of the entropy term against the reward, and \mathcal{D} represents the experience replay buffer. The environment dynamics ρ and the goal g are fixed throughout the episode while the others are time step dependent.

The Q -function network $Q_\theta(s, h, g, a, \rho)$ is optimized to improve the computed predicted returns, also called the Q value, by minimizing the difference between the predicted Q value and the target Q value, a better estimate than the predicted one. The target Q value is computed as

$$y = \mathbb{E}_{(s,h,g,a,r,s',h') \sim \mathcal{D}} [r + \gamma(Q_\theta(s', h', g, a', \rho) - \alpha \log \pi_\phi(a'|s', h', g))], \quad (13)$$

where γ is the discount factor, and unlike the next state s' and next encoded history h' , which are sampled from the replay buffer \mathcal{D} , the next action a' is computed using the current policy network π_ϕ .

The Q -network is trained similar to supervised learning using the computed target Q value. The loss function for optimizing the Q -function network is expressed as

$$\mathcal{L}_Q(\theta) = \mathbb{E}_{(s,h,g,a) \sim \mathcal{D}} \left[\frac{1}{2} (Q_\theta(s, h, g, a, \rho) - y)^2 \right]. \quad (14)$$

The SAC algorithm iteratively updates the policy π_ϕ and the action-value function Q_θ using samples drawn from the experience replay buffer \mathcal{D} . The replay buffer allows the agent to learn from a diverse set of experiences, leading to a robust and generalized policy that can effectively adapt to the variant conditions under the randomized simulation dynamics.

The Q -function network is used only during the training process and is discarded later during the live run. During the

training, it implicitly learns the model of the environment and its structure of rewards. This knowledge is then used to compute the gradients with respect to the action that results in a higher reward. Passing the environment dynamics ρ to the action-value function can improve the training by speeding up the convergence of the action-value function and might improve the performance of the agent by computing better gradients of the expected value with respect to the action. The environment dynamics ρ can be, nevertheless, omitted from the action-value function as they can be inferred from the encoded history h . This can extend the training time, but the learned policy of the agent is eventually the same.

D. Dense neural networks

The choice of the neural network architecture can have a crucial role in the performance of RL agent. Sinha *et al.* [51] showed that better benchmarks can be achieved by incorporating the state into the inputs of each layer of the policy neural network and the state-action into the inputs of each layer of the action-value neural network. We adopted an architecture that closely resembles the original DenseNet architecture [52], which enhances the flow of information and gradients throughout the network by connecting each layer of the neural network to all earlier layers. This is achieved by computing the inputs \mathbf{x} of each layer k as the concatenation of the nonlinear outputs and the inputs of the previous layer as follows:

$$\mathbf{x}_k = [\sigma(\mathbf{z}_k); \mathbf{x}_{k-1}], \quad (15)$$

where σ is the nonlinear operation, \mathbf{z} is the linear transformation of the previous layer, and \mathbf{x}_0 are the initial inputs to the neural network. Another way to perceive this architecture, as shown in Fig. 5, is that the features passed to each layer are the concatenation of the features generated from all earlier layers in addition to the input features. This architecture facilitates a more effective learning process by ensuring that both low-level features and high-level representations contribute directly to the output of each layer, thus enabling the RL agent to distinguish and utilize complex patterns in the environment more efficiently [53].

E. Reward function

The reward function is a crucial component of the reinforcement learning framework. It guides the agent toward achieving the optimization goal and it can also reinforce or suppress certain behaviors. It evaluates the performance of the agent at each time step t and provides a scalar feedback signal. The reward function used in our approach is composed of several parts, designed to encapsulate different aspects of the beam optimization process.

First, we define a term to quantify the accuracy of the beam's transverse space alignment with the desired parameters at the injection point. This is achieved by computing

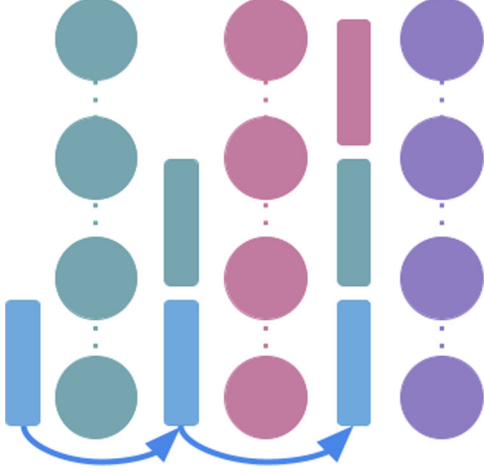


FIG. 5. Illustration of the structure of dense neural networks. The colored circles are the neurons and the color matching rectangles are the features generated by the neurons while the blue rectangle represents the input features. Each layer receives the combination of all earlier features in addition to the input features [33].

