

# Automatic Serial Femtosecond Crystallography Online Analysis with Reinforcement Learning



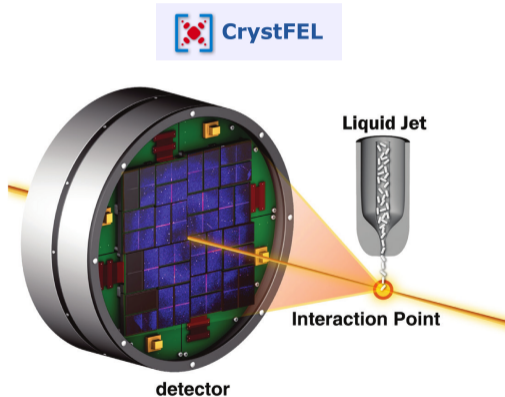
D. Ferreira de Lima   A. Davtyan   L. Gelisio

European XFEL

15 October 2021

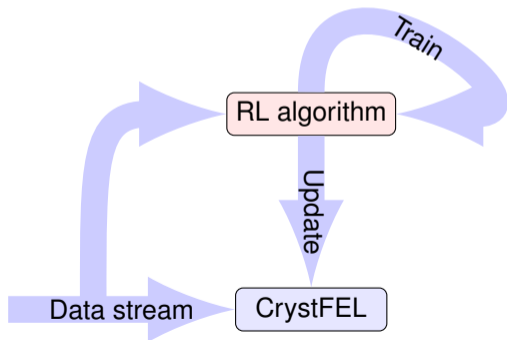
# Serial Femtosecond Crystallography analysis

- Serial Femtosecond Crystallograph's pipeline:
  - Pre-select relevant frames.
  - Identify Bragg peaks.
  - Identify Miller indices corresponding to peaks and crystal orientation.
  - Integrate intensity of the Bragg peaks.
- Requires tuning many parameters depending on the sample and experimental conditions.
  - Can we automatize the parameter tuning?



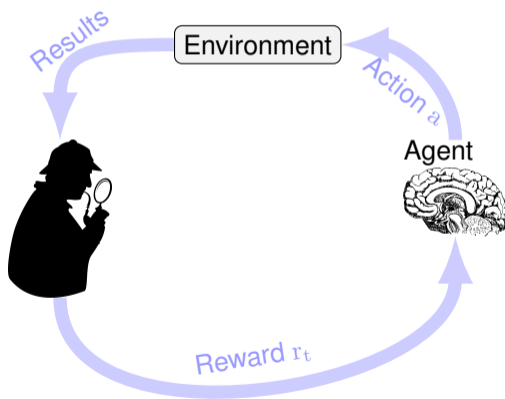
(Image from the SPB TDR, TR-2013-004)

# Automatizing SFX



- Collect a sample of the incoming data to search optimal parameters.
- While the parameter search runs, the standard pipeline runs in parallel.
  - The parameter search explores randomly other possible parameter values.
- When improved parameters are found, update the parameters in the pipeline to improve the results as data is taken.

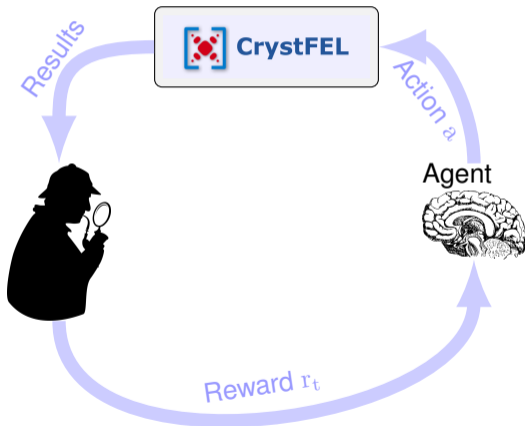
# Reinforcement Learning



- Environment  $\rightarrow$  analysis pipeline.
- Model-free RL  $\rightarrow$  environment-independent.
  - Input = last + **noise** (*perform action*).
  - Was there an improvement (*collect reward*)?
  - Update agent.
- Objective: maximize total returns  $G$  after  $T$  attempts.
  - At  $t = T$ , reset state to the best.

$$G = \sum_{t=1}^T \gamma^t r_t$$

# Online SFX tuning



- What is the environment?
  - The environment is the algorithm to be optimized: CrystFEL.
- What are the actions?
  - Actions are changes in the parameter values set by CrystFEL.
  - A current state (parameter set) is kept.
- What are the rewards?
  - Change in fraction of indexed frames.
  - Positive reward means the number of indexed frames increased.

