

Adaptive optics control with multi-agent model-free reinforcement learning: supplement

B. POU,^{1,2,*} F. FERREIRA,³ E. QUINONES,¹ D. GRATADOUR,^{3,4} AND M. MARTIN²

¹*Barcelona Supercomputing Center (BSC), C. Jordi Girona, 29, 08034 Barcelona, Spain*

²*Computer Science Department, Universitat Politècnica de Catalunya (UPC), C. Jordi Girona, 31, 08034 Barcelona, Spain*

³*LESIA, Observatoire de Paris, Université PSL, CNRS, Sorbonne Université, Univ. Paris Diderot, Sorbonne Paris Cité, 5 place Jules Janssen, 92195 Meudon, France*

⁴*Research School of Astronomy and Astrophysics, Australian National University, Canberra, ACT 2611, Australia*

*bartomeu.poumulet@bsc.es

This supplement published with Optica Publishing Group on 14 January 2022 by The Authors under the terms of the [Creative Commons Attribution 4.0 License](#) in the format provided by the authors and unedited. Further distribution of this work must maintain attribution to the author(s) and the published article's title, journal citation, and DOI.

Supplement DOI: <https://doi.org/10.6084/m9.figshare.17714879>

Parent Article DOI: <https://doi.org/10.1364/OE.444099>

Adaptive Optics control with Multi-Agent Model-Free Reinforcement Learning: supplemental document

This supplementary material provides an additional explanation on the B_{tt} modes and Soft Actor Critic. Moreover, it shows the pseudocode of the training process and additional RL hyperparameters used in our simulations.

1. ADDITIONAL EXPLANATION ON THE B_{TT} MODAL BASIS

Let's consider the influence functions, IF_i , that denote how acting on actuator i will affect the DM shape. Any phase represented by the DM, Φ_{DM} , can be obtained through a linear combination of these influence functions:

$$\Phi_{DM} = \sum_i f_i IF_i. \quad (S1)$$

Because of coupling between actuators, this basis is not orthogonal and, on top of this, can reproduce piston. While piston is not seen by a SH-WFS and does not impact image quality, it has to be removed to avoid drifting towards actuators saturation. Additionally, producing tip and tilt with the DM is usually not practical since it requires a lot of stroke, especially for actuators at the edge of the pupil. In this case as well, these two modes are usually filtered out of the degrees of freedom of the DM and compensated using a dedicated subsystem (either using a flat mirror or the DM itself mounted on a dedicated tip-tilt mount).

The B_{tt} basis uses the influence functions to build another representation of possible DM shapes through orthonormal modes, spanning the full DM actuators space, but from which piston and tip-tilt are filtered out. The method is described in appendix A of [1] and recalled hereafter. Using the spatial covariance matrix of the actuators of Eq. (S1), one can retrieve the coefficients of the best fit for piston and tip-tilt on Eq. (S1) and create a set of generators that do not produce those modes. The generators matrix is then diagonalised and normalised using the corresponding eigenvalues to obtain a matrix B describing a set of orthonormal modes from which piston and tip-tilt have been removed. Finally, the tip-tilt modes are reintroduced in the basis by adding two extra columns to the matrix with a single non-zero value on the diagonal obtained as the inverse of μ , the eigenvalues of $IF_{tt}^t IF_{tt}$ where IF_{tt} are the influence functions of tip-tilt. The final basis B_{tt} basis is expressed as:

$$B_{tt} = \begin{pmatrix} & 0 & 0 \\ & \vdots & \vdots \\ B & 0 & 0 \\ 0 & \dots & 0 & 1/\mu & 0 \\ 0 & \dots & 0 & 0 & 1/\mu \end{pmatrix}. \quad (S2)$$

The modes are normalised over the pupil surface, S , as $\frac{1}{S} \int B_{tt,i}^2 dS = 1 \mu m^2$. With this basis the phase will be computed as:

$$\Phi_{DM} = \sum_i v_i B_{tt,i}, \quad (S3)$$

where $B_{tt,i}$ are the vectors of the B_{tt} basis.

2. SOFT ACTOR CRITIC

Soft Actor Critic (SAC) [2], is considered one of the state-of-the-art methods in RL. SAC considers a stochastic policy, π , described as a probability distribution. We can introduce SAC by taking a

closer look at its name.

1. The "soft" naming is given by the author's to its method because it uses a maximum entropy version of the return, J^{MaxEnt} , where a trade-off between reward and information entropy of the policy (denoted as $H(\pi)$) is introduced:

$$J^{MaxEnt}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{T-1} \gamma^t (r_{t+1} + \beta H(\pi(\cdot|s_t))) \right], \quad (S4)$$

with the entropy of a policy defined as:

$$H(\pi(\cdot|s)) = \mathbb{E}_{a \sim \pi(\cdot|s)} [-\log(\pi(a|s))], \quad (S5)$$

where the notation $a \sim \pi$ indicates that the action, a , has been sampled from the policy π and $\beta \geq 0$ is a hyperparameter that determines the strength of the entropy against the reward. Notice that a value of $\beta = 0$ returns the original definition of J . The entropy works as a regularisation in the learning procedure which increases the robustness with different initialisations and variations of hyperparameters which many other RL algorithms are culprit of.

