



Machine Learning Activities at KIT

Andrea Santamaria Garcia, Edmund Blomley, Erik Bründermann, Michele Caselle, Stephan Kötter, Luca Scomparin, Johannes L. Steinmann, C. Xu, Anke-Susanne Müller

ACCLAIM Meeting (Jena, 3rd July 2023)



www.kit.edu

ML-related events



Creation of the Collaboration on Reinforcement Learning for Autonomous Accelerators (RL4AA)

- Kick-off with workshop organized at KIT
- Expert lectures on reinforcement learning
- Real application to accelerator tutorials
- Advanced discussion sessions

 \rightarrow Proceedings to be published soon

 \rightarrow Next year's workshop by Uni Salzburg (5-7 Feb. 2024)



https://indico.scc.kit.edu/event/3280/overview https://github.com/RL4AA/RL4AA23 https://rl4aa.github.io/

AI HERO Hackathon on Energy Efficient AI

- Training machine learning models on GPU clusters costs energy
- With the rise of larger and larger models, training should not be indiscriminate
- Many strategies: use pre-trained models, mixed precision, simplification of input, simplest model, smart hyperparameter tuning, train fewer epochs, compile your code, NVIDIA tools for fast inference, no idle cores, etc...





First RL algorithm online training and running on hardware in accelerators L. Scomparin



Reinforcement Learning





First RL algorithm online training and running on hardware in accelerators

L. Scomparin

- Agent: Vanilla PPO from Stable Baselines 3
- Actor & critic architecture: 8-16-1
- Reward: metric of the beam position (low as possible)
- **Observation:** last 8 BPM samples
- Strategy:
 - 1. Agent acts during 2000 turns (0.74 ms)
 - 2. Agent stops and is re-trained in a CPU (~2 s)
 - 3. New weights are sent to Versal board and agent starts again



Reinforcement Learning



NNs coded in Versal AIE Only forwad pass





30.06.23 Andrea Santamaria Garcia

First RL algorithm online training and running on hardware in accelerators L. Scomparin



Reinforcement Learning



First detailed comparison of BO and RL in a



real accelerator J. Kaiser, C. Xu

Reinforcement Learning

Bayesian Optimization

https://arxiv.org/abs/2306.03739

[Submitted on 6 Jun 2023]

Learning to Do or Learning While Doing: Reinforcement Learning and Bayesian Optimisation for Online Continuous Tuning

Jan Kaiser, Chenran Xu, Annika Eichler, Andrea Santamaria Garcia, Oliver Stein, Erik Bründermann, Willi Kuropka, Hannes Dinter, Frank Mayet, Thomas Vinatier, Florian Burkart, Holger Schlarb

Online tuning of real-world plants is a complex optimisation problem that continues to require manual intervention by experienced human operators. Autonomous tuning is a rapidly expanding field of research, where learning-based methods, such as Reinforcement Learning-trained Optimisation (RLO) and Bayesian optimisation (BO), hold great promise for achieving outstanding plant performance and reducing tuning times.

HELMHOLTZAI

"Machine Learning Toward Autonomous Accelerators"

- **Task**: focus and position the electron beam
- Actuators: 3 quadrupole magnets + 2 corrector magnets
- Observation: beam image on the diagnostic screen



First detailed comparison of BO and RL in a

real accelerator J. Kaiser, C. Xu



Reinforcement Learning

Bayesian Optimization





First detailed comparison of BO and RL in a real accelerator J. Kaiser, C. Xu





Use **reinforcement learning**-trained optimisation, when the tuning task is repeated often and the upfront engineering effort pays off through the better performance

Use **Bayesian optimisation**, when the tuning task is only performed a few times, e.g. during commissioning, where the upfront engineering effort does not pay-off and BO experts can be present



Bayesian optimization algorithm transferred to EuXFEL



C. Xu

Time to inject to KARA cut in half with automated tuning by BO algorithm https://doi.org/10.1103/PhysRevAccelBeams.26.034601 Emitted THz radiation at FLUTE optimized with parallel BO in simulation https://doi.org/10.18429/JACoW-IPAC2022-WEPOMS023

Transfer of algorithm to EuXFEL to tune SASE emission <u>https://www.ipac23.org/preproc/pdf/THPL028.pdf</u>





Bayesian Optimization





Lattice agnostic RL \rightarrow Generalizable RL \bigcirc .xu





Oc = OCELOT (with space charge)

Active learning to design vFFAs

A. Oeftiger, A. Santamaria Garcia, S. Hirländer, J.-B. Lagrange

Why a data-driven characterization study?

- Lattice with strong transverse coupling, nonplanar orbits
- Magnetic exponential fields = zero chroma + strong nonlinearity $B_{x,y,z} \propto B_0 \exp(m z)$
- Many lattice configurations do not yield closed orbits!

→ Betatron tune adjustment is crucial to avoid resonances (space charge)
→ Lattice design for maximum dynamic aperture has been trial-and-error in first design studies (difficult)

Now explore lattice parameter space efficiently, guided by **active learning** = iterative supervised learning

Zero momentum compaction







- Iterative training and rejection sampling give up to 85% closed orbits per iteration
- Gather 170k simulated lattices with overall >25% closed orbits

closed orbit dense regions

More simulations were launched in

Active learning to design vFFAs

A. Oeftiger, A. Santamaria Garcia, S. Hirländer, J.-B. Lagrange

Training to predict dynamic aperture

ACTIVE LEARNING





Active learning provided overall 200k lattice configurations with 25% closed orbits and well resolved high-DA regions:



S. Kötter

Goal: control the electron phase space in a linear accelerator

- Can be achieved by shaping the photoinjector laser pulse
- Modulate pulse via a spatial light modulator (SLM)
- Find the phase modulation that results in a target pulse shape
- Current status: test setup for transverse modulation

Evolve last year's experiment by ...

- 1. Switching to the Ti:Sa 800 nm photoinjector laser system
- 2. Including a compressor in the optical path
- 3. Mapping the camera intensity directly to the phase modulation via a deep neural network

Karlsruhe Institute of Technology Laser pulse shaping with Spatial Light Modulators and convolutional neural networks





Test setup for the SLM project in the FLUTE cleanroom

S. Kötter

Input:



Neural Network Architecture

- Is deep. Shallow networks did not work at all!
- Currently works with EfficientNetV1 but ResNet50 also works

Output:

- Zernike Amplitudes
- Size: 14
- Normalized to [0, 1]



[1] V. Krishna Adithya et al., EffUnet-SpaGen: An Efficient and Spatial Generative Approach to Glaucoma Detection, Journal of Imaging vol 7 (2021).



Each tuple consists of ...

S. Kötter

•

- 1. a random superpositions of the first 14 Zernike polynomials (reduces degrees of freedom)
- 2. the corresponding camera image





[1] Nschloe for wikipedia, cropped, license: https://creativecommons.org/licenses/by-sa/4.0/deed.en (2020).



S. Kötter



Preliminary results

Test data:

- Predicted phase patterns match originals reasonably well
- Inference works reasonably well!



Experimental application:

- Seems to work in general
- But is probably not useful yet



Conclusion: works but is not useful yet!

Outlook: improve experimental setup and possibly separate neural network!



Thank you for your attention! What questions do you have for me?



Karlsruhe Institute of Technology

Dr. Andrea Santamaria Garcia Al4Accelerators team leader

andrea.santamaria@kit.edu https://twitter.com/ansantam https://www.linkedin.com/in/ansantam/ https://github.com/ansantam