

# AN OPTICS CONTROLLING ENVIRONMENT AND REINFORCEMENT LEARNING BENCHMARKS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Deep reinforcement learning has the potential to address various scientific problems. In this paper, we implement an optics simulation environment for reinforcement learning based controllers. The environment incorporates nonconvex and nonlinear optical phenomena as well as more realistic time-dependent noise. Then we provide the benchmark results of several state-of-the-art reinforcement learning algorithms on the proposed simulation environment. In the end, we discuss the difficulty of controlling the real-world optical environment with reinforcement learning algorithms. We will make the code of the paper publicly available.

## 1 INTRODUCTION

In recent years, deep reinforcement learning (RL) has been used to solve challenging problems in various fields (Sutton & Barto, 2018), including self-driving car (Bansal et al., 2018) and robot control (Zhang et al., 2015). Among all of the applications, deep RL made significant progress in play games on a superhuman level (Mnih et al., 2013; Silver et al., 2014; 2016; Vinyals et al., 2017). Beyond playing games, deep RL has the potential to strongly impact the traditional control and automation tasks in the natural science, such as control problems in chemistry (Dressler et al., 2018), biology (Izawa et al., 2004), quantum physics (Bukov et al., 2018), optics and photonics (Genty et al., 2020).

In optics and photonics, there is particular potential for RL methods to drive the next generation of optical laser technologies (Genty et al., 2020). This is not only because there is increasing demand for adaptive control and automation (of tuning and control) for optical systems (Baumeister et al., 2018), but also because many phenomena in optics are nonlinear and multidimensional (Shen, 1984), with noise-sensitive dynamics that are extremely challenging to model using conventional methods. RL methods are able to control multidimensional environment with nonlinear function approximation (Dai et al., 2018). Thus, study the RL controller in optics becomes increasingly promising in optics and photonics as well as its applications in scientific research, medicine, and other industries (Genty et al., 2020; Fermann & Hartl, 2013).

Traditionally, many of the control problems in optics and photonics were implemented by stochastic parallel gradient descent (SPGD) algorithm with PID controller (Cauwenberghs, 1993; Zhou et al., 2009; Abuduweili et al., 2020a). The target is to maximize the reward (e.g. optical pulse energy) by adjusting and controlling the system parameters. The SPGD algorithm is one of the special cases of stochastic error descent method (Cauwenberghs, 1993; Dembo & Kailath, 1990). Stochastic error descent is based on the model-free distributed learning mechanism. A parameter update rule is proposed by which each individual parameter vector perturbation contributes to a decrease in error (or increase in reward). However, SPGD is typically a convex optimization solver, and many control problems in optics are non-convex. SPGD may be failed to search the global optimum of the optics control system unless the initial state of the system is near a global optimum. In general, the initial state of the optical system was adjusted by experienced experts, then utilizing SPGD to control the manually adjusted system, which becomes extremely hard with the increasing system complexity. In order to achieve efficient control and automation, deep RL was introduced to control optical systems (Tünnermann & Shirakawa, 2019; Sun et al., 2020; Abuduweili et al., 2020b). Most of the previous works implemented Deep-Q Network (Mnih et al., 2013) and Deep Deterministic Policy Gradient (Lillicrap et al., 2015), in optical control systems to achieve the comparable performance with traditional SPGD-PID control (Tünnermann & Shirakawa, 2019; Valensise et al., 2021). But

there is a lack of works on the evaluation of more RL algorithms in the more complex optical control environment.

Studying and validating RL algorithms in the real-world optical system is a challenging process because its cost is expensive and requires experienced experts to implement the optical system. Instrumenting and operating RL algorithms in a simple optical system require significant funds and manpower. An effective alternative to validate RL algorithms in optics is simulation. Simulation has been used for robotics and autonomous driving since the early days of research (Pomerleau, 1998; Bellemare et al., 2013). As learning-based robotics expands in both interest and application, the role of simulation may become ever more critical in driving research progress. But there is not any open sourced RL environment for optics control simulation by now.

In this paper, we introduce OPS (**O**ptical **P**ulse **S**tacking environment) - a scalable open simulator for controlling a typical optical system. The physics behind our OPS system is the same as many other optical problems, including coherent optical inference (Wetzstein et al., 2020) and linear optical sampling (Dorrer et al., 2003), which can be used for precise measurement, industrial manufacturing, and scientific research. A typical optical pulse stacking system directly and symmetrically stacks up the input pulses to multiply the pulse energy for output stacked pulses (Tünnermann & Shirakawa, 2017; Stark et al., 2017; Astrauskas et al., 2017; Yang et al., 2020). By providing an optical control simulation environment, we aim to encourage exploration of the application of RL on optical control tasks and furtherly explore the RL controllers in natural science. We use OPS to evaluate some important RL algorithms including twin delayed deep deterministic policy gradient (TD3, Fujimoto et al. (2018)), soft actor-critic (SAC, Haarnoja et al. (2018a)), and proximal policy optimization (PPO, Schulman et al. (2017)). After reporting the results of these RL algorithms, we discuss the difficulty of RL algorithms in the real-world optical system. With the provided simulating environment OPS and the experiments of RL algorithms, we believe that this work can promote the research on RL applications in optics as well as benefit both the machine learning and the optics community.

## 2 SIMULATION ENVIRONMENT

### 2.1 PHYSICS OF THE SIMULATION

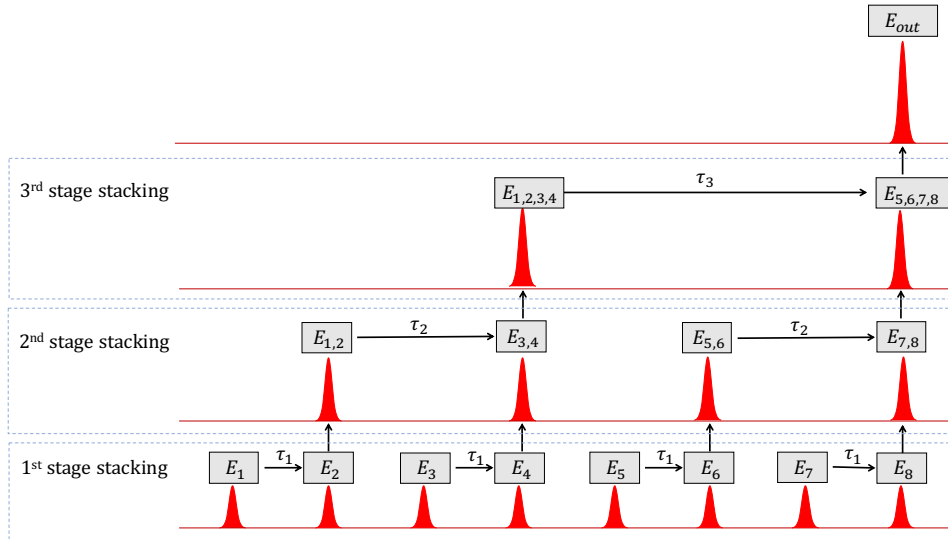


Figure 1: Illustration of the principle of optical pulse stacking. Only 3-stage pulse stacking was plotted for simplicity.

The optical pulse stacking (OPS, or also called pulse combination) system recursively stacks up the optical pulses in the time domain. The dynamics of the OPS are similar to the recurrent neural networks (RNN) or Wavenet architecture (Oord et al., 2016). We illustrate the dynamics of the OPS

in RNN style as shown in fig. 1. The input of the OPS is a periodic pulse train<sup>1</sup> with a repetition period of  $T$ . Assume the basic function of the first pulse at time step  $t$  is  $E_1 = E(t)$ , then the consecutive pulses can be described as  $E_2 = E(t + T)$ ,  $E_3 = E(t + 2T)$ ... The OPS system recursively imposes the time delay on earlier pulses for every two consecutive pulse pairs. As an example, the 1st stage time-delay controller imposes the time delay  $\tau_1$  on pulse 1 to shift the pulse 1. With the proper time delay, pulse 1 could be stacked with the next pulse  $E_2$  to create the stacked pulses  $E_{1,2} = E(t + \tau_1) + E(t + T)$ . Similarly, pulse 3 could be stacked with pulse 4 to create  $E_{3,4} = E(t + 2T + \tau_1) + E(t + 3T)$ , and so on. In 2nd stage OPS, the time delay  $\tau_2$  was further imposed to  $E_{1,2}$  to make it to stack up with  $E_{3,4}$  to create  $E_{1,2,3,4}$ . This kind of stacking is repeated in each stage OPS controller, which stacks up the pulses in geometrical progression (recursion). An  $N$ -stage OPS system simply multiplies pulse energy by  $2^N$  times by stacking up  $2^N$  pulses, in which  $N$  time delays ( $\tau_1, \tau_2, \dots, \tau_N$ ) are needed to control and stabilize. Please check the more detailed illustration and configuration of the real-world OPS experiment in appendix A.

