# Microsecond-Latency Feedback at a Particle Accelerator by Online Reinforcement Learning on Hardware

Luca Scomparin,* Michele Caselle, Andrea Santamaria Garcia, Chenran Xu,
Edmund Blomley, Timo Dritschler, Akira Mochihashi, Marcel Schuh, Johannes L. Steinmann,
Erik Bründermann, Andreas Kopmann, Jürgen Becker, Anke-Susanne Müller, and Marc Weber

*Karlsruher Institut für Technologie, Kaiserstraße 12, 76131 Karlsruhe, Deutschland*

(Dated: September 25, 2024)

The commissioning and operation of future large-scale scientific experiments will challenge current tuning and control methods. Reinforcement learning (RL) algorithms are a promising solution thanks to their capability of autonomously tackling a control problem based on a task parameterized by a reward function. The conventionally utilized machine learning (ML) libraries are not intended for microsecond latency applications, as they mostly optimize for throughput performance. On the other hand, most of the programmable logic implementations are meant for computation acceleration, not being intended to work in a real-time environment. To overcome these limitations of current implementations, RL needs to be deployed on-the-edge, i.e. on to the device gathering the training data. In this paper we present the design and deployment of an experience accumulator system in a particle accelerator. In this system deep-RL algorithms run using hardware acceleration and act within a few microseconds, enabling the use of RL for control of ultra-fast phenomena. The training is performed offline to reduce the number of operations carried out on the acceleration hardware. The proposed architecture was tested in real experimental conditions at the Karlsruhe research accelerator (KARA), serving also as a synchrotron light source, where the system was used to control induced horizontal betatron oscillations in real-time. The results showed a performance comparable to the commercial feedback system available at the accelerator, proving the viability and potential of this approach. Due to the self-learning and reconfiguration capability of this implementation, its seamless application to other control problems is possible. Applications range from particle accelerators to large-scale research and industrial facilities.

## I. INTRODUCTION

The tuning of future large-scale scientific experiments will pose great challenges. Due to the sheer amount of parameters that require adjustment, manual tuning will be an extremely demanding task. This will in turn impact the budget of large scientific endeavors as the resource requirements for commissioning and operation of the facility will potentially become unbearable. Automatic algorithms capable of finding an optimum in a parameter space, also known as optimizers, are a possible solution and have already been successfully applied to particle accelerators [1–3], with some attempts of creating standardized interfaces and algorithm libraries [4]. Even for these more refined techniques, the increase in the number of variables involved in the task will reduce their performance, a phenomenon known as the *curse of dimensionality*.

One of the possible solutions is the use of data-driven machine learning (ML) techniques [5–7]. In the case of control problems, a promising approach is the use of reinforcement learning (RL), a paradigm where an agent is trained to learn a policy in order to maximize a reward function, encoding a control goal, acting on an environment based on a set of observations. RL algorithms have already been proven to be a powerful tool through their application to control problems in several large scale facilities [8–14].

Typically, a controller must operate within time scales comparable to those of the controlled dynamics. The fact that the evaluation of a control policy needs to be performed in a sched-

uled time frame is by definition a real-time constraint [15]. Failure to meet this kind of requirement means the controller could be partially effective or even completely ineffective. As such, when the inference time reaches the microsecond timescale, the use of conventional CPUs becomes challenging, as the several layers of abstraction and caching between the application and hardware limit the achievable latency performance. Specifically, the latency achieved can not only be higher than the requirements, but its high variance could produce spurious, unpredictable violations of these constraints. Different kinds of computing hardware can circumvent these limitations by combining a higher degree of parallelization with a system of lower intrinsic overhead. Modern heterogeneous computing platforms are a good candidate for this, combining conventional CPUs with a field programmable gate array (FPGA) and an artificial intelligence (AI) accelerator. An FPGA allows the definition of custom programmable digital logic that can reach levels of parallelism that are constrained only by the number of programmable cells available on the device. Moreover, they allow precise control over the timing of each operation. The AI accelerator allows the efficient execution of more specific tasks (e.g. floating point vector multiplication), for which FPGAs are less suitable. These three components share memory, allowing a customized and optimized data-flow architecture.

One of the main issues that prevent the application of RL to many large-scale facilities is their sample efficiency, the amount of interaction with an environment required to fully train an agent. This means that the data necessary to train a successful agent is usually difficult to obtain in the real world due to their low repetition rates [16]. This issue can be mitigated by pre-training on a simulation, with the risk that the agent might perform sub-optimally in case the simulation differs signifi-

---

* luca.scomparin@kit.edu

cantly from the real world [9]. A different approach is to use more sample efficient algorithms, but those tend to complicate the tuning of the hyper-parameters [12]. The benefit of bringing RL to facilities operating on microsecond timescales, like numerous synchrotron light sources, MHz-rate free electron lasers and circular colliders, is that sufficient training data can be generated in real time, bypassing the need of training in simulation. Conversely, the constraints on the time taken to select an action require a more careful implementation of the agent.

In this work, we developed a closed-loop RL feedback system that was successfully deployed and commissioned at the Karlsruhe research accelerator (KARA), an accelerator test facility and synchrotron light source. This novel scheme enables direct RL online training in the microsecond time domain, setting the precedent as the first implementation of its kind in particle accelerators.

## II.   REAL-TIME REINFORCEMENT LEARNING

In RL the environment is modeled as a Markov decision process (MDP), a 4-tuple $(S,A,P,R)$ where $S$ is the state space, $A$ is the action space, $P$ are the transition dynamics, and $R$ is the set of all rewards. The environment provides a scalar reward value $R$ that encodes the goal of the control problem. The goal of an RL agent is to maximize the expected return $G$, defined as the cumulative reward in all future steps

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}, \qquad (1)$$

where $\gamma \in [0,1]$ is the discount factor, used to bound the cumulative reward in infinite horizon problems, extending the MDP to $(S,A,P,R,\gamma)$.

This framework models decision-making in partially stochastic control processes and incorporates the Markov property, where the conditional probability distribution $P$ of future states depends only on the present state.

Often the full state of the environment is not directly accessible and needs to be constructed from observations, giving place to partially observable Markov decision process (POMDP) defined as $(S,A,P,R,\gamma,O,\varepsilon)$, where $O$ is the space of possible observations and $\varepsilon$ the probability of transitioning to a new state $s_t$ when $o_t$ is observed $\varepsilon = \mathbb{P}(s_t|o_t)$.

Many modern RL algorithms are based on an actor-critic architecture, where an actor function, responsible for choosing the action, is paired with a critic function, estimating the expected return of a given state. The critic can either implement a *value function* $v_\pi : S \rightarrow \mathbb{R}$ or *action-value function* $Q_\pi : S \times A \rightarrow \mathbb{R}$, depending on whether it also takes into account the action that is going to be taken or not [17]. In modern deep-RL both actor and critic are approximated with a neural network (NN).

It is worth noticing that, given the interaction with an active environment, most non-simulated RL agents are subject to real-time constraints. Depending on the time scale of the

environment, this requirement could be addressed by simply allocating enough computing power to the agent. When the latency requirements become more demanding, the compound overhead of transferring data and evaluating the agent can mandate more adaptable computing architectures like FPGAs or heterogeneous platforms, in order to better control the scheduling of the data-processing pipeline.

### A.   Related work: RL on FPGA

A brief outline of the currently available implementations of RL on FPGA is described below. A more comprehensive review can be found in reference [18]. There are two main classes of implementations, depending on the nature of the value function and policy. Several works store the action-value function in a tabular form, and based on this choose the action with the maximum expected cumulative reward for each state [19–23]. The main issue of this kind of system is that its resource usage grows exponentially with the size of the observation and action space, while being not directly applicable to environments with continuous action and observation spaces. These kind of environments are widely used in particle accelerator applications and better fit the problem of this work.