the difference between the measured beam parameters at the injection point $(\mu'_x, \mu'_y, \sigma'_x, \sigma'_y)$ and the target parameters $(\mu_x, \mu_y, \sigma_x, \sigma_y)$ using the *softplus* function to ensure a smooth and differentiable measure:

$$\text{transverse}_t = \text{softplus} \left(\sum_{i \in \{x, y\}} [|\mu_i - \mu'_i| + |\sigma_i - \sigma'_i|] \right), \quad (16)$$

where the softplus function is defined as $\text{softplus}(x) = \log(1 + e^x)$. This term aims to minimize the discrepancy between the actual and desired beam parameters, encouraging precise control over the beam's alignment and shape at the point of injection. However, merely using this measure, especially in a complex IBL with high dimensionality, can lead to undesired behavior, where the agent loses some or most of the beam but still manages to match the desired transverse space. In our earlier experiments, this undesired behavior was present when only this measure was used. To enforce maintaining the efficiency of beam transmission through the IBL, we introduce a measure that evaluates the fraction of the beam successfully transmitted through the IBL:

$$\text{transmission}_t = \frac{1}{2n} \sum_{i=1}^n [\hat{d}_i + \mathbb{1}_{(\hat{d}_i=1)}], \quad (17)$$

where d_i represents the distance traveled by the i th simulated particle, \hat{x} is the scaled variable of x to be within the range $[0, 1]$, n is the total number of particles, and $\mathbb{1}_{(\text{condition})}$ is the indicator function, evaluating to 1 when *condition* is met and 0 otherwise. The precision term

precision_t combines the transverse alignment and transmission efficiency, reflecting the overall optimization performance at each time step:

$$\text{precision}_t = (1 - \widehat{\text{transverse}}_t) \cdot \text{transmission}_t. \quad (18)$$

This measure serves as the primary component of the reward for the agent, promoting the alignment of the beam's transverse space with the desired parameters while ensuring efficient transmission. To discourage unnecessary adjustments and promote stability, a penalty term is introduced. This term penalizes changes in magnet settings that result in a decrease in optimization performance, incorporating a constant k to adjust the penalty's severity:

$$\text{penalty}_t = k * \text{precision}_t * \mathbb{1}_{(m_t \neq m_{t-1})} \mathbb{1}_{(\text{precision}_t < r_{t-1})}, \quad (19)$$

where m_t represents the magnet settings at time t . The final reward at time t , r_t , is then calculated as the precision measure scaled to be negative [54,55] and adjusted for any penalties incurred:

$$r_t = \text{precision}_t - \text{penalty}_t - 1. \quad (20)$$

This reward structure is designed to finely balance the trade-off between the accuracy of beam parameters, efficiency of transmission, and the cost of adjustments, guiding the agent toward achieving the optimal set of actions for beam injection optimization at COSY.

V. IMPLEMENTATION AND TRAINING

The implementation and training of our RL framework for optimizing the injection process at COSY is based on a simulation environment that mirrors different variants of the real-world dynamics of the IBL. This simulation is the foundation for training and evaluation of the RL agent, enabling it to learn and adapt to the complexities of particle acceleration and beam optimization. At each step of the optimization run, the RL agent is required to determine the adjustments to the magnets based on its current policy. The adjustment is represented by the action a_t at time step t , which directly influences the current values of the magnets m_t . The new magnet values are calculated using the equation:

$$m_t = m_{t-1} + a_t, \quad (21)$$

where m_{t-1} represents the magnet values at the previous time step. This formulation of how the new magnet values m_t are updated reflects the practical approach used by operators who use adjustments and incremental changes to manage the dynamic and interdependent nature of the magnets. This approach provides incremental and controlled modifications that ensures maintaining stability and performance without abrupt changes that could lead to

beam loss. Additionally, using adjustments allow the RL agent to respond flexibly to the constantly changing conditions and feedback from the system.

It is crucial to maintain the magnet values within operational limits to ensure the safety and integrity of the accelerator’s components. Therefore, the magnet values are constrained to operate within the range of $[-0.7, 0.7]$, corresponding to $\pm 70\%$ of their maximum current capability. This constraint not only ensures the operational safety of the IBL but also reflects the physical limitations and operational protocols of COSY. Training within these constraints allows the RL agent to develop strategies that are viable for live deployment, ensuring that the optimization strategies are both effective and practical for real-world application.

A. Simulation environment

The simulation for training the agent was conducted using MAD-X [56], a comprehensive tool for designing and simulating particle accelerators. A virtual environment was created following the standard Gym interface to simulate the dynamics of the IBL at COSY [57,58]. Particles are simulated and tracked via the Polymorphic Tracking Code (PTC) for more comprehensive analysis [59,60]. This environment simulates the physics of beam propagation through the IBL and models the effects of the quadrupole and steerer magnets on the beam’s trajectory and characteristics.