# What is the agent?

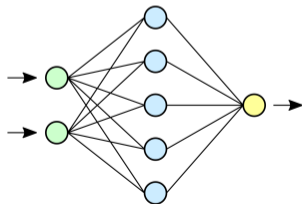
- The agent is just a parametrized function.
  - It takes the current parameter set and outputs the change in parameter values needed to optimize the return.
  - We fit the parameters of the function.
- Neural networks are universal function approximators.

**Theorem 1.** *Let  $\sigma$  be any continuous discriminatory function. Then finite sums of the form*

$$G(x) = \sum_{j=1}^N \alpha_j \sigma(y_j^T x + \theta_j) \quad (2)$$

*are dense in  $C(I_n)$ . In other words, given any  $f \in C(I_n)$  and  $\varepsilon > 0$ , there is a sum,  $G(x)$ , of the above form, for which*

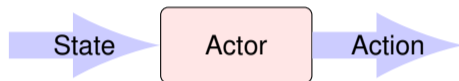
$$|G(x) - f(x)| < \varepsilon \quad \text{for all } x \in I_n.$$



G. Cybenko, Math. Control Signals Systems (1989) **2**, 303 – 314

## Actor-critic models

- One possible RL setup: actor-critic model.
- After randomly exploring actions around the current optimal:
  - Critic  $C(s)$  fits the obtained returns  $G$ , given the state  $s$ .
    - Advantage function definition:  $A(s) = C(s) - G$
    - Critic minimizes  $\mathcal{L}_C = [A(s)]^2$
  - Actor optimizes returns, based on the critic's predictions.



# Advantage Actor Critic (A3C/A2C) and ACKTR

Machine Learning, 8, 229-256 (1992)  
 © 1992 Kluwer Academic Publishers, Boston. Manufactured in The Netherlands.

## How to optimize the actor?

- Maximize returns  $\rightarrow$  move parameters in the direction of  $\nabla_{\theta} \mathbb{E}_{\pi(\tau|\theta)} [G(\tau)]$ .

## Williams (1992) showed that:

- $\mathbb{E}_{\pi(\tau|\theta)} [\nabla_{\theta} \log \pi(\tau|\theta) G(\tau)]$  is an unbiased estimator for  $\nabla_{\theta} \mathbb{E}_{\pi(\tau|\theta)} [G(\tau)]$ .
- Probability of taking a set of actions  $\tau = \{(s_1, a_1), \dots, (s_n, a_n)\}$  is  $\pi(\tau)$ .
- Parameters of the neural network  $\theta$ .

## In A3C: substitute G with A to reduce variance.

- ACKTR: Use second-order optimization algorithms to improve convergence speed.

## Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning

RONALD J. WILLIAMS [rjw@corwin.ccs.northeastern.edu](mailto:rjw@corwin.ccs.northeastern.edu)  
 College of Computer Science, 161 CN, Northeastern University, 360 Huntington Ave., Boston, MA 02115

## Asynchronous Methods for Deep Reinforcement Learning

Volodymyr Mnih <sup>1</sup>	<a href="mailto:VMNH@GOOGLE.COM">VMNH@GOOGLE.COM</a>
Adria Puigdomènech Badia <sup>1</sup>	<a href="mailto:ADRIAP@GOOGLE.COM">ADRIAP@GOOGLE.COM</a>
Mehdi Mirza <sup>1,2</sup>	<a href="mailto:MIRZAMOH@HQS.MONTREAL.CA">MIRZAMOH@HQS.MONTREAL.CA</a>
Alex Graves <sup>1</sup>	<a href="mailto:GRAVESA@GOOGLE.COM">GRAVESA@GOOGLE.COM</a>
Tim Harley <sup>1</sup>	<a href="mailto:THARLEY@GOOGLE.COM">THARLEY@GOOGLE.COM</a>
Timothy P. Lillicrap <sup>1</sup>	<a href="mailto:TLILICRA@GOOGLE.COM">TLILICRA@GOOGLE.COM</a>
David Silver <sup>1</sup>	<a href="mailto:DAVIDSILVER@GOOGLE.COM">DAVIDSILVER@GOOGLE.COM</a>
Koray Kavukcuoglu <sup>1</sup>	<a href="mailto:KORAYK@GOOGLE.COM">KORAYK@GOOGLE.COM</a>

<sup>1</sup> Google DeepMind  
<sup>2</sup> Montreal Institute for Learning Algorithms (MILA), University of Montreal

## Scalable trust-region method for deep reinforcement learning using Kronecker-factored approximation

Yuhua Wu<sup>\*</sup>  
 University of Toronto  
 Vector Institute  
[ywu@cs.toronto.edu](mailto:ywu@cs.toronto.edu)