## 2.2 CONTROL OBJECTIVE AND NOISE

The objective of the controlling OPS system is to maximize the final stacked (output) pulse energy by adjusting the time delays. For  $N$ -stage OPS system, let  $P_N$  denotes the final energy of  $N$  times stacked pulse, and  $\tau = [\tau_1, \tau_2, \dots, \tau_N]$  denote the time delays. Then the objective function for controlling  $N$ -stage OPS system is:

$$\arg \max_{\tau} P_N(\tau) = \arg \max_{\tau_1, \tau_2, \dots, \tau_N} P_N(\tau_1, \tau_2, \dots, \tau_N) \quad (1)$$

If any noise were ignored, we would analyze the exact function of the final pulse energy  $P_N$  w.r.t. the time delays  $\tau$ . Figure 2(a) shows the function of the pulse energy  $P_1(\tau_1)$  w.r.t. the first time delay  $\tau_1$  in 1-stage OPS system. And fig. 2(b) shows the function surface of  $P_2(\tau_1, \tau_2)$  w.r.t. the first and second stages time delay ( $\tau_1, \tau_2$ ) in 2-stage OPS system. As can be seen, the control function of the OPS system is non-linear and non-convex even ignoring any noises<sup>2</sup>. This is a challenging problem for any controlling algorithms to achieve the global optimum (or better local optimum) from a random initial state.

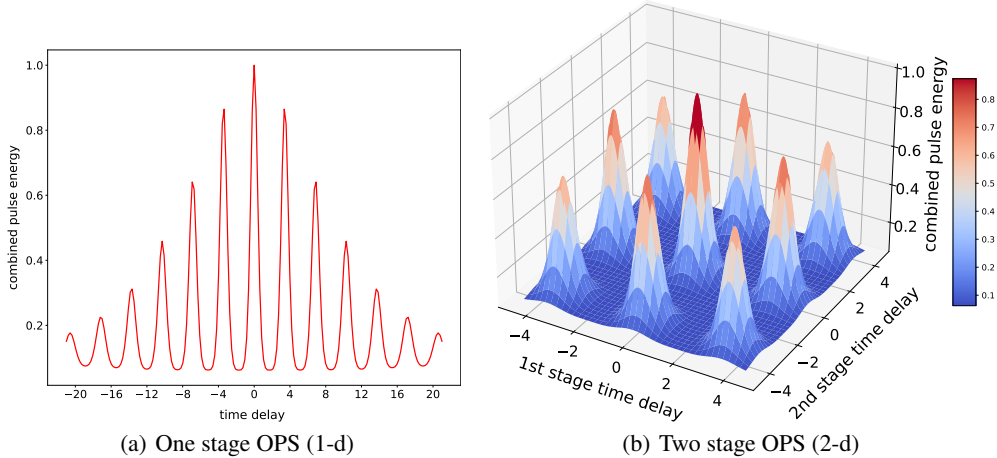


Figure 2: Function plot of the (a) 1-stage OPS: pulse energy  $P_1(\tau_1)$  w.r.t. delay line  $\tau_1$ . (b) 2-stage OPS: pulse energy  $P_2(\tau_1, \tau_2)$  w.r.t. delay lines  $(\tau_1, \tau_2)$ .

In general, noise can not be ignored and the system is quite noise-sensitive. That is because the wavelength of the pulse is in  $\mu\text{m}$  level ( $1\mu\text{m} = 10^{-6}\text{m}$ ). The noise in the environment, including vibration of optical devices and temperature drift of the atmosphere, could easily bring the shift of the time delay then change the output pulses. So the objective function in real-world practice is more complex than fig. 2, especially for higher stage (high-dimensional) OPS.

<sup>1</sup>The periodic pulse train generally emitted by lasers. The wave function of each laser pulse is almost the same except for the time delay of a period.

<sup>2</sup>Part of the reasons are the optical periodicity and nonlinearity of the coherent interference.

In this simulation, we mainly consider two kinds of noise. The first one is fast noise which comes from the vibration of devices. The noise could be formulated as a zero-mean Gaussian random noise  $\mathcal{N}(0, \sigma)$  by following the simulation noise of Tünnermann & Shirakawa (2019). The second is slow noise, which comes from slow temperature drift. The influence of the temperature drift can be formulated as a piecewise linear function (Ivanova et al., 2021). We used  $\mu_t$  to represent the time-dependent slow noise and formulated  $\mu_t$  as a slow-changed piecewise linear function.

By incorporating these two kinds of noise, overall noise  $e_t$  at time step  $t$  follow the distribution of:

$$e_t \sim \mathcal{N}(\mu_t, \sigma). \quad (2)$$

Please note that  $\mu_t$  changes very slowly with time. For episodic training for RL agents, the  $\mu_t$  could be considered as a constant value for iterations within an episode. But the value of  $\mu_t$  might differ from one episode to the next episode. It makes the distribution of the noise change with episodic training. And more importantly, it causes different noise distribution for the testing environment compared with the training environment. This kind of noise distribution shift between training and the testing environment was rarely discussed in the previous simulations. But this is quite common for noise-sensitive optical control tasks.

### 2.3 REINFORCEMENT LEARNING ENVIRONMENT

**Interactions with RL agent.** An RL agent interacts with the OPS environment in discrete time steps, as shown in fig. 3. At each time step  $t$ , the RL agent receives the current state of the OPS environment  $s_t$ . Then the RL agent chooses an action  $a_t$  to send the OPS environment. The environment conduct the action and moves to new state  $s_{t+1}$ . Then the reward  $r_t$ , which measured by the state  $s_{t+1}$ , feedback to the RL agent. The RL agent trained with the experience  $(s_t, a_t, s_{t+1}, r_t)$  to learn a policy  $\pi(a, s)$  which maximizes the expected cumulative reward.

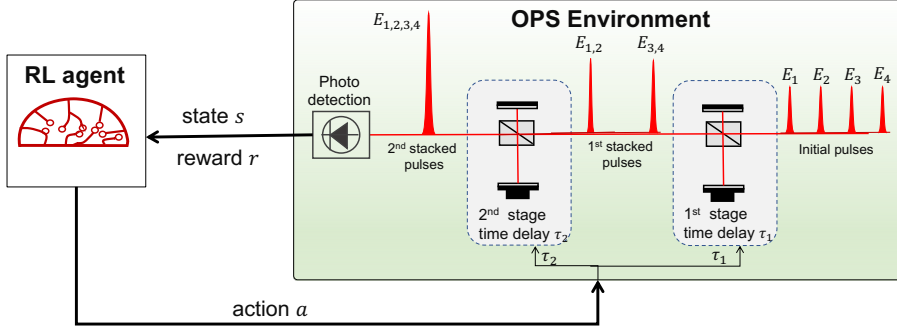


Figure 3: Illustration of the interaction between RL agent and OPS environment. Only 2-stage pulse stacking was plotted in OPS for simplicity.

**State space.** State space of OPS is a continual and multidimensional vector space. The state value  $s_t$  could be described as the pulse intensity measurement of the final stacked pulse  $s_t = \text{Intensity}(E_{\text{out}}(t))$ . So  $s_t$  is the time-domain "picture" of the final stacked pulse, which directly reflects the performance of the control. In a real-world system, pulse intensity was detected by a photo-detector then converted to digital time-series signals. In our simulation, we also implement real-time rendering of the pulse intensity to monitor the controlling process.

**Action space.** Action space of  $N$ -stage OPS environment is a continual and  $N$ -dimensional vector space. At time step  $t$ , the action  $a_t$  is a time delay value  $\tau(t)$  for  $N$ -stage OPS environment:  $a_t = \tau(t) = (\tau_1(t), \tau_2(t), \dots, \tau_N(t))$ . The time delay value  $\tau(t)$  was conducted by OPS environment to lead the next state.

**Reward.** As mentioned in section 2.2, the objective of the OPS controller is maximizing the final stacked pulse energy  $P_N(\tau)$ . We used the reward value as normalized final pulse energy:

$$r = -\frac{(P_N(\tau) - P_{max})^2}{(P_{min} - P_{max})^2}, \quad (3)$$

where  $P_{max}$  is the maximum pulse energy achieved at the global optimum, and  $P_{min}$  is the minimum pulse energy. The maximum reward 0 achieved when  $P(\tau) = P_{max}$  (peak position of

fig. 2(b)) . As the model moves better and better, the cumulative reward will be closer to zero instead of growing all the time.