2. Regarding actor critic, *actor* denotes the policy, π , which in SAC is stochastic, and *critic* denotes a function that predicts J^{MaxEnt} starting on a state s , executing action a , and then following the current policy afterwards denoted as $Q(s, a)$. The idea of the critic is to guide the actor to policies that maximise Eq. (S4).

In Deep RL, actor critic methods are functions parameterised with neural network weights, ϕ and θ respectively, henceforth we can write: $\pi_{\phi}(s)$ and $Q_{\theta}(s, a)$. The neural network parameterisation of both actor and critic done via means of neural networks is a standard procedure in RL to overcome high dimensional state spaces.

Consider a trajectory, τ , of samples $\tau = (s_1, a_1, r_1, s_2, a_2, r_2, \dots)$ obtained by interacting with the environment. Minimising the following loss function $L(\theta)$, provided by comparing critic prediction $Q_{\theta}(s, a)$ with the soft Bellman equation [2]

$$L(\theta) = \mathbb{E}_{\substack{s, a, r, s' \sim \tau, \\ a' \sim \pi_{\phi}(\cdot|s')}} \left[\frac{1}{2} \left(Q_{\theta}(s, a) - \overbrace{(r + \gamma(Q_{\theta}(s', a') - \beta \log(\pi_{\phi}(a'|s'))))}^{\text{Soft Bellman Eq.}} \right)^2 \right], \quad (S6)$$

provides an iterative procedure to learn θ . Notice that a' is sampled given a s' obtained in the trajectory.

In SAC, the actor represents a probability distribution which is usually represented as a normal distribution. The policy's neural network, π_{ϕ} , will output two values: a mean vector, μ_{ϕ} , and a standard deviation vector σ_{ϕ} dependant on parameters ϕ . This two values are the components of the Normal distribution $\mathcal{N}(\mu_{\phi}, \sigma_{\phi}^2)$. The actor should be learned in a similar iterative manner as the critic by minimising the difference between the current policy distribution and the exponential of the critic estimation which has been proved to lead to better policies [2] as:

$$L(\phi) = \mathbb{E}_{s \sim \tau, a \sim \pi_{\phi}(\cdot|s)} [\beta \log(\pi_{\phi}(a|s)) - Q_{\theta}(s, a)]. \quad (S7)$$

However, estimating the gradient of the loss above is intractable because we are sampling the actions with the same parameters that we are trying to update. To circumvent this problem, SAC uses the so-called reparameterisation trick: sampling an action with the policy parameters, ϕ , is rewritten as sampling from a deterministic function dependant on ϕ and a noise variable, ϵ , independent of ϕ , which is represented with the same distribution type as the one used to sample actions: a Normal distribution per action dimension with mean 0 and standard deviation 1, $\epsilon \sim \mathcal{N}(0, I)$. The sampled action per update becomes:

$$\hat{\mathbf{a}} = \boldsymbol{\mu}_\phi + (\sigma_\phi^2 \odot \mathcal{N}(0, I)), \quad (\text{S8})$$

where \odot indicates element-wise multiplication. Therefore, the loss expression ends up as:

$$L(\phi) = \mathbb{E}_{s \sim \tau, \hat{\mathbf{a}} \sim f(\phi, \epsilon)} [\beta \log(\pi_\phi(\hat{\mathbf{a}}|s)) - Q_\theta(s, \hat{\mathbf{a}})]. \quad (\text{S9})$$

Each loss function is used as the starting point for the backpropagation algorithm for critic and actor respectively. Once the gradients are obtained from said algorithm, they are used to update the weights via Gradient Descent optimisation method.

Finally, to fully understand SAC it is important to explain a few tricks used for stabilisation purposes. First, the tuples (s, \mathbf{a}, r, s') obtained while interacting with the environment, are saved in a memory buffer \mathcal{D} . Then, at update time, batches of this memory are used to update the weights from the actor and critic, hence now the trajectory in Eqs. (S6) and (S9) would be changed as sampling from \mathcal{D} . This is carried out because consecutive tuples are highly correlated and may lead to instability.

Second, as seen in Eq. (S6), the soft Bellman equation used to compare with the Q-function predictions depends on the same weights that are being updated θ which is another source of instability. To solve this, a copy of the original critic network with weights θ' is introduced with the name of target network and is used only on the critic term present in the Soft Bellman equation. At each update, weights θ will change and θ' will adopt a new value via Polyak averaging:

$$\theta' = \rho\theta + (1 - \rho)\theta', \quad (\text{S10})$$

where ρ is a hyperparameter to be tuned with usually low values to limit abrupt changes in the target network.