A different approach is deep-RL, where the value function and policies are approximated with a NN [24–30]. This approach scales better with increasing size of the observation and action vector, but tends to have a more complex training routine. These implementations lack the flexibility of choosing different training algorithms and changing the NN topology during operation. Distinct algorithms and their respective hyperparameter selections might exhibit different performances based on the problem at hand. Therefore, the capability of flexibly choosing these components greatly reduces the deployment effort and expedites the development process. Moreover, the policy NN implementations shown in the references [19–30] do not have real-time applications in mind, and as such the application to real-time particle accelerator controls would be limited.

For these reasons, we have chosen to implement a deep-RL system capable of performing such dynamic reconfiguration and naturally fulfilling real-time constraints.

### B.   Experience accumulator architecture

In this study, the system is implemented by employing an *experience accumulator* architecture. Depending on the complexity of the algorithm, the *training* processes in modern deep-RL can take a large amount of time compared to the time required for inference of the policy. As such, it is not possible to perform them at a microsecond level timescale. A possible solution is to only implement a general real-time policy $\pi^{(\text{edge})}$ for *inference*, using cutting-edge low-latency computing platforms. During application time, the real-time policy network only needs to perform *forward passes*, i.e. predicting the next action $a_i$ based on the observed signal $o_i$. The interactions of the policy NN with the environment are recorded, providing

state-action-reward tuples $\{(s_i, a_i, r_i)_{i=1,...,T}\}$, that can be used for asynchronously training an emulated copy of the agent $\pi^{(\text{CPU})} \leftarrow \pi^{(\text{edge})}$ on conventional computing platforms, such as a CPU. As such it is possible to achieve a greater degree of flexibility and abstraction (Figure 1). This comes at the expense of higher overhead during training, while still maintaining low inference latency. For the proximal policy optimization (PPO) algorithm agent used in this work, only the actor NN is implemented on hardware, while the critic can be evaluated at training time.

A similar approach is presented in a previous work [31], despite applying it to a system with less stringent latency requirement, thus not requiring heterogeneous computing platforms but an off-the-shelf x86 processor. Additionally, in this study, we further extend the capability of the experience accumulator architecture by introducing a post-hoc reward definition. In the definition of RL control problems, the reward is provided on a step-by-step basis by the environment. In the case of a high-speed feedback system, interacting with a real-time environment, it would then become necessary to perform this computation at pace with the environment. The need for a real-time reward function implementation increases the complexity of the system, as it needs to be implemented in the FPGA or AI accelerator, requiring specialized expertise.

In the experience accumulator architecture, a further simplification is possible. The rewards obtained by an agent are only necessary during its training. If the reward is a function only of the observation and action vectors, it is possible to skip its real-time computation and perform it at training time based on the stored observation-action tuples. This *Training Time Reward Definition* technique reduces development efforts as it allows the function to be implemented on a CPU, where the expertise required is more commonly available. In the implementation used in this work, the reward function is developed in Python and specified by the experimenter before starting the training, in this way simplifying the *reward engineering*.
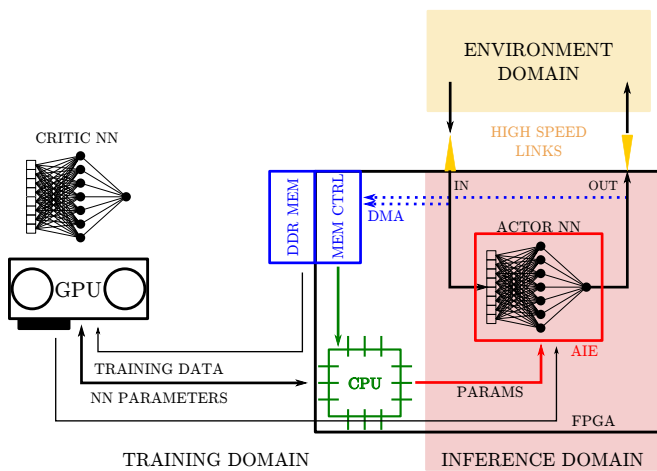


FIG. 1. Schematic drawing of the experience accumulator architecture showing the partition between real-time or inference domain (on the right) and training-time domain (on the left) on a heterogeneous platform.

## III. CLOSING THE RL LOOP: THE KINGFISHER SYSTEM

In order to simplify the deployment of RL policies in an experience accumulator architecture, all the FPGA infrastructure necessary for interfacing with the particle accelerator and storing the agent experience were implemented into an AMD-Xilinx Extensible Platform called KINGFISHER [32]. This approach has the effect of decoupling the low-level operations allowing the development of the policy and feature extraction algorithm in high-level programming languages, thanks to the aid of tools like high level synthesis (HLS).

The device of choice is the AMD-Xilinx Versal VCK190 [33], a novel heterogeneous computing platform comprising a FPGA and ARM processor in a system-on-chip architecture, with an AI engine (AIE) array capable of accelerating multiplication intensive tasks, as, for instance, NN inference.

The KINGFISHER platform receives data via a 40 Gbps Aurora 64b/66b link [34]. This data can be used to create the observation data-stream. The action data-stream is used to control a digital to analog converter (DAC) to produce an analog control signal. For the purpose of future experiments special digital output interfaces are also available. The observation-action data-streams are written to the DDR memory, without active intervention of the CPU, by means of a direct memory access (DMA) block implemented in the FPGA. The ARM processor in the Versal core then copies the data in the DDR to a file that is made available over the network. Moreover, a control server makes the parameters of the system available via the EPICS [35] control system. A trigger input is available to start the agent action at a precise time. Moreover, a counter allows the interaction with the environment only for a limited number of steps.

In order to guide the exploration of RL algorithms, a stochastic component is usually needed. Thus an intellectual property (IP) core implementing the permuted congruential generator (PCG) algorithm [36] produces a continuous stream of 32-bit floating pseudorandom-numbers uniformly distributed between 0 and 1 at a rate of 125 MSps.