**State transition function.** The environmental noise poses direct impacts to the delay lines (including the vibration and temperature shift noise of the delay line devices). So in the state transition, real conducted delay line  $\tau_{\text{real}}(t)$  is a combination of the action  $a_t = \tau(t)$  and noise  $e_t$ :

$$\tau_{\text{real}}(t) = \tau(t) + e_t = a_t + e_t, e_t \sim \mathcal{N}(\mu_t, \sigma). \quad (4)$$

Then the real time delay  $\tau_{\text{real}}(t)$  imposed to some selected pulses<sup>3</sup> by delay line devices (the device impose additional time delay for pulses) as conduct the action. The state transition is governed by state, action and noise. The exact form of the state transition follow the principle of the coherent pulse interference.

**Different game difficulty of the environment.** We implemented the OPS environment for any ( $N \in \{1, 2, 3, \dots\}$ ) stage of pulse stacking. With the increase of the number of stages, the control would become more and more difficult. In addition to the customized number of stages, we also provided three modes (easy, medium, and hard) for each stage OPS, as shown in fig. 4. The mode was determined by the initial state of the system and noise distribution:

- Easy mode. The initial state of the OPS system is near the global optimum for easy mode. Figure 6(a) shows the example initial state of the easy mode of the 3-stage OPS environment. This is a case for many traditional optics control problems: the initial state of the system is tuned by "experts" to make it easy to control for convex controllers.
- Medium mode. The initial state of the system is random, as shown in fig. 6(b), which makes the control problem becomes nonconvex. But in medium mode, the noise is time-independent and we simply set the noise distribution as  $e_t \sim \mathcal{N}(0, \sigma)$ . This is the case for many classical reinforcement learning settings. The noise distribution of each episode is the same. More importantly, the noise behavior of the training environment and testing environment is similar.
- Hard mode. Similar to medium mode and fig. 6(b), the initial state of the system is random. Different from the medium mode, the noise behavior is more complicated. The mean value of the noise distribution  $\mu_t$  is a time-dependent variable. Which slowly changes during time,  $e_t \sim \mathcal{N}(\mu_t, \sigma)$ . Because in real-world applications, we always deploy the testing environment and algorithms after training, so the noise distribution of the testing environment is different from the training environment. The hard mode is closer to real-world settings.

Mode	Initial state	Noise
easy	near the optimum	time independent; $\mu_t \equiv 0$
medium	random	time independent; $\mu_t \equiv 0$
hard	random	time dependent; $d\mu_t/dt \neq 0$

Figure 4: Comparison of the different game modes.

```

from optics_env import OPS_env
env = OPS_env(stage=5, mode="medium")
env.reset()
done = False
while not done:
    action = env.action_space.sample()
    observation, reward, done, info = env.step(action)
    env.render()

```

Figure 5: Example code of the OPS environment.

**API & sample usage.** The optical and physical principle of the simulation is based on the Nonlinear-Optics-Modeling package (Hult, 2007). The OPS environment is out of the box compatible with the widely used OpenAI Gym API (Brockman et al., 2016). We show an example code of running random agent on OPS environment as fig. 5.

**Features of the OPS environment.** We summarize the key features of the OPS environment as follows:

<sup>3</sup>Real conducted time delay is  $\tau_{\text{real}}(t) = (\tau_{1,\text{real}}(t), \tau_{2,\text{real}}(t), \dots)$ . In which, 1st stage delay line  $\tau_{1,\text{real}}(t)$  imposed to pulse-1 ( $E_1$ ), pulse-3 ( $E_3$ ) ... Then the 2nd stage delay line  $\tau_{2,\text{real}}(t)$  imposed to the previously 1st stacked pulse  $E_{1,2}$ .

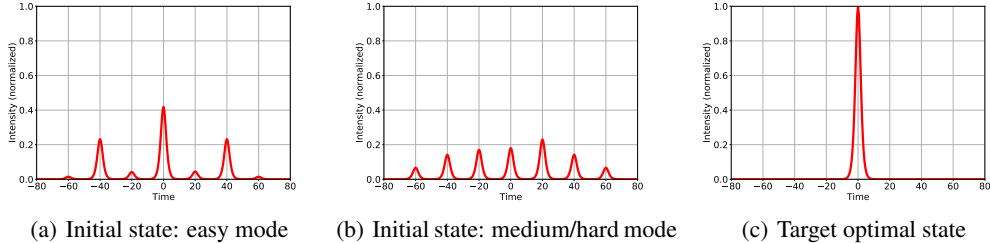


Figure 6: Rendering examples of the state of (a) initial state for easy mode, (b) initial state for medium or hard mode, (c) global optimal target state in a 2-stage OPS environment. The initial state of the easy mode has almost sacked some parts of pulses, which is more closer to the target state. The initial state of medium or hard mode is almost random and might be trapped into local optimum.

- Open source optical control environment. To the best of our knowledge, this is the first open-sourced RL environment for optics control problems. The open-source licenses enable researchers to inspect the underlying code and to modify environments if required to test new research ideas.
- Scalable and difficulty-adjustable scientific environment. Many of the recent RL environments (e.g. Atari) are easy to solve. In our OPS environment, the difficulty of the environment is flexible. And the dimension of the action space is easy to scalable with stage number  $N$ . For the  $N$ -stage OPS environment, if we choose quite larger  $N$  with hard mode, controlling the environment could become quite hard. If the hard scientific control problem could be solved effectively, which would have a broader impact on many scientific control problems.
- Realistic noise. In the hard mode of the OPS environment, we modulate the noise distribution as the time-dependent function. It made the noise distribution of the testing environment is different from the noise distribution of the training environment. This is more realistic for noise-sensitive systems. It also increases the stochasticity of the environment.

### 3 EXPERIMENTS

#### 3.1 EXPERIMENTAL SETUP

As a reference, we provide benchmark results for four state-of-the-art reinforcement learning algorithms: PPO (Schulman et al., 2017), TD3 (Fujimoto et al., 2018), and SAC (Haarnoja et al., 2018b). We implement the algorithms using stable-baseline-3 (Raffin et al., 2019). The training procedure for an RL agent is divided into a couple of episodes, each episode lasts for 200 steps. Other hyperparameters or each algorithm and training setting can be found in Appendix B.1. For each of the experimental settings, we run ten random seeds and average the results.

#### 3.2 RESULTS ON CONTROLLING 5-STAGE OPS ENVIRONMENT

In this section, we mainly report the results for the 5-stage OPS system, that stacked  $2^5 = 32$  pulses. For the results of the different stage OPS system, please check appendix B.2. In a 5-stage OPS system, we evaluate all of the four algorithms in 3 difficulty modes of the environment: easy, medium, and hard.

Training curve (plot for training reward per step w.r.t. iterations) of PPO, TD3, and SAC algorithms has been shown in fig. 7(a) for easy mode, fig. 7(b) for medium mode, and fig. 7(c) for hard mode. As can be seen, the performance of TD3 and SAC is similar and higher than PPO for all three modes. For the difficulty mode of the environment, the convergence speed is slow down and the final convergent value decreases with the increase of difficulty of the environment. As an example, in easy mode, SAC converges to the reward value of  $-0.04$  within 100,000 steps, but it takes 200,000 to converges to the reward value of  $-0.1$  for hard mode.

After training the RL agents, we evaluated the performance in the testing environment. The final return (stacked pulse power  $P_N$ ) under different iterations on easy mode, medium mode, and hard mode as shown in fig. 8(a), fig. 8(b), and fig. 8(c). As can be seen, although the training curve of the

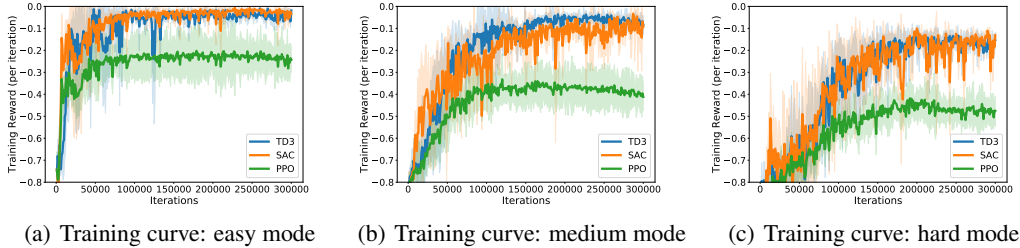


Figure 7: Training curve for SAC, TD3, and PPO on 5-stage OPS environment for (a) easy mode, (b) medium mode, and (c) hard mode. The dashed region shows the area within the standard deviation.