B. Reinforcement learning implementation

The RL model consists of two primary components: the policy network π and the action-value network Q . Both networks are implemented with dense layers, each with two hidden layers of 512 neurons. The input to the policy network includes the values of the quadrupoles and steerers in the last section of the IBL, computed statistics from the camera output $(\mu_x, \mu_y, \sigma_x, \sigma_y)$, the last seven observations of magnet values and their corresponding beam statistics, and the target beam parameters. The Q function receives the same inputs as the policy, with the addition of the environment parameters ρ , representing the dynamics of the simulation.

The training process utilized eight parallel environments to enhance learning stability and efficiency, requiring eight CPUs for environment simulation and an additional GPU for the agent’s processing and neural network training. At each step, a beam simulation with 1000 particles was performed. Each episode consisted of 32 steps, with the action of the agent having a maximum change of 15% in the magnet values per step. The discount factor γ was set to 0.95, reflecting an emphasis on future rewards. The neural networks are optimized via ADAM optimizer [61] with a learning rate of 0.0003. While training was initially set for 10,000 epochs, it was typically terminated earlier upon convergence. For the empirical analysis conducted during the development phase, agents were trained using four

parallel environments, utilizing four CPUs for environment simulation and an additional CPU for agent processing.

C. Live run setup

For the live run, an interface to the IBL was created using the same Gym environment interface, allowing for a seamless transition from simulation to real-world application. The IBL was autonomously controlled using the EPICS control system [62–64], enabling real-time adjustments to the magnet settings based on the agent’s decisions.

A fluorescent screen and a camera were deployed at the end of the IBL to provide direct feedback on the beam’s transverse space. Live image processing was implemented using the area-detector module of the EPICS control system [65]. The camera, operating at a 100 ms frame rate, captures the beam’s cross section at the injection point and provides real-time data on the beam’s characteristics. However, due to the release of particles in 20 ms bunches, the observed values of σ_x and σ_y are notably noisy, presenting an additional challenge for the optimization task.

The agent was evaluated in several different optimization tasks, each requiring adjustments to the beam parameters to meet specific targets. For the live optimization runs, the agent was limited to 16 steps per task.

VI. EXPERIMENTS AND RESULTS

The experiments conducted aimed to validate the effectiveness of the proposed RL approach in optimizing the IBL at COSY. The successful training of the agent across a diverse set of environmental dynamics in simulation, enables it to efficiently optimize the IBL during live runs. For the agent, the real-world IBL is presented as another variant of the various environment dynamics it encountered during the training and can be optimized via the learned robust policy of the trained agent.

A. Live run performance

In the live runs, the agent was tasked with optimizing the beam’s central position and its spread to meet a set of predefined tasks. These tasks involved positioning the beam at various locations on the fluorescent screen, including the center and the edge, while changing the horizontal and vertical spreads of the beam between 2 and 5 mm. The performance of the agent was quantified using the mean L2 error for the beam’s central position (μ) and the mean L1 error for its spread (σ). The agent successfully optimized the IBL in the live runs achieving a mean L2 error of 0.19 mm for the beam’s central position and an average L1 error of 0.7 mm for its spread while maintaining a 100% transmission of the incoming beam. Figure 6 shows a photograph of the beam as it propagates through the fluorescent screen during a live run optimization.

The performance of the RL agent was benchmarked against that of a human operator to evaluate its optimization

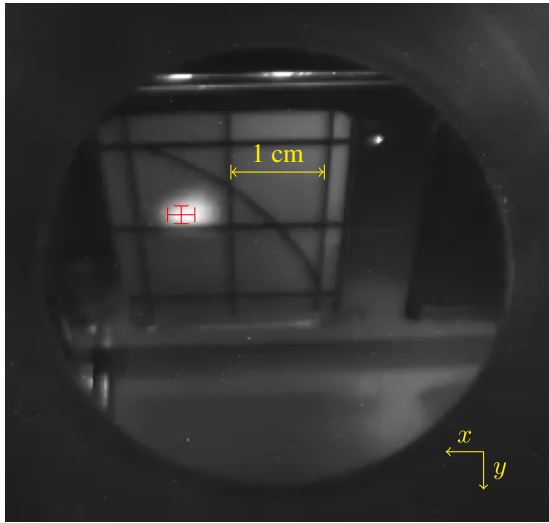


FIG. 6. A picture of the beam, illuminating through the fluorescent screen, during a real-time optimization by the RL agent. This screenshot was taken during the RL agent optimizing the beam at step 7 of the full 16 steps. The target in this instance for the beam cross section is shown in red and has the coordinates $\mu_x = 5.5$ mm, $\mu_y = -1.8$ mm, $\sigma_x = 3$ mm, and $\sigma_y = 2$ mm.

efficiency. Both the agent and the operator were given a set of optimization tasks to adjust the position of the beam using the steerers only. The operator and the agent achieved a similar optimization accuracy, see Fig. 8. However, the agent completed the optimization tasks in around 17 min while the human operator required around 1 h to complete the same tasks. While these tasks were limited to controlling the steerers, the positioning of the magnets within the last section resulted in additional complexity. This further complexity is mainly due to the placement of the last steerers before the last dipole and quadrupole resulting in nonlinear dependency of the position on the tuning.