For the control of the horizontal betatron oscillation (HBO) the specific deployment configuration is shown in Figure 2. To sample the fast-peaked signals with high MHz rates commonly generated by diagnostic instrumentation at accelerator-based light sources like synchrotrons, the KIT-developed KAPTURE system [37] is utilized. This system functions as a digitizer, enabling the sampling of four channels at a frequency equal to the accelerator's radio frequency (RF), with a selectable sampling interval down to 3 ps. The system utilizes an Highflex card with the capability of performing transfers directly to a GPU with the GPUDirect technology [38]. This setup provides bunch-by-bunch and turn-by-turn data. In order to obtain the bunch position relevant for the control of the HBO, this infrastructure is used to sample the horizontal position signal from an analog Dimtel BPMH-20-2G hybrid [39] combining the four signals from a button beam position monitor (BPM). The action analog signal is fed into a broadband amplifier and then applied to a stripline kicker used to influence the beam.
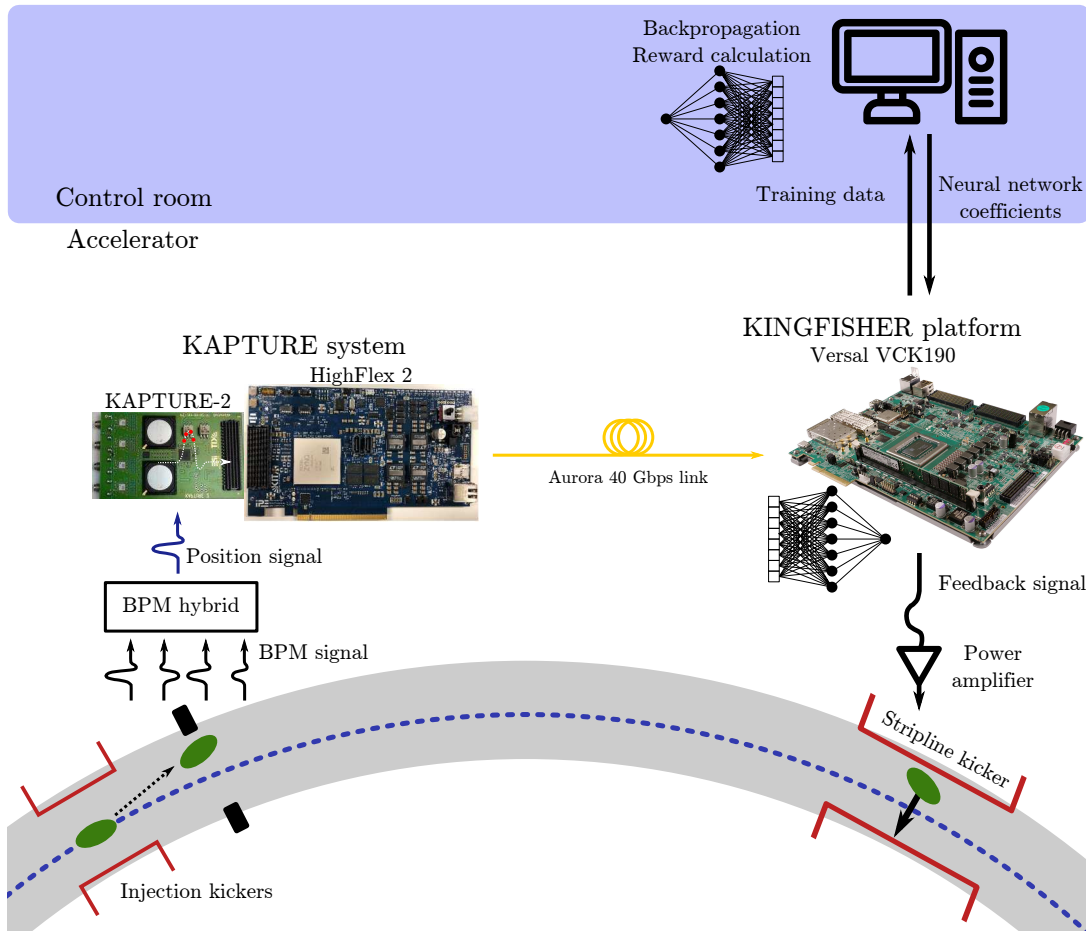
FIG. 2. Drawing of the hardware infrastructure employed in this work. The bunch position in the beam pipe is measured by processing a BPM signal with an analog hybrid. The produced analog signal is sampled with KAPTURE and then forwarded by an HighFlex 2 board through high-speed links to a Versal device. The KINGFISHER system programmed on this Versal then connects this data stream into the RL controller, while applying the output action to a stripline kicker. Every episode, comprised of 2048 interaction steps with accelerator, data is sent back to the control room for training.

### A. Implementing Reinforcement Learning Algorithms on KINGFISHER

The algorithm used in the current work is PPO [40]. This choice was dictated by its stability with respect to the change of hyperparameters. The reduced sample efficiency compared to off-policy algorithms like soft actor-critic (SAC) does not affect the current application, as it is counterbalanced by the high experience collection rate. Other RL algorithms are nonetheless easy to integrate thanks to the experience accumulator architecture.

The actor network is implemented in the AIE. The employed PPO algorithm implementation uses a NN to select the mean value of a Gaussian distribution, from which the action applied to the environment is chosen. The standard deviation of the probability distribution is a trainable parameter, that is updated together with the NN coefficients.

A schematic of the internal data processing within the actor and the KINGFISHER platform is shown in Figure 3. The first AIE tile implements a circular buffer and streams the latest eight samples to the following kernel using the cascade stream interface. These eight samples represent the observation vector, the choice of which will be described more thoroughly in Section IV B. The next kernel implements the linear layer of a NN computing the values of the sixteen hidden neurons. A rectified linear unit (ReLU) activation function is applied to the outputs. Such an activation function was selected because of its simple implementation. This function can also be turned off while the system is still running in order to implement linear filters. The output of the network is then computed with a final linear layer and passed to the last output kernel. All these kernels keep forwarding the eight input values. The last kernel takes 16 values from uniformly-distributed random number data stream and sums them, producing an approximately Gaussian-distributed sample due to the central limit theorem. This is used to add a Gaussian noise with a standard deviation selectable at run-time to the output of the network. A final data vector containing the eight input values, the random value, and the output of the network previous to the noise addition is given to the experience accumulator logic.
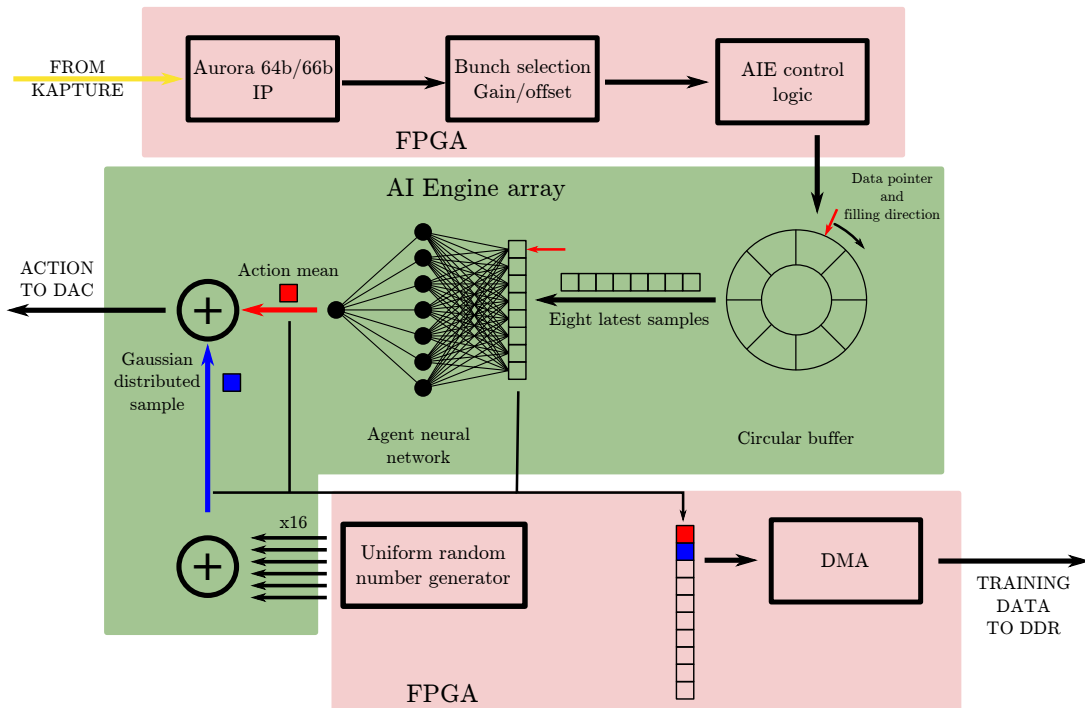
FIG. 3. Schematic representation of the data path within the Versal VCK190 firmware. In the top part, the KAPTURE data stream coming from HighFlex 2 is decoded by an Aurora IP from the protocol used for fiber-optic communication. The bunch of interest for the control experiment is selected and gain/offset correction are applied. FPGA logic takes care of forwarding data and gracefully stopping the AIE. In the middle, a buffer stores the latest eight samples and feeds them to the NN. The exploration noise is added and the action is then output to a DAC, for control of the kicker. In the bottom side, a FPGA block takes the latest data from the inference and stores them in memory through a DMA for later training.

| Hyper-parameter | Value |
|---|---|
| Learning rate $\eta$ | 0.0012 |
| Discount rate $\gamma$ | 0.99 |
| Number of steps | 2048 |
| Batch size | 64 |
| Number of epochs | 1 |

TABLE I. Hyper-parameters used for the PPO algorithm.