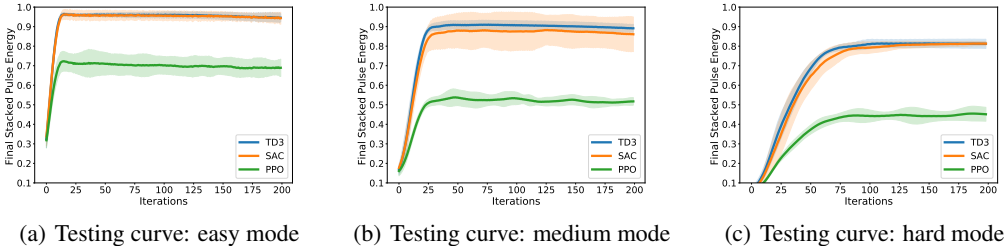


Figure 8: Evaluation of the stacked pulse power  $P_N$  (normalized) of testing environment for (a) easy mode, (b) medium mode, and (c) hard mode.

medium mode (fig. 7(b)) and hard mode fig. 7(c) is a bit of similar, the evaluation curve on testing environment of medium mode (fig. 8(b)) and hard mode fig. 8(c) is different. That is because the hard mode has a different noise distribution for the training and test environment. That makes the evaluation control on the testing environment for hard mode is slow to converge and achieved the lower final return. Please see the detailed results of the evaluation in appendix B.3, and check the animation video in supplementary video 1.

We reported the final return (combined pulse power  $P_N$ ) of the training and testing environment on the trained policy as table 1. Please note that the training environment and testing environment for easy and medium mode is similar, just like the classical Atari environment. The performance differences are mainly caused by randomness. We showed that the performance difference between the training and testing environment is much higher for hard mode. That is because of the different noise behavior of the training and testing environment, which makes the control complicated.

Mode	Evaluation on which environment	PPO	TD3	SAC
easy	training	$0.7684 \pm 0.0884$	$0.9580 \pm 0.0189$	$0.9637 \pm 0.0172$
	testing	$0.7439 \pm 0.0463$	$0.9541 \pm 0.0177$	$0.9514 \pm 0.0231$
medium	training	$0.6210 \pm 0.0828$	$0.9204 \pm 0.0351$	$0.8945 \pm 0.0501$
	testing	$0.6182 \pm 0.0229$	$0.9106 \pm 0.0217$	$0.8833 \pm 0.0838$
hard	training	$0.5473 \pm 0.0680$	$0.8524 \pm 0.0380$	$0.8515 \pm 0.0375$
	testing	$0.4461 \pm 0.0300$	$0.8130 \pm 0.0215$	$0.8071 \pm 0.0164$

Table 1: Evaluation performance of PPO, TD3, and SAC on three (easy, medium, hard) modes. Final return  $P_N$  on both the training environment and testing environment was evaluated.

### 3.3 RESULTS ON DIFFERENT STAGE EXPERIMENTS

We evaluated all of the four algorithms of the different  $N$ -stage OPS environments with hard modes. Figure 9(a) shows the training curve, and fig. 9(b) shows the testing curve of TD3 on different  $N$ -stage OPS system. As can be seen, with the increase of stage number, the training convergence became slower, and the final return  $P_N$  became smaller.



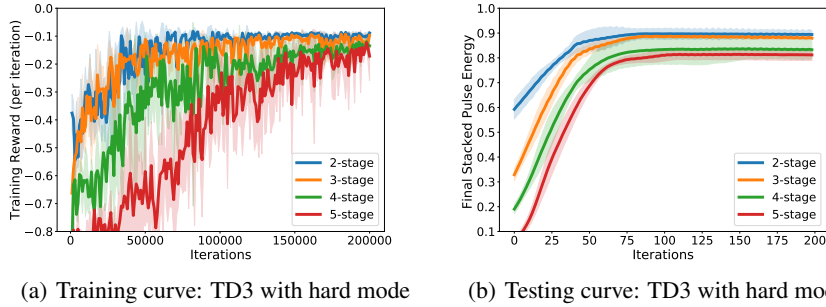


Figure 9: Comparison of the results on hard mode  $N$ -stage OPS environment with TD3 algorithms. (a) shows the training curve: (b) shows the evaluation of final return  $P_N$  of the testing environment.

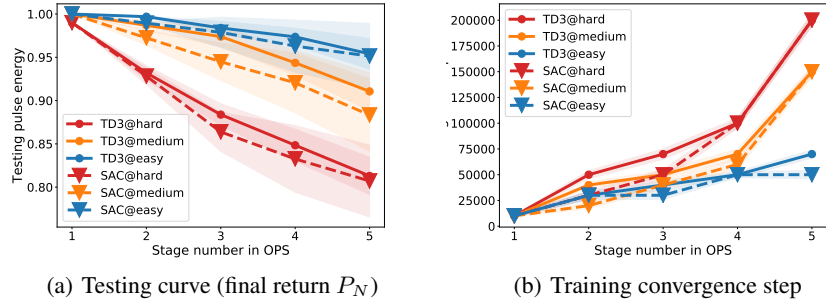


Figure 10: (a) Final return  $P_N$  of different stage OPS on testing environment controlled with TD3 or SAC. (b) Convergence steps for the training of TD3 and SAC on different stage OPS environments.

We also evaluated the trained TD3 and SAC on the different  $N$  stage testing environment, as shown in fig. 10(a) and fig. 10(b). Figure 10(a) illustrated the final return  $P_N$  under different stages OPS and different difficulty modes. For 1-stage OPS, the final return  $P_N$  could reach 1. That means TD3 and SAC are able to search the global optimum for 1-stage OPS. But for 5-stage OPS with hard mode, the SAC and TD3 only could achieve the 0.8 final return. That means the controlling algorithms were trapped into a local optimum. In the real-world experiments, this case means the 20% energy loss. Figure 10(a) illustrated the training-convergence step for different stages OPS. As the stage number increase, the number of steps to training convergence increases significantly, which will slow down the training for higher stage OPS.

## 4 DISCUSSION

### 4.1 CONNECTIONS WITH REAL-WORLD EXPERIMENT

As far as we know, the previous real-world experiments of RL algorithms on the complicated OPS system are not very successful. DQN or DDPG takes 4-hours to train a simple 1-stage OPS system, and which takes 1-2 days to train a 4-stage OPS system (Tünnemann & Shirakawa, 2019). One of the reasons for the slow training is because by deploying the RL algorithm in an optics system, we need to convert optical signal to analog signal using photo-detector, then convert the analog signal to digital signal using an analog-to-digital converter. These two conversions not only cost some additional time to process the signal but also introduce some noises. For example, as shown in fig. 8(c), TD3 algorithms almost take 300,000 steps to converge on 5-stage OPS environment. If the real-world process time of the OPS system is 0.1s per interaction, the TD3 algorithm would take more than 8 hours to train. Actually, the noise of the environment is much more complicated than the simulation environment, though the simulation considers the time-dependent property of the temperature drift. The longer training time made the temperature drift and noise more severe. It also causes some additional energy costs in the real-world optical system. So the fast training and noise-robust RL algorithms are critical to controlling tasks in optics, which is our main concern about implementing OPS simulation environments. Please see the discussion about the real-world and simulated experiments in appendix C.



## 4.2 TRANSFER TRAINED POLICY

It is possible to transfer the trained policy between different simulating environments, even transferring from simulating environment to a real-world experiment. The major difference between the simulation and real-world environments is the different noise levels. We conduct an experiment to show the transferability between different noise levels.

In our simulation environment, the noise level is dependent on the difficulty mode. Then we explore the transferability of the trained policy between "easy", "medium", and "hard" modes. We trained the policy on "hard" mode environment, then tested the trained policy on "hard", "medium", and "easy" environment, as shown in fig. 11 (a). Transfer results of "medium" trained policy and "easy" trained policy are shown in fig. 11 (b) and fig. 11 (c) respectively. As can be seen, We can perfectly transfer the harder mode trained policy to easier mode environments. When the easier mode training strategy is transferring to a harder mode environment, the performance may drop, and there is a jitter in pulse energy. The possible reason is that, when transferring to the harder environment, the trained policy could be trapped in the "bad" local minimum closer to the initial point and easily affected by noise.