Figure 7 provides a visual representation of the optimization process in action, showing the dynamic adjustments made by the agent to the steerers' current during a single optimization task. This visualization illustrates the agent's strategic approach to exploration and optimization. Initially, the agent takes targeted exploration actions and gathers information about the environment from the responses. As the optimization progresses, the agent utilizes this acquired knowledge to make more informed decisions, gradually refining its actions toward the optimal settings. This pattern of behavior illustrates the agent's ability to adapt and optimize in a complex and dynamic environment.

These results emphasize the potential of RL in enhancing the operational efficiency of particle accelerators. The degree of accuracy achieved by the agent in aligning the beam's parameters with the target values demonstrates the feasibility of employing RL for complex optimization tasks in particle acceleration facilities.

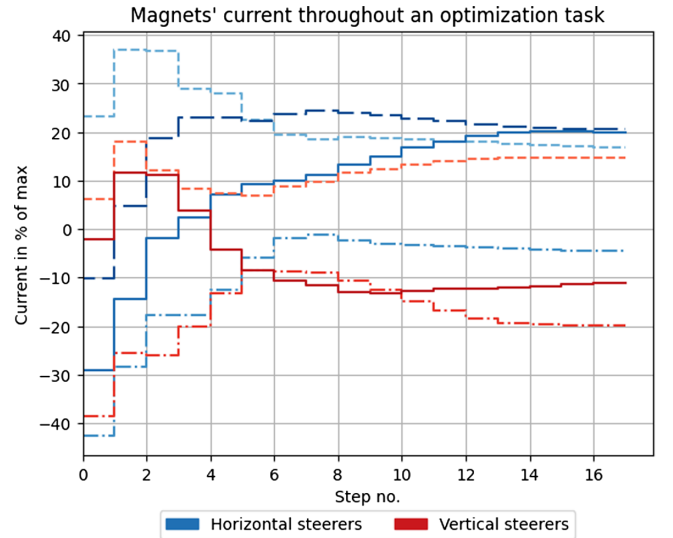


FIG. 7. The evolution of the steerers' current throughout a single optimization task induced by the agent actions. This plot illustrates a common optimization pattern, where the agent starts with targeted exploration actions and gradually converges to the optimal values as it refines its domain knowledge. Each step is 6 s long.

B. Empirical study on the architecture components

To further evaluate the impact of the architecture components on the agent's performance, an empirical study was conducted. In this study, the task was limited to controlling the beam location by optimizing the seven steerers in section 8 of the IBL to meet a series of five predefined tasks. The performance of several replicate agents, each lacking a different component of the architecture, was analyzed and compared. The examined components were the dense layers, observation noise, history, and domain randomization. The results of this study are summarized in Fig. 8.

The removal of domain randomization had the most significant impact on performance, indicating the critical role of training under varied dynamics for generalizing to real-world conditions. This reinforces the importance of domain randomization in preparing the agent for the complexities and variabilities of the live environment at COSY. The absence of observation noise, while less impactful, still resulted in a noticeable degradation in performance, underlining the importance of training under conditions that mimic the real operational environment, including sensor noise and measurement inaccuracies.

Without the inclusion of history in the agent's observation space, its ability to infer the environment's dynamics and predict the outcome of actions under partial observability was notably hindered. This limitation is evident in the increased error and illustrates the importance of temporal information for making informed decisions in environments, where the current state alone does not

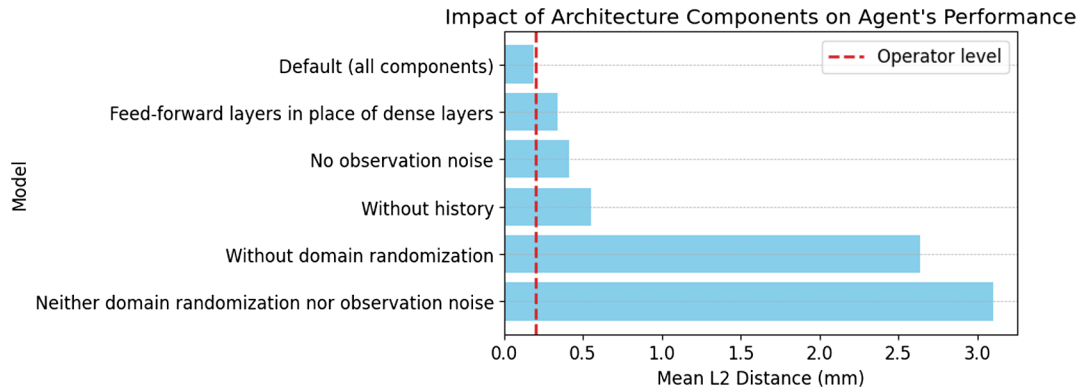


FIG. 8. Empirical analysis of the impact of architecture components on the agent’s performance. The table presents the mean L2 distance (in mm) between the optimized beam’s central position and the target position under various architectural configurations. The “Default” configuration includes all components: dense layers, observation noise, history, and domain randomization. Each subsequent bar shows the performance impact when a specific component is removed or altered, highlighting the importance of each in achieving optimal performance in the beam optimization task. The vertical line marks the optimization level of a human operator.

provide complete information about the system. Utilizing recent advances in machine learning research, by adopting dense layers, allowed the agent to take more precise actions by forming more complex representations of the observations.