A computer in the control room then fetches the data and uses a modified version of the Stable-baselines3 library [41] PPO [40] implementation to train the policy using the hyper-parameters in Table I. The new parameters are then loaded onto the AIE kernel at run-time and new data for training is gathered. The whole inference loop has a latency of 2.8 μs. For this work, a number of 2048 action steps was chosen. This number has been manually selected in such a way that the agent would not have enough time to cause beam loss and disrupt operation.

The critic network was chosen to be identical to the actor network. The remaining hyper-parameters are available in Table I.

## IV. DEMONSTRATION OF REAL-TIME RL CONTROL AT THE ACCELERATOR

### A. Task description: Horizontal betatron oscillations at KARA

The Karlsruhe research accelerator (KARA) is a ramping electron storage ring with 0.5 GeV injection energy, with various operation modes between 0.5 GeV and 2.5 GeV. Its lattice is based on four sectors of double-bend achromats, with a total circumference of 110 m.

In synchrotrons, the horizontal beam dynamics is dominated at first order by the linear betatron motion. The bunch displacement from the reference orbit $x$ can be expressed with Hill's equation

$$\frac{d^2x}{ds^2} + K_x(s)x = 0 \tag{2}$$

where $s$ is the length along the ring and $K_x$ are the periodic focusing functions.

The general solution to this equation is

$$x(s) = x_0 \sqrt{\beta_x(s)} \cos[\theta(s) + \theta_0]$$
$$\theta(s) = \int_0^s \frac{ds}{\beta_x(s)} \tag{3}$$

where $\beta_x$ are the beta-functions and $(x_0, \theta_0)$ set the initial conditions. The number of full oscillations per revolution is defined

as the horizontal tune $Q_x$. In injection mode, used for this experiment, $Q_x \approx 6.76$. An observer fixed in a given position of the accelerator will only observe the fractional part, corresponding to an oscillation frequency in the order of $700\,\text{kHz}$, depending both on the operation mode, the beam current and more generally the state of the machine. This oscillation frequency sets the timescale at which the RL agent must be able to act, in this case in the order of a few microseconds.

A stripline kicker, based on the design from [42], is capable of affecting the HBOs by applying a bunch-by-bunch and turn-by-turn horizontal force to the beam based on the signal provided into an analog input.

The injection in the KARA storage ring makes use of three strong kicker and one septum magnet in order to properly merge the beam already in the machine with the one that is being injected. These kicker magnets can only activate at a rate of $1\,\text{Hz}$, for a few revolutions, and can move the beam on a displaced orbit in the horizontal plane. When the kicker switches back off, the displaced bunches will start performing betatron oscillations around the reference orbit. The goal of the RL agent is to damp this oscillation as quickly as possible. Notably, the strength of the injection kickers is orders of magnitude stronger than the stripline kickers used for feedback, thus an agent cannot damp an oscillation in a single kick, and turn-by-turn control is required.

A classical controller for this problem exists and is already available in commercial solutions. These bunch-by-bunch (BBB) feedback systems [43] are usually based on a finite impulse response (FIR) filter that takes the input signal and applies a $\pi\,\text{rad}$ phase at the frequency of the instability. In this way, a linear kick is produced with an opposite sign compared to the displacement, in this way damping the oscillations. The output of this kind of filter can be computed as

$$y(t) = \sum_{i=0}^{N} c_i x(t-i) \quad (4)$$

where $c_i$ are the coefficients of the filter, and $N$ is its order. The tuning of the filter coefficients directly impacts the performance of the controller, as it defines its behavior with respect to external noise and the bandwidth over which a suitable phase offset is produced. So far this is usually hand-tuned. A review on the topic is provided in [44].

### B. Formulation as an RL task

For the current problem of controlling the HBO, the RL environment was modeled as follows. Given the HBO dynamics at a fixed position in the storage ring can be approximated by an harmonic oscillator, the position $x$ and its derivative $\dot{x}$ are sufficient to have full knowledge of the state of the system. In a discrete-time setting, the derivative can be computed from the time difference of two consecutive samples.

$$\dot{x}(t) = \frac{x(t) - x(t-1)}{\Delta t} \quad (5)$$

This in turn means that the two latest position samples, $x(t)$ and $x(t-1)$, are also a full representation of the system's state.

In practice, though, only having two values is subject to measurement noise. Thus, the observation vector was defined as the last eight positions $\boldsymbol{o}_t = (x(t), x(t-1), ..., x(t-7))^T$. The signal is sampled at the revolution frequency, i. e. ca. $2.7\,\text{MHz}$, thus eight samples span roughly two periods of the betatron oscillation.

The action is a force that is applied to the bunch through a stripline kicker. In the harmonic oscillator model, this corresponds to a driving force. Under this definition, the system is a MDP.

Several different reward definitions were chosen, and the respective performance of the final agents are compared in Section IV E. All rewards studied in this paper penalize the agent when the $x$ position differs from zero, corresponding to the reference orbit.

### C. Simulation study

In order to study the interaction of an agent with the HBO, an environment based on the Gymnasium library [45] was developed. The dynamics was modeled as a damped harmonic oscillator with user selectable undamped angular frequency $\omega_0$ and damping ratio $\Gamma$. The environment stores the actions $a_i$ performed on the system and convolves this vector with the Green's function $B(t,t')$ of the damped harmonic oscillator

$$B(t,t') = \Theta(t-t') \frac{e^{-\Gamma(t-t')}}{\omega} \sin \omega(t-t'), \quad (6)$$

with $\omega = \sqrt{\omega_0^2 - \Gamma^2}$. As such, the position $x(t)$ is computed as

$$x(t) = \sum_i B(t, (i+\Delta\tau)T_{\text{rev}}) a_i, \quad (7)$$

where $\Theta(t)$ is the Heaviside step function and $T_{\text{rev}}$ is the revolution time of the accelerator. An additional user selectable delay $\Delta\tau$ is added to the argument of the function to study the effect of latency. Gaussian noise is added to the samples, reflecting the behavior of real-life data. A kick of intensity one order of magnitude higher than what the agent can perform is applied at a random time to simulate the external kicker.

In order to guide the selection of an RL algorithm and agent structure, the environment was used to test the training performance with the algorithms available in the Stable-baselines3 library [41]. PPO and the observation vector definition of Section IV B were thus validated in simulation before testing the complete system on the accelerator. This is necessary in order to disentangle a hardware platform failure from an issue with the RL problem formulation, in the case control could not achieved during the tests at KARA.

### D. Reinforcement Learning Control

The system discussed in this study was allowed to interact with the accelerator for 2048 revolutions (corresponding to

784 µs). During this period an external kicker excited the oscillations as shown in Figure 4. After each of these episodes, a training step was performed, updating the coefficients of the NN. The new set of weights and biases were uploaded to the agent and the operation was repeated.

In order to study the evolution of the oscillation amplitude, the amplitude of the oscillation was obtained as the absolute value of the Hilbert transform of the raw oscillation signal $x(t)$. As shown in Figure 4 an increase in the damping rate of the oscillations is an indication that the agent in question is achieving control of the environment. Moreover, it is possible to examine the trend of the cumulative reward obtained during each training episode as shown in Figure 5. A clear increase in the obtained cumulative reward is visible as more episodes are used for training. This clearly shows that the agent improves with experience, as it is expected.

Several different training configurations were tested, each one with a different reward definition and the number of neurons in the hidden layer of the actor. This led to agents with different performances. An example of training configuration is L2, 12 N, meaning the L2 norm defined in Table II is used, together with an actor having 12 neurons in the hidden layer.