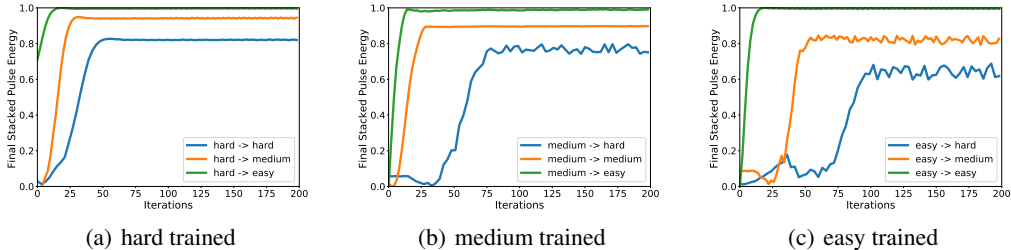


Figure 11: Demonstration of the transfer performance of the trained policy on (a) hard mode training environment; (b) medium mode training environment; (c) easy mode training environment.

Figure 11 shows that it is possible to transfer the trained policy between different noise levels, but it is more useful to train the policy in a harder environment than tested on the easier environment. Thus, we could explore fast and robust controlling algorithms in more harder simulation environment (with introducing more noise and more uncertainty in the simulation) then deploy the trained policy to real-world physical systems.

## 4.3 FUTURE WORK

One important work is to explore fast and robust RL algorithms in the simulation environment then deploy them to the real-world control system. Another important topic is that explore the additional information about the function of the OPS (fig. 2) and incorporating it with RL controllers. For controlling the optical system, we are not interested in "generic" nonconvex problems, but rather, optical systems (or other scientific control problems) typically provide us with much richer structural information. So it is promising to incorporate additional structural information behind the function of the optical control into the RL controllers.

## 5 CONCLUSION

In this paper, we introduce OPS – an open-sourced simulator for controlling the pulse stacking system. To the best of our knowledge, this is the first open-sourced RL environment for optics control problems. Then we evaluated the SAC, TD3, and PPO on our proposed simulation environment. By providing an optical control simulation environment and RL benchmarks, we aim to encourage exploration of the application of RL on optical control tasks and furtherly explore the RL controllers in natural science. We believe that this work can promote the research on RL applications in optics as well as benefit both the machine learning and the optics community.

#### ETHICS STATEMENT

I certify that all co-authors of this work have read and commit to adhering to the ICLR Statement on Ethics, Fairness, Inclusivity, and Code of Conduct.

#### REPRODUCIBILITY STATEMENT

The source code of the paper is attached to the supplementary materials. I certify that all of the experiments in the paper can be reproduced with the provided source code.

#### REFERENCES

- Abulikemu Abuduweili, Bowei Yang, and Zhigang Zhang. Modified stochastic gradient algorithms for controlling coherent pulse stacking. In *Conference on Lasers and Electro-Optics*, pp. STh4P.1, 2020a.
- Abulikemu Abuduweili, Bowei Yang, and Zhigang Zhang. Control of delay lines with reinforcement learning for coherent pulse stacking. In *Conference on Lasers and Electro-Optics*, pp. JW2F.33, 2020b.
- Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76, 2021.
- Ignas Astrauskas, Edgar Kaksis, Tobias Flöry, Giedrius Andriukaitis, Audrius Pugžlys, Andrius Baltuška, John Ruppe, Siyun Chen, Almantas Galvanauskas, and Tadas Balčiūnas. High-energy pulse stacking via regenerative pulse-burst amplification. *Optics letters*, 42(11):2201–2204, 2017.
- Mayank Bansal, Alex Krizhevsky, and Abhijit Ogale. Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst. *arXiv preprint arXiv:1812.03079*, 2018.
- Thomas Baumeister, Steven L Brunton, and J Nathan Kutz. Deep learning and model predictive control for self-tuning mode-locked lasers. *JOSA B*, 35(3):617–626, 2018.
- Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47: 253–279, 2013.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- Marin Bukov, Alexandre GR Day, Dries Sels, Phillip Weinberg, Anatoli Polkovnikov, and Pankaj Mehta. Reinforcement learning in different phases of quantum control. *Physical Review X*, 8(3): 031086, 2018.
- Gert Cauwenberghs. A fast stochastic error-descent algorithm for supervised learning and optimization. *Advances in neural information processing systems*, 5:244–251, 1993.
- Yudong Chen and Yuejie Chi. Harnessing structures in big data via guaranteed low-rank matrix estimation: Recent theory and fast algorithms via convex and nonconvex optimization. *IEEE Signal Processing Magazine*, 35(4):14–31, 2018. doi: 10.1109/MSP.2018.2821706.
- Yuejie Chi, Yuxin Chen, and M. Yue Lu. Recent advances in nonconvex methods for high-dimensional estimation. *ICASSP tutorial*, 2018.
- Bo Dai, Albert Shaw, Lihong Li, Lin Xiao, Niao He, Zhen Liu, Jianshu Chen, and Le Song. Sbed: Convergent reinforcement learning with nonlinear function approximation. In *International Conference on Machine Learning*, pp. 1125–1134. PMLR, 2018.
- Amir Dembo and Thomas Kailath. Model-free distributed learning. *IEEE Transactions on Neural Networks*, 1(1):58–70, 1990.
- C Dorrer, DC Kilper, HR Stuart, G Raybon, and MG Raymer. Linear optical sampling. *IEEE Photonics Technology Letters*, 15(12):1746–1748, 2003.

- Oliver J Dressler, Philip D Howes, Jaebum Choo, and Andrew J deMello. Reinforcement learning for dynamic microfluidic control. *ACS omega*, 3(8):10084–10091, 2018.
- Qiang Du, Tong Zhou, Lawrence R Doolittle, Gang Huang, Derun Li, and Russell Wilcox. Deterministic stabilization of eight-way 2d diffractive beam combining using pattern recognition. *Optics letters*, 44(18):4554–4557, 2019a.
- Simon S Du, Yuping Luo, Ruosong Wang, and Hanrui Zhang. Provably efficient  $q$ -learning with function approximation via distribution shift error checking oracle. *arXiv preprint arXiv:1906.06321*, 2019b.
- Gabriel Dulac-Arnold, Daniel Mankowitz, and Todd Hester. Challenges of real-world reinforcement learning. *arXiv preprint arXiv:1904.12901*, 2019.
- Martin E Fermann and Ingmar Hartl. Ultrafast fibre lasers. *Nature photonics*, 7(11):868–874, 2013.
- Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pp. 1587–1596. PMLR, 2018.
- Goëry Genty, Lauri Salmela, John M Dudley, Daniel Brunner, Alexey Kokhanovskiy, Sergei Kobtsev, and Sergei K Turitsyn. Machine learning and applications in ultrafast photonics. *Nature Photonics*, pp. 1–11, 2020.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pp. 1861–1870. PMLR, 10–15 Jul 2018a.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018b.
- Johan Hult. A fourth-order runge–kutta in the interaction picture method for simulating supercontinuum generation in optical fibers. *J. Lightwave Technol.*, 25(12):3770–3775, Dec 2007.
- Yoanna M Ivanova, Hannah Pallubinsky, Rick Kramer, and Wouter van Marken Lichtenbelt. The influence of a moderate temperature drift on thermal physiology and perception. *Physiology & Behavior*, 229:113257, 2021.
- Jun Izawa, Toshiyuki Kondo, and Koji Ito. Biological arm motion through reinforcement learning. *Biological cybernetics*, 91(1):10–22, 2004.
- Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Sobhan Miryoosefi, Kianté Brantley, Hal Daumé III, Miroslav Dudík, and Robert Schapire. Reinforcement learning with convex constraints. *arXiv preprint arXiv:1906.09323*, 2019.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- D Pomerleau. An autonomous land vehicle in a neural network. *Advances in Neural Information Processing Systems; Morgan Kaufmann Publishers Inc.: Burlington, MA, USA*, 1998.
- Antonin Raffin, Ashley Hill, Maximilian Ernestus, Adam Gleave, Anssi Kanervisto, and Noah Dornmann. Stable baselines3. <https://github.com/DLR-RM/stable-baselines3>, 2019.

- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Yuen-Ron Shen. The principles of nonlinear optics. *New York*, 1984.
- David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy gradient algorithms. In Eric P. Xing and Tony Jebara (eds.), *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pp. 387–395, Beijing, China, 22–24 Jun 2014. PMLR. URL <http://proceedings.mlr.press/v32/silver14.html>.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- Henning Stark, Michael Müller, Marco Kienel, Arno Klenke, Jens Limpert, and Andreas Tünnermann. Electro-optically controlled divided-pulse amplification. *Optics express*, 25(12):13494–13503, 2017.
- Chang Sun, Eurika Kaiser, Steven L Brunton, and J Nathan Kutz. Deep reinforcement learning for optical systems: A case study of mode-locked lasers. *Machine Learning: Science and Technology*, 1(4):045013, 2020.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Henrik Tünnermann and Akira Shirakawa. Delay line coherent pulse stacking. *Opt. Lett.*, 42(23):4829–4832, Dec 2017.
- Henrik Tünnermann and Akira Shirakawa. Deep reinforcement learning for coherent beam combining applications. *Opt. Express*, 27(17):24223–24230, Aug 2019.
- Carlo M Valensise, Alessandro Giuseppe, Giulio Cerullo, and Dario Polli. Deep reinforcement learning control of white-light continuum generation. *Optica*, 8(2):239–242, 2021.
- Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, et al. Starcraft ii: A new challenge for reinforcement learning. *arXiv preprint arXiv:1708.04782*, 2017.
- Gordon Wetzstein, Aydogan Ozcan, Sylvain Gigan, Shanhui Fan, Dirk Englund, Marin Soljačić, Cornelia Denz, David AB Miller, and Demetri Psaltis. Inference in artificial intelligence with deep optics and photonics. *Nature*, 588(7836):39–47, 2020.
- Bowei Yang, Guanyu Liu, Abuduweili Abulikemu, Yan Wang, Aimin Wang, and Zhigang Zhang. Coherent stacking of 128 pulses from a ghz repetition rate femtosecond yb: fiber laser. In *Conference on Lasers and Electro-Optics*, pp. JW2F.28, 2020.
- Fangyi Zhang, Jürgen Leitner, Michael Milford, Ben Uprocft, and Peter Corke. Towards vision-based deep reinforcement learning for robotic motion control. *arXiv preprint arXiv:1511.03791*, 2015.
- Pu Zhou, Zejin Liu, Xiaolin Wang, Yanxing Ma, Haotong Ma, Xiaojun Xu, and Shaofeng Guo. Coherent beam combining of fiber amplifiers using stochastic parallel gradient descent algorithm and its application. *IEEE Journal of Selected Topics in Quantum Electronics*, 15(2):248–256, 2009.

## A ADDITIONAL INFORMATION OF OPTICAL PULSE STACKING

### A.1 SYSTEM CONFIGURATION

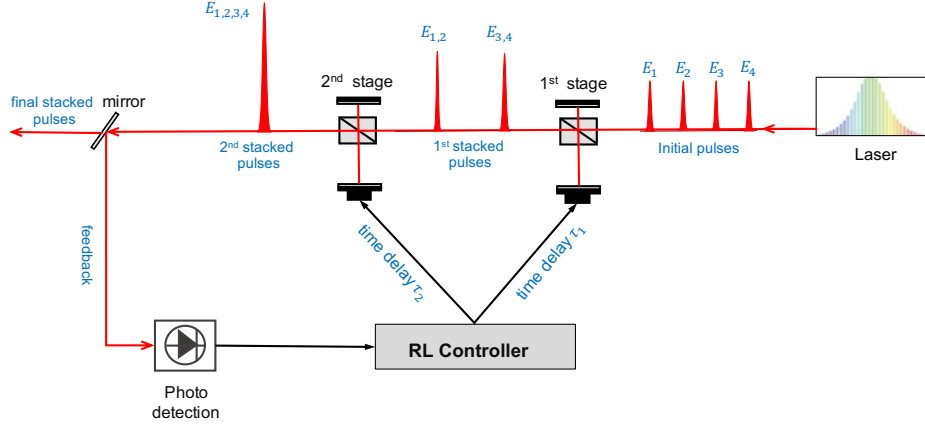


Figure 12: Configuration of optical pulse stacking (OPS) system. Only a 2-stage OPS system was plotted for simplicity.

The configuration of the optical pulse stacking (OPS) system is shown in fig. 12. The source laser delivers a train of periodic optical pulses. Given the base wave function of the laser pulse  $E(t)$  and period  $T$ , each laser pulse in fig. 12 can be described as:

$$E_1 = E(t), E_2 = E(t + T), E_3 = E(t + 2T), E_4 = E(t + 3T). \quad (5)$$

Then the laser pulses were sent to the  $n$ -stage OPS system. In each OPS time delay stacking module, a time delay should be given to former pulses for every two consecutive pulse pairs.

The demonstration of each OPS time delay stacking module is shown in fig. 13. Please check the animation in supplementary video 2. Figure 13(a) shows the initial state of the two pulses before processing by this stage OPS time delay. Figure 13(b) ~ fig. 13(e) show the (chronological) procedures of stacking two pulses by imposing additional time delay. The former pulse  $E_1$  was refracted by the splitter to undergo additional delay lines (vertical path between "mirror" and "time delay controller and mirror" in fig. 13(a)). The latter pulse directly transmitted the splitter. If the displacement of the additional delay line is  $d_1$ , then the additional time delay imposed to  $E_1$  is  $\tau_1 = d_1/c$ , where  $c$  is the light speed. Thus, in experimental implementation, the time delay of the pulse was imposed by additional delay line displacement. The value of the delay line displacement is controlled by the RL controller.

In OPS system, the 1st stage time-delay controller imposes the time delay  $\tau_1$  on pulse  $E_1$  to stack with pulse  $E_2$  to create the stacked pulses  $E_{1,2} = E(t + \tau_1) + E(t + T)$ . Similarly, After imposing time delay, pulse  $E_3$  could be stacked with pulse  $E_4$  to create  $E_{3,4} = E(t + 2T + \tau_1) + E(t + 3T)$ . In the 2nd stage OPS, the time delay  $\tau_2$  was further imposed to  $E_{1,2}$  to make it stacking up with  $E_{3,4}$  to create  $E_{1,2,3,4}$ :

$$E_{1,2,3,4} = E(t + \tau_1 + \tau_2) + E(t + T + \tau_2) + E(t + 2T + \tau_1) + E(t + 3T). \quad (6)$$

If noise was ignored, when  $\tau_1 = T, \tau_2 = 2T$ ,  $E_{1,2,3,4}$  achieved the maximum value  $4E(t + 3T)$ . Furtherly, for  $N$ -stage OPS system, 1st, 2nd, ...,  $N$ -th time delay  $\tau_1, \tau_2, \dots, \tau_N$  matches to  $2^0, 2^1, \dots, 2^{N-1}$  times of the pulse period  $T$ , the  $2^N$  pulses will be perfectly stacked, and the power of the output pulse reaches the global maximum. In practice, noise can not be ignored, so the exact value of time delay  $\tau_1, \tau_2, \dots, \tau_N$  could be adjusted according to the feedback.

### A.2 REAL PHYSICAL SYSTEM

The real physical OPS system is shown in fig. 14. RL algorithm computes the value of each time delay  $\tau_1, \tau_2, \dots, \tau_N$  and sends the values to each stage delay line controller. (RL controller connected with the electric signal line of 1st, 2nd, 3rd delay line located at the bottom of the fig. 14.) Real-world OPS control experiments are quite costly and slow.

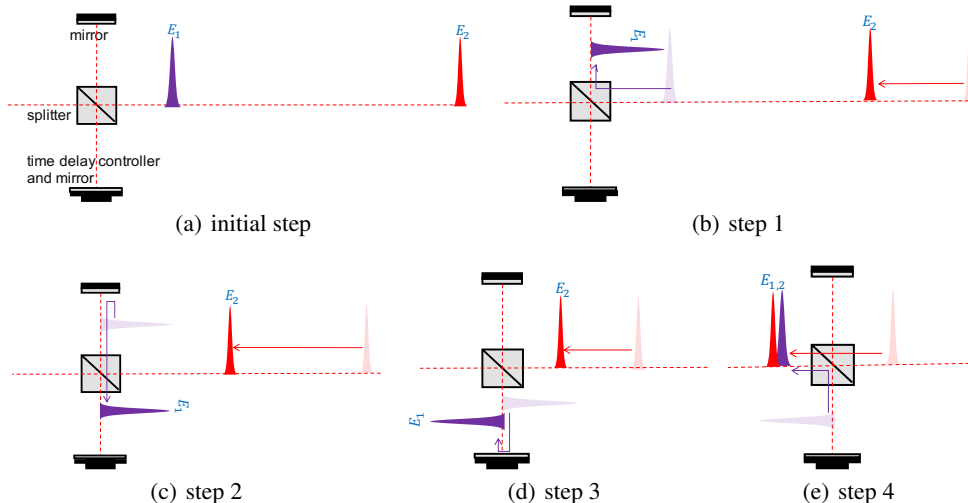


Figure 13: Demonstration of stacking two pulses with additional time delay. (a) shows the initial state of the 2 pulses before stacking. (b)-(e) show the (chronological) procedure of the stacking 2 pulses with imposing additional time delay. The former (latter) pulse was plotted with purple (red). The solid (transparent) plotted pulse shows the pulse position at the current (last) step. The arrow denotes the shifting value of a pulse.