This empirical study demonstrates that each component of the architecture contributes to the overall effectiveness of the RL approach in optimizing the IBL at COSY. The integration of dense layers, observation noise, history, and domain randomization into the RL framework is essential for achieving high performance in complex optimization tasks.

VII. SUMMARY AND CONCLUSIONS

This research introduces a framework for the application of reinforcement learning (RL) to optimize the injection process at accelerator facilities. By employing an RL agent tailored to the specific challenges of optimizing the injection beam line, we developed an agent trained solely in simulation that is capable of optimizing the beam’s cross section dynamically to meet predefined targets during a live operation at the COSY facility.

Live run evaluations demonstrated that the RL agent could effectively optimize the beam’s cross section by aligning the beam’s central position and spread with the target values given to the agent by operators. The agent’s optimization accuracy was on par with that of a human operator, yet it completed the tasks in a time reduced by a factor of 3 of the time required by the human operator.

The empirical study of the architecture components validated the importance and significance of each element in the RL framework. Especially, the results highlighted the critical role of domain randomization and observation noise in preparing an agent trained solely in simulation for real-world control operation. In addition, incorporating historical data into the training process was shown to be essential

for dealing with partial observability and enhancing the agent’s decision-making capabilities.

In conclusion, this research illustrates the potential of machine learning methods, namely, RL, to enhance the efficiency of particle accelerator operations. The successful application of an RL agent for beam injection control at COSY demonstrates the potentials in adopting broader applications of machine learning techniques in particle accelerators and other fields requiring efficient and precise control of complex systems.

ACKNOWLEDGMENTS

The authors thank the COSY crew, particularly V. Kamerzhiev, for their great assistance in preparation and during the beam shifts, and for the support from JEDI team. In addition, the authors thank Hans-Joachim Stein for the insightful and helpful discussions. Simulations were performed with computing resources granted by RWTH Aachen University under Project No. rwth0905.

-
- [1] Jonas Degraeve, Federico Felici, Jonas Buchli, Michael Neunert, Brendan Tracey, Francesco Carpanese, Timo Ewalds, Roland Hafner, Abbas Abdolmaleki, Diego de las Casas, Craig Donner, Leslie Fritz, Cristian Galperti, Andrea Huber, James Keeling, Maria Tsimpoukelli, Jackie Kay, Antoine Merle, Jean-Marc Moret, Noury Seb *et al.*, Magnetic control of tokamak plasmas through deep reinforcement learning, *Nature (London)* **602**, 414 (2022).
 - [2] Arlene J. Lennox, Overview of accelerators in medicine, *Conf. Proc. C* **930517**, 1666 (1993).
 - [3] Ken Peach, Puthenparampil Wilson, and Bleddyn Jones, Accelerator science in medical physics, *Br. J. Radiol.* **84**, S4 (2011).

