CryoRL: Reinforcement Learning Enables Efficient Cryo-EM Data Collection

Quanfu Fan^{1,6}, Yilai Li^{2,6}, Yuguang Yao³, John Cohn¹, Sijia Liu³, Seychelle M. Vos^{4,7}, and Michael A. Cianfrocco^{2,5,7}

¹ MIT-IBM Watson AI Lab, Cambridge, MA USA

² Life Sciences Institute, University of Michigan, Ann Arbor, MI USA

³ Department of Computer Science and Engineering, Michigan State University, East Lansing, MI USA

⁴ Department of Biology, Massachusetts Institute of Technology, Cambridge, MA USA

⁵ Department of Biological Chemistry, Michigan Medicine, University of Michigan,

Ann Arbor, MI USA

⁶ Equal contributions

⁷ For correspondence: M.A.C. mcianfro@umich.edu & S.M.V. seyvos@mit.edu

Abstract. Single-particle cryo-electron microscopy (cryo-EM) has become one of the mainstream structural biology techniques because of its ability to determine high-resolution structures of dynamic bio-molecules. However, cryo-EM data acquisition remains expensive and labor-intensive, requiring substantial expertise. Structural biologists need a more efficient and objective method to collect the best data in a limited time frame. We formulate the cryo-EM data collection task as an optimization problem in this work. The goal is to maximize the total number of good images taken within a specified period. We show that reinforcement learning offers an effective way to plan cryo-EM data collection, successfully navigating heterogenous cryo-EM grids. The approach we developed, cryoRL, demonstrates better performance than average users for data collection under similar settings.

1 Introduction

Single-particle cryo-electron microscopy (cryo-EM) has become one of the mainstream structural biology techniques due to its ability to solve the structures of many bio-molecules with moderate heterogeneity and without the need for crystallization. In recent years, continued software development has led to automation in both data collection and image processing [2]. Moreover, with the improvement of the detectors and microscopes techniques, data acquisition has been dramatically accelerated [7,35].

Cryo-EM serves as a critical tool in the development of vaccines and therapeutics to combat COVID-19 by SARS-CoV-2 (Fig. 1). Within weeks of the release of the genomic sequence of SARS-CoV-2, cryo-EM determined the first SARS-CoV-2 spike protein structure [38]. Since this original publication, cryo-EM was used to determine additional SARS-CoV-2 structures such as spike protein bound to

antibody fragments [20,28], remdesivir bound to SARS-CoV-2 RNA-dependent RNA polymerase [6,41,13], and reconstructions of intact SARS-CoV-2 virions [39,11].



Fig. 1: Cryo-EM structure of the SARS-CoV-2 spike protein.

Despite these advances, cryo-EM data collection remains *ad-hoc*, rudimentary, and subjective. Due to variations in sample quality across a cryo-EM grid, users collect images at different magnifications ranging from resolutions of 0.66 mm to 500 Å. Significant user expertise enables experts to define and refine locations suitable for data collection. To provide objective

feedback, "on-the-fly" image processing [16,?] can confirm high-quality regions on the cryo-EM sample. Despite this information, data collection remains highly subjective. Cryo-EM is also an expensive technique, further compounding challenges faced by users. Purchasing, preparing, and installing a top cryo-electron microscope can cost about \$10 to 20 million USD, and the daily operational cost can be around \$10,000 USD [9]. Therefore, structural biologists need methods that can help collect the best data possible in a limited amount of time.

In this paper, we formulate the data collection problem as an optimization task where the goal is to learn intelligent strategies from data to guide the microscope movement, possibly via manual suggestion or robotic manipulation. We model the optimization problem as a Markov decision process and propose to solve it by combing supervised classification and deep reinforcement learning (RL) [31]. We present a new data acquisition algorithm, cryoRL, which enables data collection with no subjective decisions, no user intervention, and increased efficiency. To address the potential enormously large action space in our problem, we further propose to eliminate irrelevant or sub-optimal actions based on the classification results to enable effective policy exploration, which improves cryoRL in both efficiency and accuracy. As compared with human subjects, cryoRL achieves better performance than average users. To the best of our knowledge, cryoRL is the first AI-based algorithm in cryo-EM data acquisition such that a policy is learned and can directly help the user steer the microscope.