### E. Training Time Reward Definition

Figure 6 compares the performance of different reward functions employed during training. To do so, the oscillation amplitude was fitted with an exponential function

$$f(t; A, \lambda) = A e^{-t\lambda} \tag{8}$$

and the damping rate $\lambda$ was employed as a reward independent metric. Provided $x$ is the position obtained from the BPM, the reward functions definitions are listed in Table II.

| Reward name | Definition |
|---|---|
| L1 | $-|x|$ |
| L2 | $-x^2$ |
| Tanhsq | $-\tanh\left(x^2\right)$ |

TABLE II. Definition of the different reward functions used experimentally, where $x$ denotes the transverse horizontal position of the beam read by the BPM.

The agents trained with all of these three reward choices reached a final performance better than the FIR controller and the baseline with the untrained agents.

### F. Online Network Structure Reconfiguration

In order to further increase the level of flexibility of the system, the possibility of dynamically modifying the NN structure without the need of re-implementing or re-packaging the firmware was implemented. This was achieved by embedding a smaller network into one with a greater number of neurons and layers by appropriately switching off different weights.

Additionally, maintaining the number of computations constant allows to have identical latency between different training trials, thus removing this variability when comparing different agents.

Agents with several different layer sizes have been trained and their performance is shown in Figure 6. The performance of all agents increased with training, outperforming the traditional FIR controller. The best performing agent, with 12 neurons in the hidden layer and trained with an L2-norm reward (in short notation: L2, 12 N; cf. also Table II), is used throughout the rest of this work for comparison with classical control techniques.

### G. Training stability and robustness

The training procedure was repeated several times, with the same setting but a different beam current, to study the stability of the agents produced. One would expect an effect for two reasons. First of all the BPM signal is not normalized, so the amplitude will vary with current. Second of all the HBO tune is current dependent [46]. Despite a 20% reduction in beam current due to the natural decay of the beam, all resulting agents achieve a very similar final reward, as shown in Figure 7. This shows the robustness of the RL agent against variations of current.

One of the main components of the KARA storage ring injection line is the injection septum magnet. Its impulse activation is necessary to guarantee the injection of the electron bunch coming from the booster into the main ring. The leaking magnetic field, though, also affects the beam that is already in the storage ring. This effect is visible in Figure 8, which corresponds to a shift in the position of the beam. Such an effect was not present in the simulation, but it was nonetheless possible to train an agent capable of correctly handling this new phenomenon. This is an example of the versatility and adaptability of RL algorithms, that are sometimes able to autonomously learn from situations they are not originally designed for.

### H. Improvement during cumulative reward plateau

As can be seen in Figure 7, the cumulative reward reaches a plateau around step number 50. Nonetheless, if one studies the trend of the damping rate measured at a BPM in a different part of the ring, it is still possible to observe an increase of the damping rate even around step 100 as shown in Figure 9. This is due to the fact that noise in the input data adds an offset to the cumulative reward that hides small improvements. Such a phenomenon needs to be considered in future experiments as it could potentially hinder further improvement of the agent.

### V. DISCUSSION

We trained several working agents and their performance can be compared and evaluated.
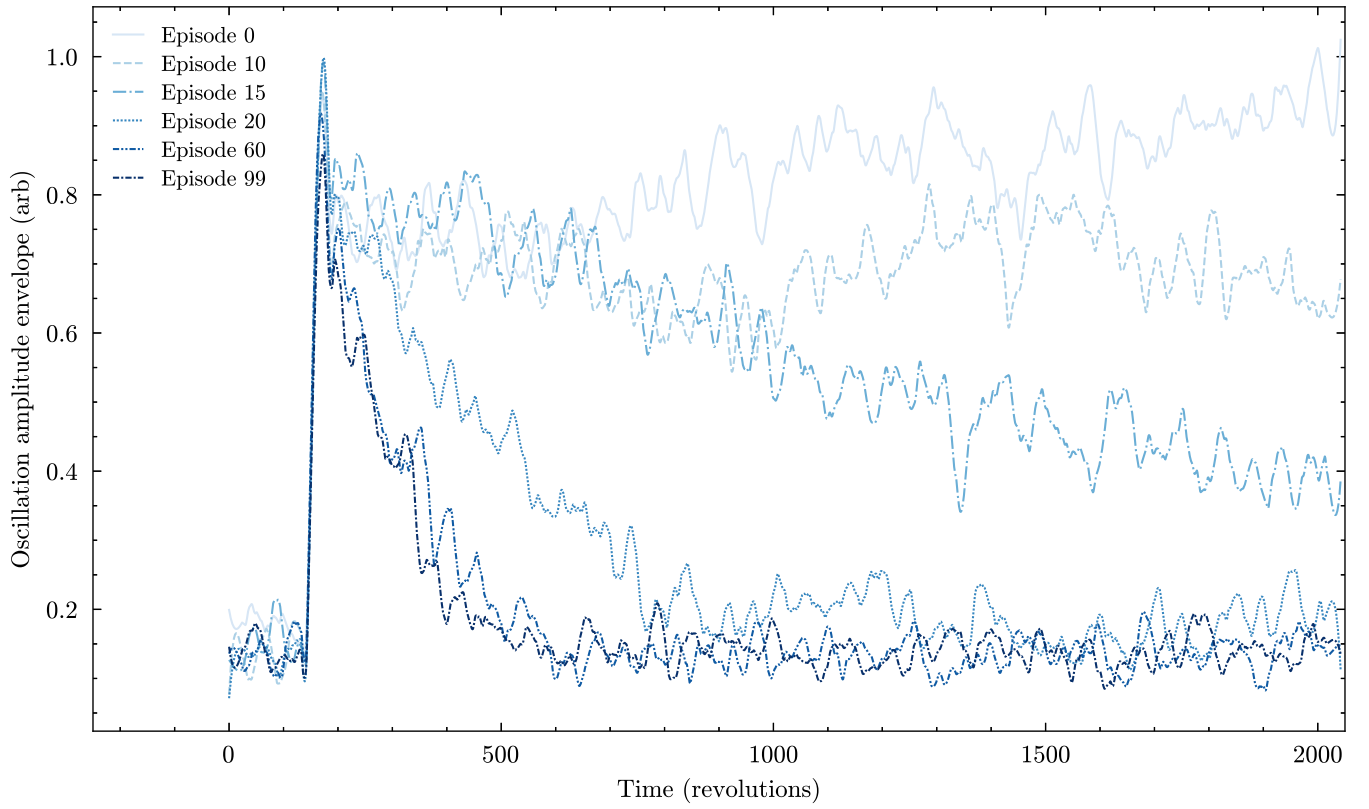
FIG. 4. Smoothed envelope of oscillations measured by a BPM. The sharp increase around $t = 150$ revolutions is due to the injection kicker emulating an instantaneous external excitation. Notice how at step 0 the randomly selected agent is destabilizing the beam, leading to an increase of the oscillation amplitude. Moreover, the rate of damping increases with the number of training steps, i.e. with the agent experience.
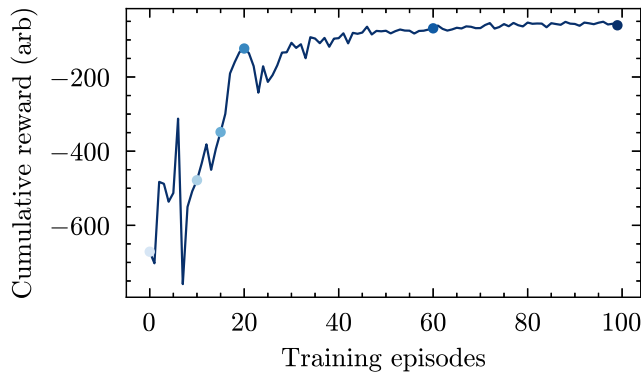


FIG. 5. Cumulative reward obtained by an agent as a function of the number of training episodes. It can be seen how experience is gained (i.e. more episodes are used for training) and eventually plateaus. The colored points correspond to the episodes depicted in Figure 4.