- [4] R. Maier, U. Bechstedt, J. Dietrich, S. Martin, D. Prasuhn, A. Schnase, H. Schneider, H. Stockhorst, and R. Tölle, Cooler synchrotron COSY, *Nucl. Phys.* **A626**, 395 (1997).
- [5] U. Bechstedt, J. Dietrich, R. Maier, S. Martin, D. Prasuhn, A. Schnase, H. Schneider, H. Stockhorst, and R. Tölle, The cooler synchrotron COSY in Jülich, *Nucl. Instrum. Methods Phys. Res., Sect. B* **113**, 26 (1996).
- [6] S. Karanth, E. J. Stephenson, S. P. Chang, V. Hejny, S. Park, J. Pretz, Y. K. Semertzidis, A. Wirzba, A. Wrońska *et al.*, First search for axionlike particles in a storage ring using a polarized deuteron beam, *Phys. Rev. X* **13**, 031004 (2023).
- [7] P. Adlarson, W. Augustyniak, W. Bardan, M. Bashkanov, F. S. Bergmann, M. Berłowski, H. Bhatt, M. Büscher, H. Calén, I. Ciepał, H. Clement, D. Coderre, E. Czerwiński, K. Demmich, E. Doroshkevich, R. Engels, A. Erven, W. Erven, W. Eyrich, P. Fedorets *et al.*, Evidence for a new resonance from polarized neutron-proton scattering, *Phys. Rev. Lett.* **112**, 202301 (2014).
- [8] Colin Wilkin, The legacy of the hadron physics programme at COSY, *EPJ Web Conf.* **130**, 01007 (2016).
- [9] Lucio Rossi, Technical challenges for future accelerators, in *Bruno Touschek 100 Years*, edited by Luisa Bonolis, Luciano Maiani, and Giulia Pancheri (Springer International Publishing, Cham, 2023), pp. 175–190, 10.1007/978-3-031-23042-4_14.
- [10] Vladimir Shiltsev, Stephen Gourlay, and Raubenheimer Tor, Challenges of future accelerators for particle physics research, *Front. Phys.* **10**, 6 (2022).
- [11] Claudio Emma, Auralee Edelen, Adi Hanuka, Brendan O’Shea, and Alexander Scheinker, Virtual diagnostic suite for electron beam prediction and control at FACET-II, *Information* **12**, 61 (2021).
- [12] A. Sanchez-Gonzalez, P. Micaelli, C. Olivier, T. R. Barillot, M. Ilchen, A. A. Lutman, A. Marinelli, T. Maxwell *et al.*, Accurate prediction of X-ray pulse properties from a free-electron laser using machine learning, *Nat. Commun.* **8**, 15461 (2017).
- [13] C. Emma, A. Edelen, M. J. Hogan, B. O’Shea, G. White, and V. Yakimenko, Machine learning-based longitudinal phase space prediction of particle accelerators, *Phys. Rev. Accel. Beams* **21**, 112802 (2018).
- [14] J. Duris, D. Kennedy, A. Hanuka, J. Shtalenkova, A. Edelen, P. Baxevanis, A. Egger, T. Cope, M. McIntire, S. Ermon, and D. Ratner, Bayesian optimization of a free-electron laser, *Phys. Rev. Lett.* **124**, 124801 (2020).
- [15] Alexander Scheinker, Auralee Edelen, Dorian Bohler, Claudio Emma, and Alberto Lutman, Demonstration of model-independent control of the longitudinal phase space of electron beams in the linac-coherent light source with femtosecond resolution, *Phys. Rev. Lett.* **121**, 044801 (2018).
- [16] S. C. Leemann, S. Liu, A. Hexemer, M. A. Marcus, C. N. Melton, H. Nishimura, and C. Sun, Demonstration of machine learning-based model-independent stabilization of source properties in synchrotron light sources, *Phys. Rev. Lett.* **123**, 194801 (2019).
- [17] A. Scheinker, Adaptive machine learning for time-varying systems: Low dimensional latent space tuning, *J. Instrum.* **16**, P10008 (2021).
- [18] Chris Tennant, Adam Carpenter, Tom Powers, Anna Shabalina Solopova, Lasitha Vidyaratne, and Khan Iftekharuddin, Superconducting radio-frequency cavity fault classification using machine learning at Jefferson Laboratory, *Phys. Rev. Accel. Beams* **23**, 114601 (2020).
- [19] E. Fol, R. Tomás, J. Coello de Portugal, and G. Franchetti, Detection of faulty beam position monitors using unsupervised learning, *Phys. Rev. Accel. Beams* **23**, 102805 (2020).
- [20] Xiaoying Pang, Sunil Thulasidasan, and Larry Rybarcyk, Autonomous control of a particle accelerator using deep reinforcement learning, arXiv:2010.08141.
- [21] D. Jia *et al.*, Design of an E x B chopper based on permanent magnets, in *Proceedings of the 14th International Particle Accelerator Conference (IPAC’23), Venice, Italy* (JACoW, Geneva, Switzerland, 2023), pp. 1446–1449.
- [22] A. Schuett *et al.*, Optimizing non-linear kicker injection parameters using machine learning, in *Proceedings of the 15th International Particle Accelerator Conference (IPAC’24), Nashville, TN* (JACoW, Geneva, Switzerland, 2024), pp. 3571–3574.
- [23] Verena Kain, Simon Hirlander, Brennan Goddard, Francesco Maria Velotti, Giovanni Zevi Della Porta, Niky Bruchon, and Gianluca Valentino, Sample-efficient reinforcement learning for CERN accelerator control, *Phys. Rev. Accel. Beams* **23**, 124801 (2020).
- [24] M. Cai *et al.*, Improving the performance of the SXFEL through Proximal Policy Optimization, in *Proceedings of the 14th International Particle Accelerator Conference (IPAC’23), Venice, Italy* (JACoW, Geneva, Switzerland, 2023), pp. 1916–1919.
- [25] Jan Kaiser, Chenran Xu, Annika Eichler, Andrea Santamaria Garcia, Oliver Stein, Erik Bründermann, Willi Kuroпка, Hannes Dinter, Frank Mayet, Thomas Vinatier, Florian Burkart, and Holger Schlarb, Reinforcement learning-trained optimisers and Bayesian optimisation for online particle accelerator tuning, *Sci. Rep.* **14**, 07 (2024).
- [26] Jonathan Ho and Stefano Ermon, *Generative Adversarial Imitation Learning*, Advances in Neural Information Processing Systems Vol. 29, edited by D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (Curran Associates, Inc., Red Hook, NY, 2016).
- [27] Chelsea Finn, Sergey Levine, and Pieter Abbeel, Guided cost learning: Deep inverse optimal control via policy optimization, in *Guided Cost Learning: Deep Inverse Optimal Control via Policy Optimization*, edited by Maria Florina Balcan and Kilian Q. Weinberger (PMLR, New York, NY, 2016), pp. 49–58.
- [28] Richard S. Sutton and Andrew G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. (The MIT Press, Cambridge, MA, 2018).
- [29] A. Awal, J. Hetzel, R. Gebel, V. Kamerdzhiev, and J. Pretz, Optimization of the injection beam line at the Cooler Synchrotron COSY using Bayesian optimization, *J. Instrum.* **18**, P04010 (2023).
- [30] Sabrina Appel, Vera Chetvertkova, Wolfgang Geithner, Frank Herfurth, Udo Krause, Stephan Reimann, Mariusz Sapinski, Petra Schütt, and Daniel Österle, Automated