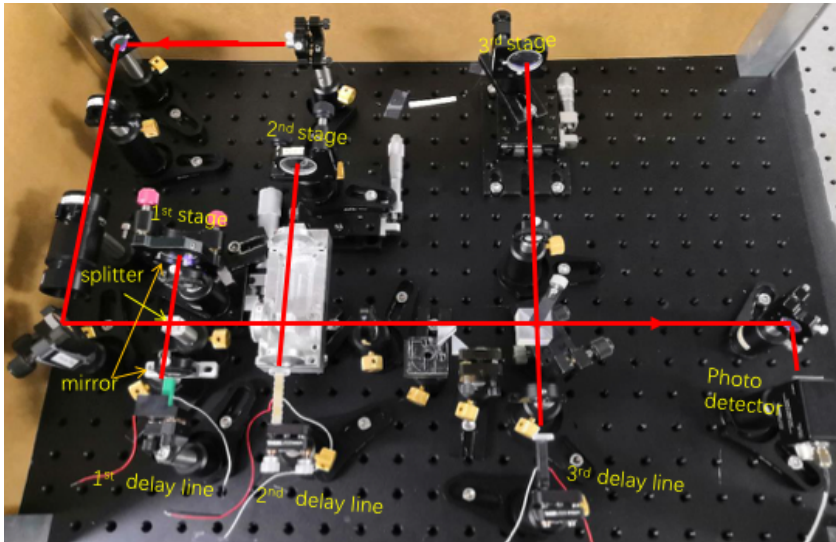


Figure 14: Real world optical pulse stacking system.

## B ADDITIONAL EXPERIMENTS

### B.1 EXPERIMENTAL SETTING

We evaluated the performance of PPO, TD3, and SAC in our OPS environment. For each of PPO, TD3, and SAC, we performed hyperparameter search to achieve better performance. For the search, we trained on 5-stage OPS environment with medium difficulty. Each of the hyperparameter sets was repeated with 3 random seeds. For each algorithm, the best hyperparameter set was decided based on the final performance in the testing environment. After the search, each of the best hyperparameter sets was used to run experiments with 10 different random seeds on all scenarios. The hyperparameter range and selected value of TD3 can be found in table 2, hyperparameter range and selected value of SAC can be found in table 3, and hyperparameter range and selected value of PPO can be found in table 4.



Hyperparameter	Range	Best-selected
Size of the replay buffer	{1000,10000,100000}	10000
Step of collect transition before training	{100, 1000, 10000}	1000
Unroll Length/ $n$ -step	{1,10, 100}	100
Training epochs per update	{1,10, 100}	100
Discount factor ( $\gamma$ )	{0.98, 0.99, 0.999}	0.98
Noise type	{'normal', 'ornstein-uhlenbeck', None}	'normal'
Noise standard value	{0.1, 0.3, 0.5, 0.7, 0.9}	0.7
Learning rate	{0.0001, 0.0003, 0.001,0.003,0.01}	0.001
Policy network hidden layer	{1, 2, 3}	2
Policy network hidden dimension	{64, 128, 256}	256
Optimizer	Adam	Adam

Table 2: TD3: ranges used during the hyperparameter search and the final selected values.

Hyperparameter	Range	Best-selected
Size of the replay buffer	{1000,10000,100000}	10000
Step of collect transition before training	{100, 1000, 10000}	1000
Unroll Length/ $n$ -step	{1,10, 100}	1
Training epochs per update	{1,10, 100}	1
Discount factor ( $\gamma$ )	{0.98, 0.99, 0.999}	0.98
Generalized State Dependent Exploration (gSDE)	{True, False}	True
Soft update coefficient for "Polyak update" ( $\tau$ )	{0.002,0.005, 0.01, 0.02}	0.005
Learning rate	{0.0001, 0.0003, 0.001,0.003,0.01}	0.001
Policy network hidden layer	{1, 2, 3}	2
Policy network hidden dimension	{64, 128, 256}	256
Optimizer	Adam	Adam

Table 3: SAC: ranges used during the hyperparameter search and the final selected values.

Hyperparameter	Range	Best-selected
Unroll Length/ $n$ -step	{128,256,512,1024,2048}	1024
Training epochs per update	{1,5,10}	10
Clipping range	{0.1,0.2,0.4}	0.2
Discount factor ( $\gamma$ )	{0.98, 0.99, 0.999}	0.98
Entropy Coefficient	{0, 0.001, 0.01, 0.1}	0.01
GAE ( $\lambda$ )	{0.90, 0.95, 0.98, 0.99}	0.95
Value function coefficient	{0.1,0.3,0.5,0.7,0.9}	0.5
Learning rate	{0.0001, 0.0003, 0.001,0.003,0.01}	0.001
Gradient norm clipping	{0.1, 0.5, 1.0, 5.0}	0.5
Policy network hidden layer	{1, 2, 3}	2
Policy network hidden dimension	{64, 128, 256}	256
Optimizer	Adam	Adam

Table 4: PPO: ranges used during the hyperparameter search and the final selected values.

## B.2 RESULTS ON CONTROLLING OPS ENVIRONMENT

We reported the training curve (training reward w.r.t. iterations) and testing curve (return (stacked pulse power  $P_N$ ) w.r.t. testing iterations) on the 4-stage OPS environment in fig. 15, and on the 6-stage OPS environment in fig. 16. As can be seen, the performance of TD3 and SAC is higher than PPO. Compared with fig. 15 (4-stage), and fig. 7 (5-stage) to fig. 16 (6-stage), with the increase of stage number, the training convergence became slower, and the final return  $P_N$  became smaller, especially for medium mode and hard mode difficulty.

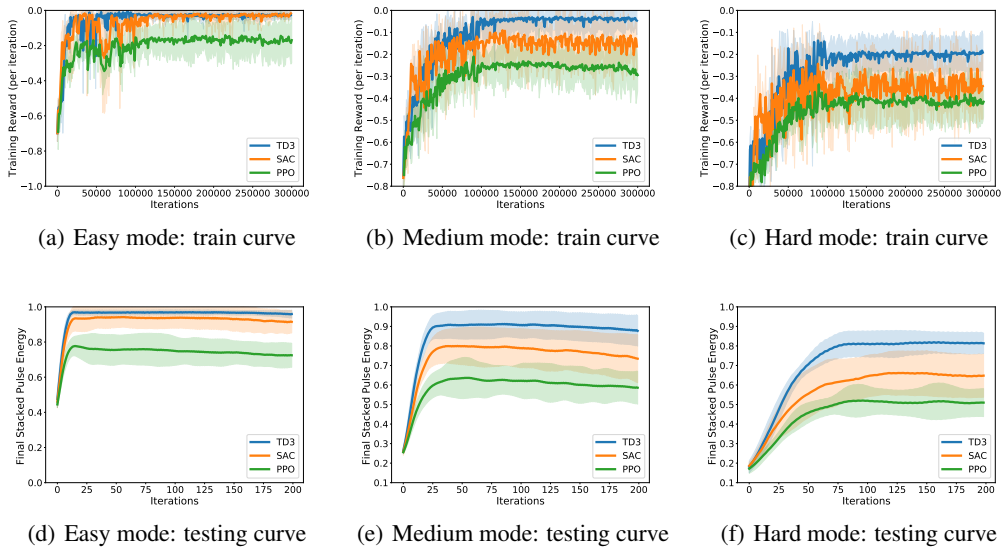


Figure 15: 4-stage OPS experiments. Training reward was plotted for (a) easy mode, (b) medium mode, and (c) hard mode. Evaluation of the stacked pulse power  $P_4$  (normalized) of testing environment was plotted for (d) easy mode, (e) medium mode, and (f) hard mode.