We collected datasets on different grid types to design, implement, and test cryoRL. The first of its kind, our data collection involves no user decision; instead, we selected the areas from a systematic pattern of data collection to obtain images of all holes and micrographs (See Section 3 for details of our data). Our datasets will be released to the public to serve as a critical benchmark for evaluating cryo-EM data collection algorithms.

To summarize, our high level conclusions and contributions include:

CryoRL achieves better performance than other popular optimization techniques such as Genetic Algorithm [36] and Simulated Annealing [12], and demonstrates good generalization capability;

CryoRL: Reinforcement Learning Enables Efficient Cryo-EM Data Collection

- CryoRL with our proposed invalid action elimination runs 2~3 times faster than the vanilla DQN baseline while enabling more robust policy learning;
- CryoRL offers a new approach to cryo-EM data collection that demonstrates promising results by outperforming average users in a human performance study;
- We are providing a first-kind-of cryo-EM dataset that is critical for algorithm development and benchmarking.

2 Related Work

There are currently no automated, 'intelligent' cryo-EM data collection approaches. Instead, subjective decision-making drives cryo-EM data acquisition. To guide user-driven data collection, on-the-fly image analysis provides feedback on data quality, including Appion [16], Warp[32], and cryoSPARC Live. To provide more objective measures of data quality to users, researchers have developed a pretrained deep learning-based micrograph assessment model [21] and downstream on-the-fly data processing [30]. However, despite these efforts, on-the-fly processing requires a sizeable number of micrographs before providing useful feedback. Data collection requires user training to develop expertise to guide data collection in the most efficient manner possible.

Reinforcement learning (RL) has been widely applied to address practical optimization problems such as network planning [1], vehicle routing [23], online recommendation [44], robot trajectory optimization [14] as well as game playing [29]. Some practical applications have adopted RL for enabling fast data collection or processing. For example in [33,43], RL-based methods are proposed to optimize unmanned aerial vehicle flight trajectories for efficient data collection. In [15], a deep RL agent is used to condition the state of the probe for autonomous Scanning Probe Microscopy (SPM). Nevertheless, automating data collection by machine learning techniques in real-world scenarios remains an under-explored problem.

Large action spaces are a common problem to deal with in RL. Existing techniques include action masking [4,40] to mask out invalid actions, action elimination [42] to remove inferior actions, and action reshaping [10] to transform a discrete action space to a simpler one or a continuous one. Our proposed action elimination is in a similar spirit to Action Elimination Network (AEN) [42], but instead relies on the estimated quality of a hole to exclude invalid actions directly rather than learning to reduce the action space.

3 Cryo-EM Data Collection

The general practice of data acquisition in cryo-EM is abstracted in Fig. 2. Typically, a purified biological sample is dispensed and vitrified onto a grid comprised of gold or copper support bars. A grid contains a mesh of squares, and each square has a lattice of regularly-spaced holes. Ideally, within each hole,



Fig. 2: Overview of cryo-EM data collection. A purified sample is prepared and vitrified on the support grid. The atlas image provides a low magnification overview by stitching multiple "grid-level" images into a single montage. Next, users will select specific squares to image at medium magnification. After inspection, the user selects "patch" areas on the square to inspect holes with higher magnification, using the patch image to decide holes to collect for micrographs. The micrographs contain high-resolution images for downstream data processing.

there are vitrified single-particles related to the sample of interest, where data collection amounts to users recording images of each hole as micrographs.

Cryo-EM samples exhibit heterogeneity across the specimen. Whereas there are many local correlations between squares and holes on the grid, many holes are empty, contain aggregates, or contain non-vitreous ice contamination. The user has no prior knowledge of such distribution until the square-level or hole-level images are acquired, which can be captured by the microscope by changing to different magnifications. Note that each greater magnification requires significant time for the microscope to move and settle. Moreover, because the time on the microscope is precious and limited, data collection can typically cover less than 1% of the total grid. The user needs to navigate through the "grid-square-hole" hierarchy and collect the best micrographs in a limited time.

In this paper, we suppose that a user preselects a set of squares and patch-level images by a quick atlas survey. We formulate the data acquisition task to find the highest quality holes and plan the overall data collection route with cryoRL. Although this is not the traditional way people collect cryo-EM data, we believe such prerequisites provide a global understanding of the hole quality distribution on a grid by taking a series of low to medium magnification square and patch level images.