Elman Mansimov<sup>\*</sup>  
 New York University  
[mansimov@cs.nyu.edu](mailto:mansimov@cs.nyu.edu)

Shun Liao  
 University of Toronto  
 Vector Institute  
[sliao3@cs.toronto.edu](mailto:sliao3@cs.toronto.edu)

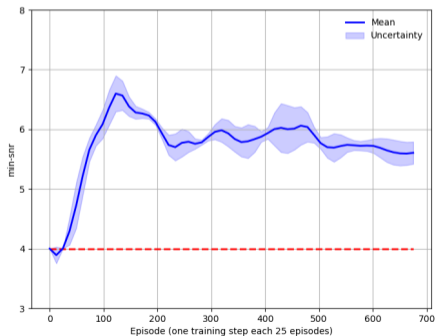
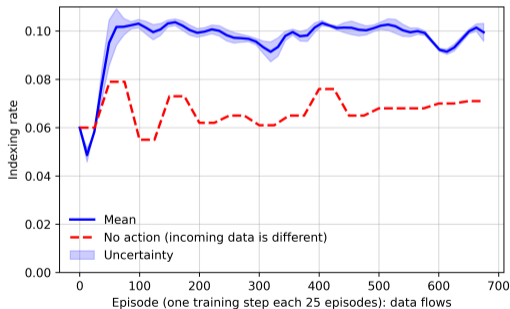


# Tests

- Tested the method with a dataset of:
  - hen egg white lysozyme;
  - collected on an AGIPD detector at the SPD experiment.
- Used a cache of 1000 images to explore parameter space in the RL algorithm.
- Optimization on 10 parallel cores.
- RL algorithm takes only  $\sim 1$  second out of 1 minute: most time spent on the analysis pipeline itself.

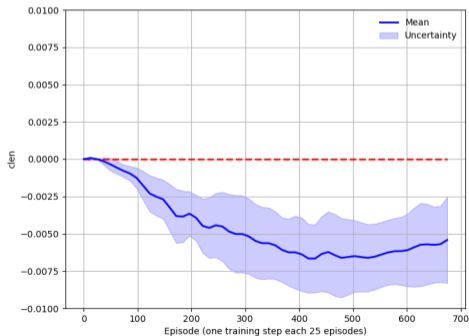
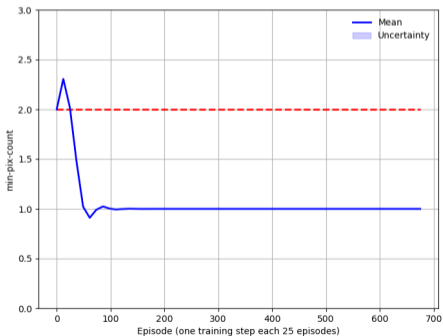
## Results

- Indexed fraction increases when compared to maintaining the initial setup.
- The minimum signal-to-noise ratio parameter is tuned away from the initial value.



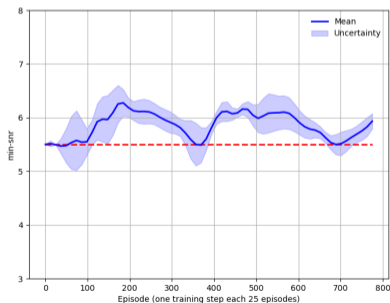
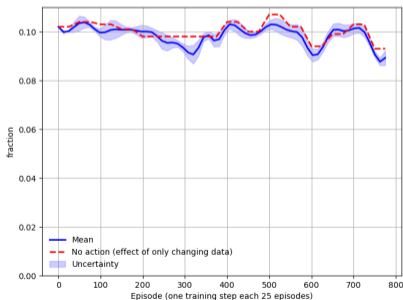
## Other tuned variables

- The detector to sample distance and minimum number of hits can also be tuned.
- Uncertainty band: root-mean-squared-error of 4 random NN initializations.



## Stability test

- What if we already start close to the ideal conditions?
- Focus on speed → using a small number of images to explore parameters.
  - Could increase the number of frames for smaller uncertainty on the rewards.



# Summary

- Apply on-policy reinforcement learning algorithm to optimize SFX pipeline online.
  - Clear goal: maximize number of indexed frames.
- Main concept:
  - Set initial parameters and try to change it randomly.
  - Fit the obtained rewards with the critic.
  - Fit actor to optimize rewards.
  - Automatically adapt parameters as the data flows.
- System tested offline with further tests on new datasets coming.