Third, on training time, the value of β is automatically tuned via optimisation where the entropy is constrained to a certain value, \bar{H} , in expectation (in [2] it is simply $\bar{H} = -\text{Dim}(\mathbb{A})$). This can be expressed as a loss in the following equation:

$$L(\beta) = \mathbb{E}_{s \sim \tau, \mathbf{a} \sim \pi_\phi(\cdot|s)} [-\beta \log \pi_\phi(\mathbf{a}|s) - \bar{H}]. \quad (\text{S11})$$

With this update, the additional hyperparameter β is eliminated.

Fourth, critic estimations are known to be overoptimistic [3] which may lead to wrong actions. The established solution to this problem is to use two Q-networks and select the minimum over both predictions. Consequently, SAC has a total of four Q-networks if we take into account that both online and target networks need two networks each to avoid the overestimation of the critic.

Finally, in SAC, the policy, π_ϕ , outputs two vectors: a mean vector, $\boldsymbol{\mu}_\phi$, and a standard deviation vector σ_ϕ dependant on parameters ϕ . This two values are the parameters of a Normal distribution $\mathcal{N}(\boldsymbol{\mu}_\phi, \sigma_\phi^2)$. While training, SAC samples from the Normal distribution to obtain the action i.e. $\mathbf{a}_t \sim \mathcal{N}(\boldsymbol{\mu}_\phi, \sigma_\phi^2)$. This is done for exploratory purposes in order to discover actions that may lead to more reward in the future. However, in testing time, once the policy is trained, the action is simply the mean of the Normal distribution, $\mathbf{a}_t = \boldsymbol{\mu}_\phi$, as ineffective exploratory actions should be avoided to achieve maximum performance.

3. PSEUDOCODE OF THE TRAINING LOOP

Pseudocode S1 shows the training process of a MARL controller in the AO context. Due to time concerns, in our implementation, we have parallelised the update but not the action selection. In practice, the only part that may increase the delay is the action selection, which must be parallelised to the finest degree.

4. ADDITIONAL RL HYPERPARAMETERS

In Table S1, we show the hyperparameters used for SAC in all the AO experiments. Most of the hyperparameter values are default ones (as reported in [2]) as they usually work well for any problem. We experimented with the value of γ and found out that lower values of $\gamma \in [0, 0.5]$ lead to better performance. We recall here that γ is an hyperparameter that weighs future rewards and, in the AO case, the commands only have a short-term effect, therefore, a small value of γ is required for high performance.

Algorithm S1. MARL-AO training loop

```

1: Given a trained autoencoder,  $N$  policies ( $\pi^1, \dots, \pi^N$ ),  $N$  empty memories ( $\mathcal{D}^1, \dots, \mathcal{D}^N$ ) and a
   vector of  $M^i$  modes to be controlled per agent  $i$ .
2: for episode=1 to end do
3:   simulation.reset() # Resets simulation and DM shape
4:   for  $t=0$  to  $T$  do
5:     raytrace() # Raytracing operation for the target image and WFS
6:     wfs.compute_image() # Computes WFS image based on last raytracing operation
7:     autoencoder.denoise_wfs_image()
8:     wfs.do_centroids()
9:      $\hat{c}_t = \text{integrator.do\_control}()$  # Computes current residual commands
10:    for agent  $i=1$  to  $N$  in parallel do
11:       $s_t^i = \text{create\_state}(i)$  # Creates state for agent  $i$  based on  $\hat{c}_t$  and previous  $\hat{C}$ 
12:       $a_t^i \sim \pi^i(\cdot | s_t^i)$ 
13:    Concatenate the actions into a joint action:  $a_t = (a_t^1, \dots, a_t^N)$ 
14:     $C_t = B_{tt} \cdot (\hat{C}_{t-1} - g \cdot \hat{c}_t + a_t)$ 
15:    applyCommand( $C_t$ ) # Applies  $C$  simulating the delay in the DM
16:    calculate_psf_image()
17:    for agent  $i$  to  $N$  in parallel do
18:       $r_t^i = -\frac{\kappa}{|M^i|} \sum_{m \in M^i} (\hat{c}_t^m)^2$ 
19:    if  $t \geq d$  then
20:      for agent  $i=1$  to  $N$  in parallel do
21:        Add tuple ( $s^i = s_{t-d}^i, a^i = a_{t-d}^i, s'^i = s_t^i, r^i = r_t^i$ ) to memory  $\mathcal{D}^i$ 
22:        Train SAC with mini-batches of tuples as reported in Section 2 and [2].

```

γ	0.1	Memory size	10^6
Num. hidden size	256	Batch size	256
Num. hidden layers	2	Gradient descent optimiser	ADAM
ρ	0.005	β	Automatic
Learning rate	0.0003	Neural Network Layers	Fully connected

Table S1. Soft Actor Critic hyperparameters.

REFERENCES

1. F. Ferreira, E. Gendron, G. Rousset, and D. Gratadour, "Numerical estimation of wavefront error breakdown in adaptive optics," *Astron. & Astrophys.* **616**, A102 (2018).
2. T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine, "Soft actor-critic algorithms and applications," *arXiv preprint arXiv:1812.05905* (2018).
3. H. Van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 30 (2016), pp. 2094–2100.