Fig. 3: cryoRL-guided cryo-EM data collection. cryoRL consists of an offline hole-level classifier to estimate the quality of holes and a deep Q-network to learn effective strategies for steering the microscope for data collection. The hole classifier outputs serve as features of the RL network. The agent (or user) provides feedback (rewards) to the system according to the microscope movement and the measured quality of the micrographs taken for the holes recommended by the system. In addition, to address the potential issue of large action spaces in our problem, we developed an efficient method to eliminate invalid actions based on the quality of holes.

Each micrograph has an objective measure of data quality, which is the goodness-of-fit for the frequency domain when estimating the defocus of the micrograph. We introduce the term "CTFMaxRes" to be the maximum resolution (Å) for the fit of the contrast transfer function (CTF) to a given micrograph using the program CTTFIND4 [26]. CTFMaxRes is calculated from the 1D power spectrum of the micrograph and estimates the maximum resolution for the detected CTF oscillations [5]. The field of cryo-EM utilizes CTFMaxRes to provide an indirect metric for data quality. In general, the lower this value, the higher the quality of the micrograph. CryoRL will predict the quality of each hole from the patch-level image using an image classifier (Section 4). For simplicity, we define CTFMaxRes as the CTF value for this paper.

4 RL-based Approach

4.1 Overview

During data collection, a user needs to make decisions based on the quality of images taken at different magnification levels: grid, square, and patch-level. Given that the data are visually similar and there are significant costs (time) of moving to other grid areas and refocusing, there is no easy planning that a user can make manually in a regular data collection. As a result, the user explores only a small portion of a grid, making the data collection process inefficient and subjective.

In this work, we formulate the data collection problem as planning an optimal path for operating the microscope. The goal is to move the microscope to explore desired places on a grid in a given amount of time, with the operational cost

constraints taken into account (Section 4.2). We propose to solve the path planning problem by RL, a technique that has demonstrated success in many vision applications [18]. Compared to other widely used optimization solvers such as Genetic Algorithm (GA) [36] and Simulated Annealing (SA) [12], RL is more suitable for modeling sequential problems and possibly less heuristic in system design.

As illustrated in Fig. 3, our proposed approach combines an image classifier and an RL network to enable automatic planning of microscope movement. The supervised classifier categorizes a hole into low or high quality based on its CTF value. Efficient hierarchical feature representations for cryo-EM images at different magnification levels are generated from the classification results. These features, along with the observation history, are exploited to train a deep Q-network (DQN) [22] to assess the status of all the unvisited holes and suggest the best holes to look at next. We further design a rewarding mechanism to drive the learning of DQN. The design in general values small microscope movements to avoid wasted time. For example, moving to a different patch on the same grid-level image receives a higher score than changing to an entirely new gridlevel image (Section 4.3). Finally, to handle the potentially large action space in our problem, we propose a method to eliminate invalid actions, which not only results in a significant speedup of CryoRL by $2\sim3$ times, but also improves the robustness of the approach (Section 4.3).

As mentioned earlier, a human user can usually cover a small portion of the grid during a data collection session. In contrast, one significant advantage of our proposed approach is that it allows for a substantially larger exploration of the grid by the microscope as the approach learns to focus on promising regions with high-quality data. We demonstrate in Section 5 that our system is highly effective, achieving comparable performance to human subjects.

4.2 Problem Formulation



Fig. 4: A schematic illustration of a path showing the microscope movement planned in data collection. Different microscopic operations are associated with different costs, which are indicated by the edge width.

As previously described, crvo-EM data collection is steering the microscope hierarchically at different magnification levels to explore a grid to identify highquality micrographs. This sequential process involves several mechanical operations to allow microscope navigation to different regions of a grid. The process of data collection involves area switching (changing to a new grid-level image), square switching (changing to a new square-level image), and patch switching (changing to a new patchlevel image). Since the data distribution is non-uniform on a grid and it takes time to prepare the microscope

for imaging at different levels, an automatic method to guide the data exploration more intelligently will improve data quality and efficiency for data collection.