FIG. 6. Damping rates of an instantaneous external excitation that is achieved for different kinds of training parameters. The $y$-axis denotes the reward function used for the specific training according to Table II and the number of neurons in the hidden layer. The trained and untrained (baseline) RL agents are compared against an FIR controller. Note that the negative baseline values are caused by the agents with random coefficients actually exciting the instability. All agents outperform the state-of-the-art FIR controller showing higher damping rates.

As can be seen in Figure 6, the performance of the FIR controller is not constant. This behavior is mainly due to variations in the beam current that, due to its linear nature, affects the action signal amplitude. This is not the case for the RL agent, as it is able to automatically adapt to the variation in beam current. Thus, the trained agent outperformed the FIR controller
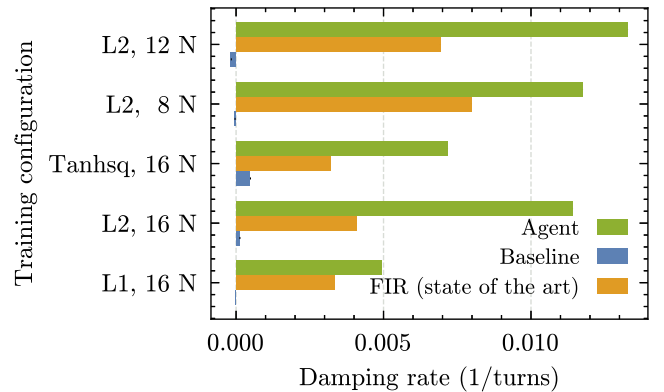
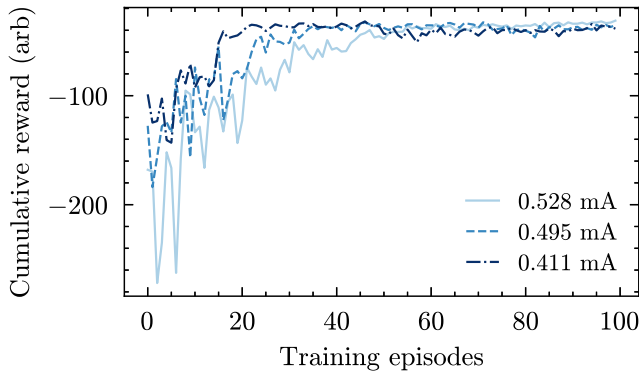in all of our evaluations. Additionally, the untrained RL agent

FIG. 7. Comparison of different training episodes with different beam current conditions, where the final agents achieve almost identical performances.
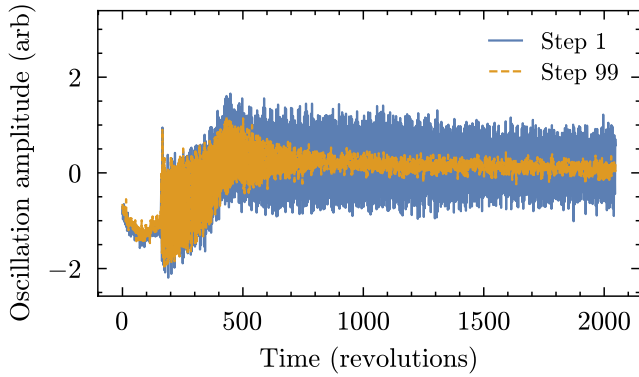


FIG. 8. Evolution of the horizontal position of the beam after an instantaneous external excitation at time $t \approx 175$ revolutions, influenced by an RL agent before and after training. The septum magnet was active, inducing a baseline shift superimposed on the usual exponential decay. Nonetheless, the trained agent (orange) is capable of damping the oscillation compared to the untrained one (blue).
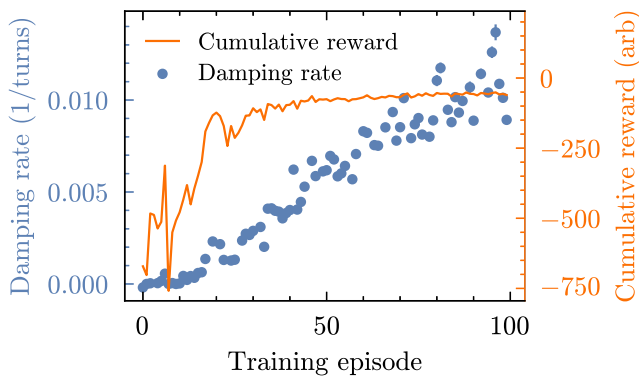


FIG. 9. Damping rate (blue) measured with a different BPM system as a function of the training step, compared with the reward function (orange). Notice the absence of a plateau in the damping rate curve.

is shown as the baseline, with clearly inferior performances when compared to the trained agent and the FIR controller.

Compared to the FIR controllers conventionally used, NN agents are capable of exhibiting non-linear output response. This allows the implementation of more complex policies. Particularly shallow NNs with ReLU activation functions have been shown to behave linearly in some cases. When pure sinusoidal input is fed into a black-box, the non-linear behavior produces harmonics of the fundamental sinusoidal input. The total harmonic distortion plus noise (THD+N), is defined as

$$\text{THD+N} = \frac{\sqrt{A_N^2 + \sum_{i=2}^{\infty} A_i^2}}{A_1}, \qquad (9)$$

where $A_i$ is the amplitude of the $i$-th harmonic, where $i = 1$ is the fundamental, and $A_N$ is the noise amplitude. This metric expresses the amount of non-linear components in the output of a given device. For a linear controller, like a FIR filter, in the case where noise is negligible, THD+N is approximately zero. The THD+N was computed for different amplitudes of a sinusoidal input at the betatron oscillation frequency. The behavior for the L2, 12 N agent is shown in Figure 10. The amount of non-linear harmonic content is consistent, with a steep increase at a level compatible with the noise floor of the signal provided to the agent. This might indicate that the agent is learning to apply more complex actions in the case of high-amplitude, and thus highly penalized, observations.
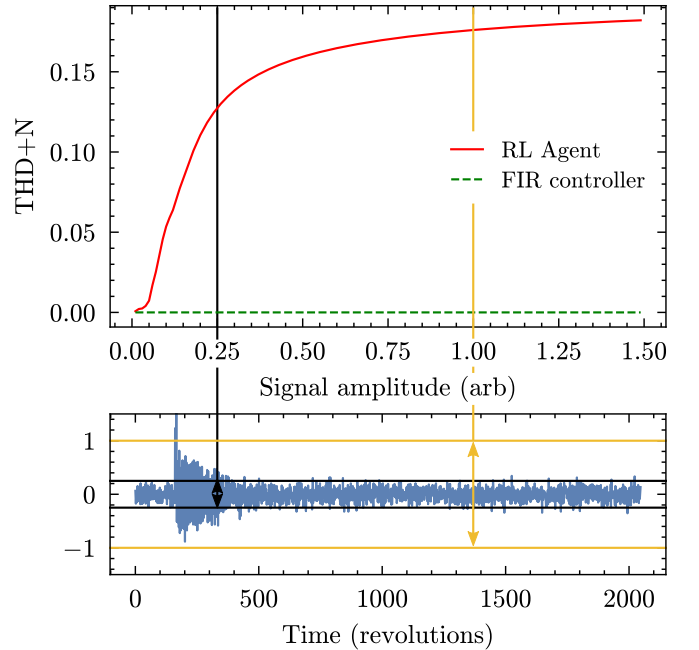


FIG. 10. Total harmonic distortion plus noise (THD+N) as a function of the input signal amplitude for an (L2-norm, 12 neurons) agent. Note that the FIR filter with negligible noise will have a (THD+N) close to 0. In the bottom plot, a training signal is shown, allowing to determine which part of the training episode employs a higher non-linear behavior. The noise level is indicated with the black lines, while an higher amplitude is shown with the yellow lines. Signals above the noise level tend to produce more non-linear actions.