- optimization of beam lines using evolutionary algorithms, in *Proceedings of the 8th International Particle Accelerator Conference, Copenhagen, Denmark, 2017* (JACoW, Geneva, Switzerland, 2017), pp. 3941–3944.
- [31] L. Catani, Beam transport and optimization tools based on evolutionary strategies, in *Workshop on Automated Beam Steering and Shaping (ABS)* (CERN, 1999), pp. 56–61, [10.5170/CERN-1999-007.56](https://cds.cern.ch/record/10.5170/CERN-1999-007.56).
- [32] Kenneth Li, Aspen K. Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg, Emergent world representations: Exploring a sequence model trained on a synthetic task, in *Proceedings of the 11th International Conference on Learning Representations* (Curran Associates, Inc., Red Hook, NY, 2023).
- [33] Awal Awal, Machine learning methods for injection optimization at the Cooler Synchrotron COSY, Ph.D. dissertation, RWTH Aachen University, 2024.
- [34] W. Brautigam, R. Brings, R. Gebel, R. Maier, A. Schnase, and H. Jungwirth, Operation of the cyclotron JULIC as injector for the cooler synchrotron COSY-Juelich, in *Proceedings of the 15th International Conference on Cyclotrons and Their Applications* (IOP, Bristol, 1999), pp. 654–657, <https://cds.cern.ch/record/318011>.
- [35] Beñat Alberdi Esuain, Optimization of injection in COSY, Master’s thesis, RWTH Aachen University, Aachen, Germany, 2019, https://collaborations.fz-juelich.de/ikp/jedi/public_files/theses/main.pdf.
- [36] Majid Ghasemi and Dariush Ebrahimi, Introduction to reinforcement learning, [arXiv:2408.07712](https://arxiv.org/abs/2408.07712).
- [37] Richard Bellman, *Dynamic Programming* (Dover Publications, Mineola, NY, 1957).
- [38] Richard Bellman, A Markovian decision process, *J. Math. Mech.* **6**, 679 (1957).
- [39] Leslie Pack Kaelbling, Michael L. Littman, and Anthony R. Cassandra, Planning and acting in partially observable stochastic domains, *Artif. Intell.* **101**, 99 (1998).
- [40] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller, Playing Atari with deep reinforcement learning, [arXiv:1312.5602](https://arxiv.org/abs/1312.5602).
- [41] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin A. Riedmiller, Andreas Kirkeby Fidjeland, Georg Ostrovski, Stig Petersen, Charlie Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis, Human-level control through deep reinforcement learning, *Nature (London)* **518**, 529 (2015).
- [42] Nicolas Heess, Jonathan J. Hunt, Timothy P. Lillicrap, and David Silver, Memory-based control with recurrent neural networks, [arXiv:1512.04455](https://arxiv.org/abs/1512.04455).
- [43] Doo Re Song, Chuanyu Yang, Christopher McGreavy, and Zhibin Li, Recurrent network-based deterministic policy gradient for solving bipedal walking challenge on rugged terrains, [arXiv:1710.02896](https://arxiv.org/abs/1710.02896).
- [44] Ramya Ramakrishnan, Ece Kamar, Debadeepta Dey, Eric Horvitz, and Julie Shah, Blind spot detection for safe sim-to-real transfer, *J. Artif. Intell. Res.* **67**, 191 (2020).
- [45] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel, Domain randomization for transferring deep neural networks from simulation to the real world, in *Proceedings of the 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (2017), pp. 23–30.
- [46] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and P. Abbeel, Sim-to-real transfer of robotic control with dynamics randomization, in *Proceedings of the 2018 IEEE International Conference on Robotics and Automation (ICRA)* (IEEE, New York, 2018), pp. 3803–3810, [10.1109/ICRA.2018.8460528](https://arxiv.org/abs/10.1109/ICRA.2018.8460528).
- [47] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver, Universal value function approximators, in *Proceedings of the 32nd International Conference on Machine Learning, Lille, France*, edited by Francis Bach and David Blei (PMLR, New York, NY, 2015), pp. 1312–1320.
- [48] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine, Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor, in *Proceedings of the 35th International Conference on Machine Learning*, edited by Jennifer Dy and Andreas Krause (PMLR, New York, NY, 2018), pp. 1861–1870.
- [49] Thomas Degris, Martha White, and Richard Sutton, Off-policy actor-critic, in *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, edited by John Langford and Joelle Pineau (Omnipress, New York, NY, 2012), pp. 457–464.
- [50] Ziyu Wang, Victor Bapst, Nicolas Heess, Volodymyr Mnih, Remi Munos, Koray Kavukcuoglu, and Nando de Freitas, Sample efficient actor-critic with experience replay, [arXiv:1611.01224v2](https://arxiv.org/abs/1611.01224v2).
- [51] Samarth Sinha, Homanga Bharadhwaj, Aravind Srinivas, and Animesh Garg, D2RL: Deep dense architectures in reinforcement learning, [arXiv:2010.09163](https://arxiv.org/abs/2010.09163).
- [52] Gao Huang, Zhuang Liu, and Kilian Q. Weinberger, Densely connected convolutional networks, [arXiv:1608.06993](https://arxiv.org/abs/1608.06993).
- [53] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, Tom Goldstein, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Visualizing the loss landscape of neural nets, in *Advances in Neural Information Processing Systems* (Curran Associates, Inc., Red Hook, NY, 2018), Vol. 31.
- [54] Parand Alizadeh Alamdari, Toryn Q. Klassen, Rodrigo Toro Icarte, and Sheila A. McIlraith, Be considerate: Avoiding negative side effects in reinforcement learning, in *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems, AAMAS ’22, Richland, SC* (International Foundation for Autonomous Agents and Multiagent Systems, 2022), pp. 18–26.
- [55] Aditya Gudimella, Ross Story, Matineh Shaker, Ruofan Kong, Matthew Brown, Victor Shnayder, and Marcos Campos, Deep reinforcement learning for dexterous manipulation with concept networks, [arXiv:1709.06977](https://arxiv.org/abs/1709.06977).
- [56] H. Grote and F. Schmidt, MAD-X: An upgrade from MAD8, *Conf. Proc. C* **030512**, 3497 (2003).