As shown in Fig. 4, an effective data collection session aims at finding a sequence of holes where there is a considerable portion of high-quality micrographs. Let $\mathcal{H} = \{h_l | l = 1 \cdots n_h\}$ be a sequence of holes in a set of patches \mathcal{P} sampled from different square-level and grid-level images (\mathcal{S} and \mathcal{G}) by the user. We denote \mathcal{P}_{h_l} , \mathcal{S}_{h_l} and \mathcal{G}_{h_l} as the corresponding patch-level, square-level and grid-level images of h_l , respectively. Also, $ctf(h_l)$ is a function representing the CTF value of a hole h_l . Our goal is to identify a maximum subset of holes from \mathcal{H} with low-CTF values in a given amount of time τ . Mathematically, this is equivalent to optimizing an object function as follows,

$$\max \sum_{l=0}^{n_h - 1} \left(\rho(h_l) - c(t(h_l)) \right) \quad \text{s.t.} \quad \sum_{l=0}^{n_h - 1} t(h_l) \le \tau \tag{1}$$

where $\rho(h_l)$ be such an indicator function for a hole h that

$$\rho(h_l) = \begin{cases} 1 & \text{if } ctf(h_l) \le 6.0\\ 0 & \text{otherwise} \end{cases} \tag{2}$$

and c is a cost associated with the corresponding microscope operation and determined by the total amount of time $t(h_l)$ spent on h_l . In this work, we define $t(h_l)$ in minutes by the movement of the microscope, i.e,

$$t(h_l) = \begin{cases} 2.0 & \text{if } \mathcal{P}_{h_{l-1}} = \mathcal{P}_{h_l} \text{ (same patch)} \\ 3.0 & \text{if } \mathcal{P}_{h_{l-1}} \neq \mathcal{P}_{h_l}, \mathcal{S}_{h_{l-1}} = \mathcal{S}_{h_l} \text{ (same square)} \\ 5.0 & \text{if } \mathcal{S}_{h_{l-1}} \neq \mathcal{S}_{h_l}, \mathcal{G}_{h_{l-1}} = \mathcal{G}_{h_l} \text{ (same grid)} \\ 10.0 & \text{if } \mathcal{G}_{h_{l-1}} \neq \mathcal{G}_{h_l} \text{ (different grid)} \end{cases}$$

Note that the time t above is set in a way so that it highly corresponds to the natural time of the microscope movement in real-world scenarios. Nevertheless, in practice, it can be more precisely calculated based on the distance of the microscope movement and other factors.

By setting $r(h_l) = \rho(h_l) - c(t(h_l)))$, we can further rewrite Eq. 1 as

$$\max \sum_{l=0}^{n_h - 1} r(h_l) \quad \text{s.t.} \ \sum_{l=0}^{n_h - 1} t(h_l) \le \tau \tag{3}$$

Eq. 3 has the same form as the standard accumulative reward (without a discount factor) that is maximized in RL [31]. In what follows, we describe how to design a RL system to solve the path optimization problem in Eq. 3.

4.3 Path Optimization by Reinforcement Learning

We study the cryo-EM data acquisition task by RL, where an agent interacts with environment (i.e. the grid here) by sequentially selecting holes for taking

micrographs over a sequence of time steps, with an objective to maximize the cumulative reward described in Eq. 3. We briefly describe the basic components of our system as follows.

Environment: the atlas or grid.

Agent: a robot or user steering the microscope.

States. Let $u_i \in \{0, 1\}$ be a binary variable denoting the status of hole, i.e. visited or unvisited. Then a state s in our setting can be represented by a sequence of holes and their corresponding statuses $s = \langle (h_1, u_1), (h_2, u_2), ..., (h_{n_h}, u_n) \rangle$ where n is the total number of holes.

Actions. An action a_i of the agent in our system is to move the microscope to the next target hole h_i for imaging. Note that in our case, any unvisited hole has a chance to be picked by the agent as a target, thus the action space is large. Also, during tests, the number of holes (i.e actions) is unknown. Instead of adopting more sophisticated methods to handle continuous action space as proposed in [17,19], we simply modify the Q-network to estimate the Q-value for every single hole rather than all of them at once. We show this suffices for handling the large action space in our case.