The trade-off between training through interaction with a simulated or real-world environment is an important aspect of

| Training platform | Interaction time | Training time |
|---|---|---|
| Simulation CPU | 17.6 s | 137 s |
| Simulation GPU | 17.6 s | 227 s |
| Online CPU | 0.076 s | 260 s |

TABLE III. Comparison of the time necessary to perform 100 training steps for different training platforms.

the application of RL to large-scale facilities. A simulation-driven approach is sometimes necessary to ensure safety, both of the facility and its personnel, while in other cases, it is dictated by the time necessary to obtain the training dataset. One advantage of real-world training is that trained agents are directly transferable into operation. This is not the case for agents trained on simulation, as their transfer to the real world is potentially hindered by non-modeled phenomena and, more in general, differences between the training and real-world environments. This work presents the opportunity to compare these two approaches. To do so, the time necessary to train an agent through interaction with the accelerator and simulation is reported in Table III. It is important to consider that the online case comprised waiting for a 1 Hz trigger, controlling the kickers. These numbers comprise the overhead of transferring training data, the time necessary for NN back-propagation and the interaction time with the environment. Given both approaches used the same model, the back-propagation time can be assumed to be equal. The interaction time, on the other hand, is 17.6 s for the simulation, while the trial on the accelerator takes 0.076 s or 76 ms. It is relevant to notice how the simulation being employed, and discussed in greater detail in the method section, is lightweight while still performing two orders of magnitudes slower than gathering data on the machine. As such, approximately 50% of the real-world training time is consumed by data-access overhead. This could be reduced in future implementations of the system with several approaches: using higher speed network links paired with lower overhead network protocols, or by accelerating the training directly on the FPGA and AIE array sharing memory with the experience accumulator hardware. In conclusion, for systems with dynamics that is computationally intensive to simulate, the techniques described in this article will greatly improve the time necessary for training. An environment-driven training procedure, i.e. training directly on the real-world task, becomes thus not only possible but also more flexible than a simulation study as the total deployment time would be reduced.

In certain scenarios characterized by rapid dynamics and computationally intensive simulations demanding high-performance computing clusters, it may be conceivable that training directly on the accelerator consumes less energy than utilizing simulations. Such a possibility could significantly influence the sustainability of ML methodologies.

## VI. CONCLUSIONS

In this study, we introduced the experience accumulator architecture and the training time reward definition, marking a significant step in implementing real-world on-the-edge environment learning for RL-based controllers. Particularly in scenarios where simulation costs are prohibitive and data generation rates are high, this methodology emerges as a promising solution, enabling the deployment of RL controllers. Our application of this approach to real-time control of particle accelerator dynamics yielded inference latency of a few microseconds.

Utilizing advanced heterogeneous computing platforms such as KINGFISHER, our presented implementation facilitates the deployment of plug-and-play RL systems operating at microsecond latency scales. This opens the door to intelligent control of ultra-fast non-linear dynamics [47] in systems such as particle accelerators and fusion experiments. The efficacy of these capabilities was demonstrated through the control of the HBO at the accelerator test facility KARA, resulting in a functional controller comparable and outperforming state-of-the-art FIR controllers used during standard operation phases of synchrotron light sources.

## VII. ACKNOWLEDGEMENTS

## VIII. CODE AND DATA AVAILABILITY

All data and code used for this work are provided by the authors upon reasonable request.

[1] B. Mustapha and P. Ostroumov, Optimization algorithms for accelerator physics problems, Proceedings of ICAP09, San Francisco, CA FR1IOPK01 (2009).

[2] A. Scheinker, X. Pang, and L. Rybarcyk, Model-independent particle accelerator tuning, Physical Review Special Topics - Accelerators and Beams **16**, 10.1103/physrevstab.16.102803 (2013).

[3] A. Hofler, B. Terzić, M. Kramer, A. Zvezdin, V. Morozov, Y. Roblin, F. Lin, and C. Jarvis, Innovative applications of genetic algorithms to problems in accelerator physics, Physical Review Special Topics - Accelerators and Beams **16**, 10.1103/physrevstab.16.010101 (2013).

[4] Roussel, Ryan, Bartnik, Adam, Edelen, Auralee, and Mayes, Christopher, Xopt: A simplified framework for optimiza-

tion of accelerator problems using advanced algorithms 10.18429/JACOW-IPAC2023-THPL164 (2023).

[5] A. Edelen, N. Neveu, M. Frey, Y. Huber, C. Mayes, and A. Adelmann, Machine learning for orders of magnitude speedup in multiobjective optimization of particle accelerator systems, Physical Review Accelerators and Beams **23**, 10.1103/physrevaccelbeams.23.044601 (2020).

[6] R. Roussel, A. Hanuka, and A. Edelen, Multiobjective bayesian optimization for online accelerator tuning, Physical Review Accelerators and Beams **24**, 10.1103/physrevaccelbeams.24.062801 (2021).

[7] J. Kaiser, C. Xu, A. Eichler, A. Santamaria Garcia, O. Stein, E. Bründermann, W. Kuropka, H. Dinter, F. Mayet, T. Vinatier, F. Burkart, and H. Schlarb, Reinforcement learning-trained optimisers and bayesian optimisation for online particle accelerator tuning, Scientific Reports **14**, 10.1038/s41598-024-66263-y (2024).

[8] J. Degrave, F. Felici, J. Buchli, M. Neunert, B. Tracey, F. Carpanese, T. Ewalds, R. Hafner, A. Abdolmaleki, D. de las Casas, C. Donner, L. Fritz, C. Galperti, A. Huber, J. Keeling, M. Tsimpoukelli, J. Kay, A. Merle, J.-M. Moret, S. Noury, F. Pesamosca, D. Pfau, O. Sauter, C. Sommariva, S. Coda, B. Duval, A. Fasoli, P. Kohli, K. Kavukcuoglu, D. Hassabis, and M. Riedmiller, Magnetic control of tokamak plasmas through deep reinforcement learning, Nature **602**, 414 (2022).

[9] J. Kaiser, C. Xu, A. Eichler, A. Santamaria Garcia, O. Stein, E. Bründermann, W. Kuropka, H. Dinter, F. Mayet, T. Vinatier, F. Burkart, and H. Schlarb, Learning to do or learning while doing: Reinforcement learning and bayesian optimisation for online continuous tuning (2023).

[10] N. Madysa, F. Velotti, N. Biancacci, R. Alemany-Fernández, V. Kain, and B. Goddard, Automated intensity optimisation using reinforcement learning at LEIR, JACoW IPAC **2022**, 941 (2022).

[11] N. Bruchon, G. Fenu, G. Gaio, M. Lonza, F. H. O'Shea, F. A. Pellegrino, and E. Salvato, Basic reinforcement learning techniques to control the intensity of a seeded free-electron laser, Electronics **9**, 781 (2020).

[12] V. Kain, S. Hirlander, B. Goddard, F. M. Velotti, G. Z. Della Porta, N. Bruchon, and G. Valentino, Sample-efficient reinforcement learning for cern accelerator control, Phys. Rev. Accel. Beams **23**, 124801 (2020).

[13] F. O'Shea, N. Bruchon, and G. Gaio, Policy gradient methods for free-electron laser and terahertz source optimization and stabilization at the fermi free-electron laser at Elettra, Physical Review Accelerators and Beams **23**, 122802 (2020).