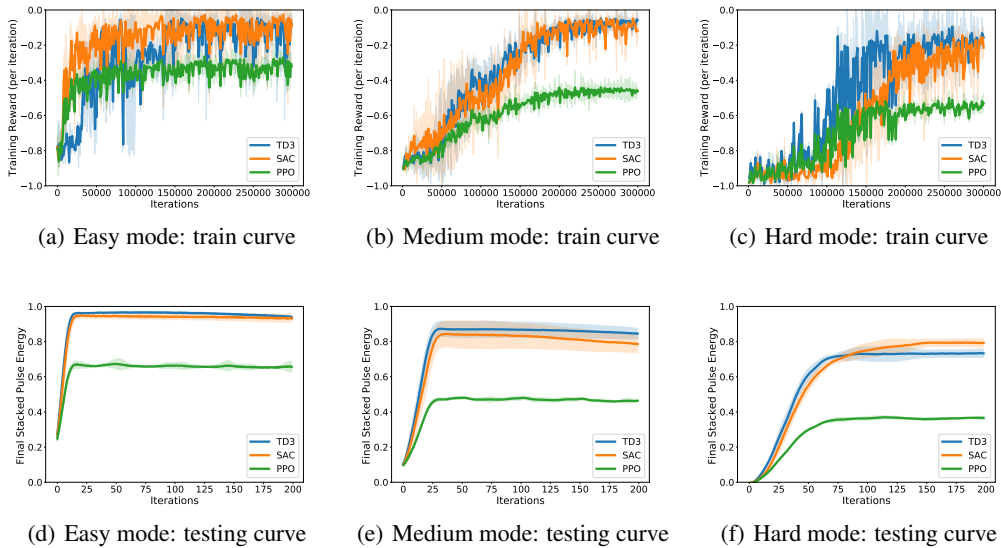


Figure 16: 6-stage OPS experiments. Training reward was plotted for (a) easy mode, (b) medium mode, and (c) hard mode. Evaluation of the stacked pulse power  $P_6$  (normalized) of testing environment was plotted for (d) easy mode, (e) medium mode, and (f) hard mode.

### B.3 DEMONSTRATION OF THE CONTROLLING OPS ENVIRONMENT

Figure 17 shows the pulse trains on a 5-stage hard mode OPS system controlled by TD3 from the random initial state. The animation of this process can be seen in supplementary video 1. It is seen that TD3 algorithm could achieve (local) maximum power within 40 iterations.

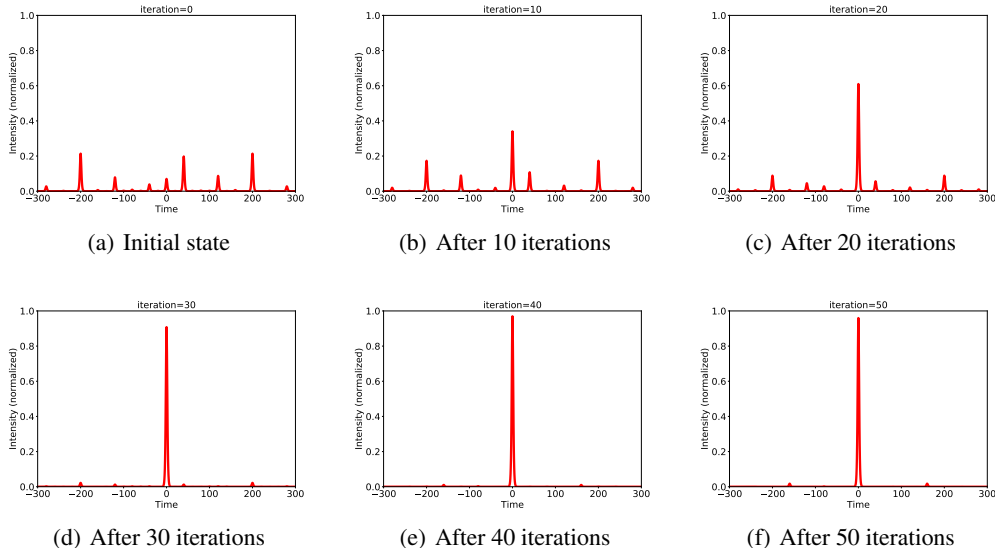


Figure 17: Demonstration of the controlling 5-stage OPS hard mode testing environment by TD3 algorithm after training. (a): initial state of pulses; (b) pulse state after 10 control iterations; (c) pulse state after 20 control iterations; (d) pulse state after 30 control iterations; (e) pulse state after 40 control iterations; (f) pulse state after 50 control iterations.

## C DISCUSSION

### C.1 REAL-WORLD ENVIRONMENT AND SIMULATION ENVIRONMENT

Deploying the RL algorithm in the real-world optics system requires converting optical signal to electrical analog signal using photo-detector (PD), then converting the analog signal to digital signal using an analog-to-digital converter (ADC). These two conversions cost some additional time to process the signal and cause feedback delay. At a conservative estimate, the regular PD and ADC processing takes 0.01s per step. Then it would be possible to implement deep reinforcement algorithms on FPGAs to create control output by the feed of digital observation signal. In a proper implementation, FPGA computing time would be less than 0.01s per step. Including signal converting, neural-network inference time, and time-delay of the optical-mirror driver (controller), the time cost per control step is in the magnitude of 0.1s<sup>4</sup>. But in our simulation system, we could speed up the control step by at least 10 times (with GPUs). More importantly, for RL training on real-world OPS systems, it needs to manually tune the optical devices when the optical beams are totally misaligned caused by the exploring process of RL. The initial alignment of the complex OPS system is usually tuned by experts to take several hours even several days, that the time-cost is depending on the system complexity<sup>5</sup>. But in our simulation system, the initial alignment could be done by simply "reset" the environment. So the value of the simulation environment can be summarized as:

- Faster control process than a real-world experiment.
- Easy to "reset" (and initial align) the environment, while it takes a lot of works to reset or initial align a real-world experiment.
- It is safer and cheaper. In real-world experiments, it has potential risk when the optical beams are totally misaligned, because refracted light is non-predictable and may shed on experimenters.

<sup>4</sup>With expensive high-speed PD, AD/DA cards, and optical mirror driver, as well as efficient FPGA implementation, the control-speed time would be reduced to 0.01s. But it will increase the budget of the devices.

<sup>5</sup>We cannot detect any stacking signal when the optical beams are totally misaligned. So the RL algorithms would fail. It needs to align manually in this scenario.

## C.2 REAL-WORLD EVALUATION

The impact of the simulation must be valued by the real measurement. Part of the correctness of our simulation has been evaluated by the simplified beam combining experiments (Tünnermann & Shirakawa, 2019; Yang et al., 2020). Specifically, Tünnermann & Shirakawa (2019) implemented a simple real experiment and the same simulation, the authors found the simulation is valuable. Our simulation and experimental settings are complicated than Tünnermann & Shirakawa (2019), but the physics behind them is the same. Actually, if we set stage number =1, our simulation is almost the same as Tünnermann & Shirakawa (2019). We will do detailed real experiments and justification in the near future.

## C.3 POTENTIAL IMPACT AND ADDITIONAL RELATED WORKS

**Machine learning community.** High-dimensional real-world reinforcement learning problems are extremely challenging (Dulac-Arnold et al., 2019). In our simulation environment, if we choose a quite large N-stage number with hard mode, controlling the environment could become high-dimensional and difficult. Few recent works studied the distribution shift in RL (Agarwal et al., 2021; Du et al., 2019b). In the hard mode of the OPS environment, the noise distribution of the testing environment is different from the noise distribution of the training environment. Therefore our simulation environment is beneficial to solve the hard and realistic reinforcement learning problems. In recent years, statistical procedures have been developed to promote low-dimensional structures using convex relaxations, rather than directly solving the nonconvex problems (Chen & Chi, 2018; Chi et al., 2018). As shown in fig. 2, we know the function of the OPS objective (if ignoring noise). The function typically provides us with much richer structural information and physical constraints. So it is possible to explore the additional information about the function of the OPS and incorporating it with RL algorithms. In many of the real-world cases, we are not interested in "generic" nonconvex problems, but rather, we focus on more specific nonconvex control with physical constrain or some known objective function (Miryoosefi et al., 2019). Exploring the nonconvex and periodic objective of OPS would benefit the real-world RL problems that including some structural information.

**Optics community.** High pulse energy lasers can be used in laser accelerators, large-scale material processing, and medicine (Fermann & Hartl, 2013). Optical (coherent) pulse stacking is one of the easiest and promising ways to scale the pulse energy (Tünnermann & Shirakawa, 2017). However, the conventional control algorithms for OPS are not very effective (Du et al., 2019a). RL methods are able to control this kind of multi-dimensional and nonlinear environment. Similar to our OPS control system, all of the optical control problems are affected by the nonlinearity and periodicity of the light inference (as shown in fig. 2), including coherent optical inference (Wetzstein et al., 2020) and linear optical sampling (Dorrer et al., 2003), which can be used for precise measurement, industrial manufacturing, and scientific research. We believe our simulation is one of the important and typical optical control environments. Beyond OPS, RL methods have the potential to drive the next generation of optical laser technologies even the next generation of scientific control technologies (Genty et al., 2020). This is because many phenomena in optics are nonlinear and multidimensional, with noise-sensitive dynamics that are extremely challenging to model using conventional methods.