Rewards. We assign a positive reward 1.0 to the agent if an action results in a target hole with a CTF value less than 6.0Å and 0.0 otherwise. The agent also receives a negative reward depending on the operational cost associated with a hole visit. Specifically, we model the negative reward as $c(h_l) = 1.0 - e^{-\beta(t(h_l) - t_0)}(\beta > 0, t_0 \ge 0)$. We empirically set β and t_0 to 0.185 and 2.0, which define the final reward function for our RL system as,

$$r(a_{i}) = \begin{cases} 1.0 & \text{if } ctf(h_{l}) < 6.0 \& \mathcal{P}_{h_{i-1}} = \mathcal{P}_{h_{i}} \\ 0.57 & \text{if } ctf(h_{l}) < 6.0 \& \mathcal{P}_{h_{i-1}} \neq \mathcal{P}_{h_{i}} \& \mathcal{S}_{h_{i-1}} = \mathcal{S}_{h_{i}} \\ 0.23 & \text{if } ctf(h_{l}) < 6.0 \& \mathcal{S}_{h_{i-1}} \neq \mathcal{S}_{h_{i}} \& \mathcal{G}_{h_{i-1}} = \mathcal{G}_{h_{i}} \\ 0.09 & \text{if } ctf(h_{l}) < 6.0 \& \mathcal{G}_{h_{i-1}} \neq \mathcal{G}_{h_{i}} \\ 0.0 & otherwise \end{cases}$$

Note that the design principle of these rewards is to reward more small microscope movement. As shown later in the experiments (Section 5.3), CryoRL is not sensitive to the changes of the rewards as long as the design described above is followed.

Deep Q-learning We apply the deep Q-learning approach proposed in [22] to learn our policy for cryo-EM data collection. The goal of the agent is to select a sequence of actions (i.e. holes) based on a policy to maximize future rewards (i.e the total number of low-CTF holes). In Q-learning, this is achieved by maximizing the action-value function $Q^*(s, a)$, i.e. the maximum expected return achievable by any strategy (or policy) π , given an observation (or state) s and some action a to take. In other words, $Q^*(s, a) = \max_{\pi} E[R_t|s_t = s, a_t = a, \pi]$ where $R_t = \sum_t^{\infty} \gamma^{t-1} r_t$ is the accumulated future rewards with a discount factor γ . Q^{*} can be found by solving the Bellman Equation [31] as follows,

$$Q^*(s,a) = E_{s'}[r + \gamma \max_{a'} Q^*(s',a')|s,a]$$
(4)

In practice, the state-action space can be enormous, thus in [22], a deep neural network parameterized by θ is applied to approximate the action-value function. The network is referred to as Deep Q-Network (DQN) in the original paper. DQN can be trained by minimizing the following loss functions $L(\theta)$,

$$L(\theta) = E_{s,a,r,s'}[(y - Q(s,a;\theta)^2]$$
(5)

where $y = E_{s'}[r + \gamma \max_{a'} Q(s', a')|s, a]$ is the target for the current iteration. The derivatives of the loss function $L(\theta)$ are expressed as follows:

$$\nabla_{\theta} L(\theta) = E_{s,a,r,s'}[(r + \gamma \max_{a'} Q(s',a';\theta') - Q(s,a;\theta))\nabla_{\theta} Q(s,a;\theta)]$$
(6)

Experience replay is further adopted in [22] to store into memory the transition at each time-step, i.e (s_t, a_t, r_t, s_{t+1}) , and then sample the stored samples for model update during training.





Fig. 5: The architecture of DQN. The network has only one single output node to estimate the Q-value for an action-state pair.

Fig. 6: Examples of hole images and their CTF values. A hole with a CTF ≤ 6 is considered good in our paper.

DQN with Action Elimination via Patch Ranking In a regular scenario such as playing Atari [22] where the action space is small and fixed, a network can be trained to predict all the actions at once. However, this is not suitable for our case as our action space is not fixed and can grow large depending on the training data size. To deal with this issue, we modify the Q-network to predict the Q-value for each hole (i.e., action) using one single output, as shown in Fig. 5. The Q-value for all the actions can then be batch processed and the ϵ -greedy scheme is applied for action selection. The DQN used in our work is a 3-layer fully connected network. The size of each layer is 128, 256 and 128, respectively.

The potentially enormously large action space in our problem makes policy exploration quite inefficient as most actions sampled from such a space are not useful. To avoid executing too many sub-optimal actions in learning, we propose an effective method to reduce the action space by restricting the valid actions to a small portion of holes predicted by the classifier as low-CTFs. We start by ranking all the grid-level images by their numbers of low-CTF predictions, from high to low, and then the patches in the same way. The high-ranked patches are likely to contain more valid holes and should be visited more frequently during the learning. The pre-specified duration (i.e. τ in Eq. 1) as well as the switching costs $t(h_l)$ defined in Section 4 allow us to obtain an upper limit N_{max} of good