[14] J. St. John, C. Herwig, D. Kafkes, J. Mitrevski, W. A. Pellico, G. N. Perdue, A. Quintero-Parra, B. A. Schupbach, K. Seiya, N. Tran, M. Schram, J. M. Duarte, Y. Huang, and R. Keller, Real-time artificial intelligence for accelerator control: A study at the Fermilab booster, Phys. Rev. Accel. Beams **24**, 104601 (2021).

[15] K. Shin and P. Ramanathan, Real-time computing: a new discipline of computer science and engineering, Proceedings of the IEEE **82**, 6 (1994).

[16] S. Hirlaender, L. Lamminger, G. Zevi Della Porta, and V. Kain, Ultra fast reinforcement learning demonstrated at CERN AWAKE, JACoW **IPAC2023**, THPL038 (2023).

[17] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction* (A Bradford Book, Cambridge, MA, USA, 2018).

[18] M. Rothmann and M. Porrmann, A survey of domain-specific architectures for reinforcement learning, IEEE Access **10**, 13753 (2022).

[19] A. R. Baranwal, S. Ullah, S. S. Sahoo, and A. Kumar, Relaccs: A multilevel approach to accelerator design for reinforcement learning on fpga-based systems, IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems **40**, 1754 (2021).

[20] S. S. Sahoo, A. R. Baranwal, S. Ullah, and A. Kumar, Memorel: A memory-oriented optimization approach to reinforcement learning on fpga-based embedded systems, in *Proceedings of the 2021 on Great Lakes Symposium on VLSI*, GLSVLSI '21 (Association for Computing Machinery, New York, NY, USA, 2021) p. 339–346.

[21] L. M. D. Da Silva, M. F. Torquato, and M. A. C. Fernandes, Parallel implementation of reinforcement learning q-learning technique for fpga, IEEE Access **7**, 2782 (2019).

[22] Y. Meng, S. Kuppannagari, R. Rajat, A. Srivastava, R. Kannan, and V. Prasanna, Qtaccel: A generic fpga based design for q-table based reinforcement learning accelerators, in *2020 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)* (2020) pp. 107–114.

[23] S. Spanò, G. C. Cardarilli, L. Di Nunzio, R. Fazzolari, D. Giardino, M. Matta, A. Nannarelli, and M. Re, An efficient hardware implementation of reinforcement learning: The q-learning algorithm, IEEE Access **7**, 186340 (2019).

[24] S. Shao, J. Tsai, M. Mysior, W. Luk, T. Chau, A. Warren, and B. Jeppesen, Towards hardware accelerated reinforcement learning for application-specific robotic control, in *2018 IEEE 29th International Conference on Application-specific Systems, Architectures and Processors (ASAP)* (2018) pp. 1–8.

[25] C.-W. Hu, J. Hu, and S. P. Khatri, Td3lite: Fpga acceleration of reinforcement learning with structural and representation optimizations, in *2022 32nd International Conference on Field-Programmable Logic and Applications (FPL)* (2022) pp. 79–85.

[26] H. Cho, P. Oh, J. Park, W. Jung, and J. Lee, FA3c, in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems* (ACM, 2019).

[27] J. Su, J. Liu, D. B. Thomas, and P. Y. Cheung, Neural network based reinforcement learning acceleration on FPGA platforms, ACM SIGARCH Computer Architecture News **44**, 68 (2017).

[28] H. Watanabe, M. Tsukada, and H. Matsutani, An fpga-based on-device reinforcement learning approach using online sequential learning, in *2021 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)* (2021) pp. 96–103.

[29] M.-J. Li, A.-H. Li, Y.-J. Huang, and S.-I. Chu, Implementation of deep reinforcement learning, in *Proceedings of the 2nd International Conference on Information Science and Systems*, ICISS '19 (Association for Computing Machinery, New York, NY, USA, 2019) p. 232–236.

[30] Y. Meng, S. Kuppannagari, and V. Prasanna, Accelerating proximal policy optimization on cpu-fpga heterogeneous platforms, in *2020 IEEE 28th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)* (2020) pp. 19–27.

[31] G. Book, A. Traue, P. Balakrishna, A. Brosch, M. Schenke, S. Hanke, W. Kirchgässner, and O. Wallscheid, Transferring online reinforcement learning for electric motor control from simulation to real-world experiments, IEEE Open Journal of Power Electronics **2**, 187 (2021).

[32] L. Scomparin *et al.*, KINGFISHER: A Framework for Fast Machine Learning Inference for Autonomous Accelerator Systems, JACoW **IBIC2022**, 151 (2022).

[33] Xilinx VCK190 Evaluation Board, `https://www.xilinx.com/products/boards-and-kits/vck190.html`.

[34] Aurora 64B/66B protocol specification (SP011), `https://docs.amd.com/v/u/en-US/aurora_64b66b_protocol_spec_sp011`.

[35] Experimental physics and industrial control system (epics), `https://epics.anl.gov/`.

[36] M. E. O'Neill, *PCG: A Family of Simple Fast Space-Efficient Statistically Good Algorithms for Random Number Generation*, Tech. Rep. HMC-CS-2014-0905 (Harvey Mudd College, Claremont, CA, 2014).

[37] M. Caselle, KAPTURE-2. a picosecond sampling system for individual THz pulses with high repetition rate, J. Instrum. **12**, C01040.

[38] M. Caselle, L. A. Perez, M. Balzer, T. Dritschler, A. Kopmann, H. Mohr, L. Rota, M. Vogelgesang, and M. Weber, A high-speed daq framework for future high-level trigger and event building clusters, Journal of Instrumentation **12** (03), C03015.

[39] Dimtel BPMH-20-2G hybrid, `https://www.dimtel.com/products/bpmh`.

[40] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, Proximal policy optimization algorithms (2017), arXiv:1707.06347 [cs.LG].

[41] A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann, Stable-baselines3: Reliable reinforcement learning implementations, Journal of Machine Learning Research **22**, 1 (2021).

[42] M. Dehler, V. Schlott, D. Bulfone, M. Lonza, and R. Ursic, Current status of the elettra/sls transverse multibunch feedback, in *Proc. EPAC* (2000) pp. 1894–1896.

[43] Dimtel iGp12 bunch-by-bunch feedback system, `https://www.dimtel.com/products/igp12`.

[44] M. Lonza and Presented By H. Schmickler, Multi-bunch feedback systems, CERN Yellow Reports: School Proceedings , Vol 3 (2017): Proceedings of the CAS (2017).

[45] M. Towers, J. K. Terry, A. Kwiatkowski, J. U. Balis, G. d. Cola, T. Deleu, M. Goulão, A. Kallinteris, A. KG, M. Krimmel, R. Perez-Vicente, A. Pierré, S. Schulhoff, J. J. Tai, A. T. J. Shen, and O. G. Younis, Gymnasium (2023).

[46] E. Blomley, *Investigation and Control of Beam Instabilities at the Karlsruhe Research Accelerator using a 3-D Digital Bunch-by-Bunch Feedback System*, Ph.D. thesis, Karlsruher Institut für Technologie (KIT) (2021), 54.11.11; LK 01.

[47] Scomparin, Luca, Mochihashi, Akira, Santamaria Garcia, Andrea, Kopmann, Andreas, Mueller, Anke-Susanne, Xu, Chenran, Blomley, Edmund, Bruendermann, Erik, Steinmann, Johannes, Becker, Jürgen, Weber, Marc, Schuh, Marcel, Caselle, Michele, and Dritschler, Timo, Preliminary results on the reinforcement learning-based control of the microbunching instability, in *Proc. 15th International Particle Accelerator Conference*, IPAC'24 - 15th International Particle Accelerator Conference No. 15 (JACoW Publishing, Geneva, Switzerland, 2024) pp. 1808–1811.