- [57] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba, OpenAI Gym, [arXiv:1606.01540](https://arxiv.org/abs/1606.01540).
- [58] Mark Towers, Jordan K. Terry, Ariel Kwiatkowski, John U. Balis, Gianluca de Cola, Tristan Deleu, Manuel Goulão, Andreas Kallinteris, K. G. Arjun, Markus Krimmel, Rodrigo Perez-Vicente, Andrea Pierré, Sander Schulhoff, Jun Jet Tai, Andrew Tan Jin Shen, and Omar G. Younis, Gymnasium: A standard interface for reinforcement learning environments, [arXiv:2407.17032](https://arxiv.org/abs/2407.17032).
- [59] F. Schmidt, Mad-X PTC integration, in *Proceedings of the 2005 Particle Accelerator Conference* (IEEE, New York, 2005), pp. 1272–1274, [10.1109/PAC.2005.1590731](https://doi.org/10.1109/PAC.2005.1590731).
- [60] P. K. Skowronski, F. Schmidt, and E. Forest, Advances in MAD-X using PTC, *Conf. Proc. C* **070625**, 3381 (2007).
- [61] Diederik P. Kingma and Jimmy Ba, Adam: A method for stochastic optimization, [arXiv:1412.6980](https://arxiv.org/abs/1412.6980).
- [62] Leo R. Dalesio, Jeffrey O. Hill, Martin Kraimer, Stephen Lewis, Douglas Murray, Stephan Hunt, William Watson, Matthias Clausen, and John Dalesio, The experimental physics and industrial control system architecture: Past, present, and future, *Nucl. Instrum. Methods Phys. Res., Sect. A* **352**, 179 (1994).
- [63] Matthias Clausen and Leo Dalesio, EPICS: Experimental physics and industrial control system, ICFA Beam Dyn. Newslett. **47**, 56 (2008), <https://inspirehep.net/literature/808273>.
- [64] Heidelberg Ion-Beam Therapy Center (HIT). cpymad: A Cython binding to MAD-X (2014), <https://github.com/hibtcpymad> [accessed March 7, 2025].
- [65] areaDetector - EPICS area detector software, <https://areadetector.github.io/areaDetector/index.html> [accessed March 7, 2025].