Feature Type	Definition	Value
hole	is it low-CTF? is it visited?	$\{0,1\}$
patch/square/grid	# of unvisited holes # of unvisited lCTFs # of visited holes # of visited lCTFs	$0 \sim 150^*$
microscope movement	a new patch-level image? a new square-level image? a new grid-level image?	$\{ \begin{matrix} 0,1 \\ \{ 0,1 \\ \{ 0,1 \} \\ \{ 0,1 \} \end{matrix}$

*: the maximum number of holes allowed in a grid-level image in our setting Table 1: Input features to DQN

holes if all the holes are assumed to be low CTF and visited in the sorted order described above. We then select a minimum set of patches P with a total number of low-CTF predictions $\geq \beta N_{max}$ ($\beta > 0$), and all the holes in P define a reduced new space for Q-learning. Here, β is a user-defined parameter to control the size of the valid action set. Our approach is not sensitive to β , and any number between 1.0 ~2.5 works reasonably well. We thus empirically set β to 1.5 in tests and a larger number in training to enlarge the exploration space for CryoRL.

The details of our algorithm can be found in the appendix. Unlike the Elimination Network (EAN) proposed in [42], our approach redefines the action space before Q-learning, so it can be applied to any policy learners without modification of them. We show later in the experiments that our approach results in a significant speedup of $2\sim3$ times over the vanilla DQN and improves other policy learners such as A2C [24] and C51 [3] remarkably.

Features to DQN The quality of a hole is directly determined by its CTF value. Similarly, the number of low-CTF holes (lCTFs) in a hole-level image indicates the quality (or value) of the image, and a good RL policy should always consider prioritizing high-quality patches first in planning. The same holds true for square-level and grid-level images. Based on this, we design hierarchical input features to the DQN according to the quality of images at different levels. We also consider the information of microscope movement as it tells whether the microscope is exploring a new region or staying at the same region. The details of these features can be found in Table 1. Finally, a sequence of these features for the last k-1 visited holes as well as the current one to be visited are concatenated together to form the input to DQN. In our experiments, k is empirically set to 4. Hole-level Classification We trained the hole-level classifier offline by cropping out the holes in our data using the location provided in the meta data. Fig. 6 illustrates a few examples of hole images. These images are actually visually ambiguous, confounding the task of building generalized hole classifiers, as shown in Section 5.2. Using an offline classifier enables fast learning of the Q function as only the Q-network is updated in training and its input features can be computed

efficiently. However, it is possible to jointly learn the classifier and DQN to further improve performance. We leave this possibility for future work.

11

5 Experiments

5.1 Experimental Setup

Dataset To design and evaluate the performance of cryoRL, we collected an "unbiased" cryo-EM dataset (**Y1**) to provide a systematic overview all squares, patches, holes, and micrographs within a defined region of a cryo-EM grid. Specifically, aldolase at a concentration of 1.6 mg/ml was dispensed on a support grid and prepared using a Vitrobot. Instead of picking the most promising squares and holes, we randomly selected 31 squares across the whole grid and imaged almost all the holes in these selected squares. This resulted in a dataset of 4017 micrographs from holes in these 31 squares. Overall, the data quality was poor, given that only 33.4% of the micrographs have a CTF below 6 Å. However, this makes the dataset very suitable for developing and testing algorithms for data collection algorithms, because 1) a perfect algorithm will aim to find the best data from mostly bad micrographs, and 2) the "unbiasedness" of this dataset ensures that when an algorithm selects a hole, the corresponding micrograph, and its metric can be provided as feedback.

In addition, we collected another different dataset $(\mathbf{Y2})$ of 3969 micrographs with a different sample and grid type. We split both datasets into training and validation sets by a ratio of 2:1. In the experiments below, we evaluate our approach mainly based on Y1 while using Y2 to test the transferribility of CryoRL.

Training and Evaluation We used the Tianshou reinforcement learning framework [37] to learn cryoRL. Each model was trained with 20 epochs, using the Adam optimizer and an initial learning rate of 0.01. We set the duration in our system to 240 minutes for training, and evaluate the system at 120, 240, 360 and 480 minutes, respectively.

5.2 Main Results

Comparison with Baselines. We first developed a greedy-based method purely based on the hole classification results. This method performs a primary sorting on the grid-level images by their quality (i.e., the total number of low CTF holes), followed by a secondary sorting on the patches within a grid by the quality of patches. The sorted patches are then scanned in order, with only the holes classified as low CTFs visited. While simple, this greedy approach serves as a strong baseline when the hole-level classifier is strong.

We also compare our approach with two other widely used optimization techniques in practice: Genetic Algorithms (GA) [36] and Simulated Annealing (SA) [12]. In these two solvers, solutions are sampled at the patch level rather than at the hole level for efficiency, and the fitness of the solutions are assessed

according to the objective function proposed in this paper, i.e Eq. 1. Since GA and SA are largely based on heuristic, the best solutions determined by them are scanned in a similar way to the greedy-based approach described above during the evaluation.

u					
	Methods	$\tau = 120$	$\tau{=}240$	$\tau{=}360$	$\tau{=}480$
	Random Greedy Genetic Alg. (GA) [36] Simulated Annealing (SA) [12] offline path planing [31]	$ \begin{vmatrix} 2.6 \pm 1.4 \\ 41.8 \pm 2.5 \\ 28.3 \pm 6.5 \\ 39.4 \pm 6.5 \\ 44.3 \pm 0.9 \end{vmatrix} $	5.1 ± 1.6 69.3 ± 3.2 72.3 ± 6.8 73.3 ± 7.0 84.6 ± 6.1	$7.3\pm2.3 \\ 104.9\pm4.9 \\ 115.7\pm7.8 \\ 104.7\pm8.9 \\ 121.4\pm6.7$	9.8 ± 2.2 147.9 ±5.1 150.4 ±6.8 147.9 ±9.6 166.6 ±4.9
	$\begin{array}{c} {\rm CryoRL-DQN} \ ({\rm ours}) \\ {\rm CryoRL-DQN}^{\dagger} \ ({\rm ours}) \end{array}$	$ \begin{array}{c} 41.7 \pm 3.1 \\ \textbf{47.4} \pm 0.5 \end{array} $	86.6±3.0 89.0 ±3.1	132.0 ±2.3 131.8±1.8	$\begin{array}{c} 171.4{\pm}2.0 \\ \textbf{172.6}{\pm}2.0 \end{array}$
	human	31.9 ± 10.6	77.4 ± 6.2	-	-

Table 2: Perf. comparison of CryoRL with baseline approaches on Y1

Table 2 reports the total number of low-CTF holes (#lCTF) found by each approach. For fair comparison, all the results are averaged over 50 trials starting from random picked holes. Here, ResNet50 is used as the offline classifier, which achieves an accuracy around 83% in low-CTF hole classification (see Table 4). The results based on ResNet18 can be found in the appendix. As shown in the table, our approach (cryoRL-DQN) is clearly superior to all the baseline methods, producing quite promising results. With action elimination, the fast version of CryoRL (cryoRL-DQN[†]) improve the performance further. Note that while offline path planning yields comparable performance to our method, it is prohibitively costly in computation.

To further illustrate the advantage of our approach, we plot for each approach the percentage of low-CTF holes over the total number of holes visited by time in Fig. 7. Our approach demonstrates high efficacy in data collection, finding $\sim 95\%$ of the holes in good quality. As a comparison, the percentage of low-CTF holes in Y1 is 33.4% and the classification accuracy of low CTFs is only 83.9%.





Fig. 7: Percentage of lCTF images visited during data collection.

Fig. 8: Runtime comparison of CryoRL-DQN and its fast version with action elimination (CryoRL-DQN^{\dagger}).

We also experimented with several other RL variants including dueling DQN [34], DQN with prioritized replay [27], A2C [24] and C51 [3]. As seen from Table 3, the DQN family overall perform better than A2C and C51. Interestingly, A2C and C51 benefit substantially from action elimination and gain significant performance boosts, suggesting that restricting the actions to smaller

CryoRL: Reinforcement Learning Enables Efficient Cryo-EM Data Collection

Methods	$\tau = 120$	$\tau{=}240$	$\tau {=} 360$	$\tau{=}480$
CryoRL-A2C	35.5 ± 7.8	74.0 ± 9.0	111.3 ± 8.8	147.0±8.8
CryoRL-C51	39.4 ± 4.2	76.3 ± 3.1	109.6 ± 2.0	141.0 ± 2.7
CryoRL-DQN	41.7 ± 3.1	86.6 ± 3.0	132.0 ± 2.3	171.4 ± 2.0
CryoRL-DQN (dueling)	44.6 ± 3.3	89.3 ± 4.4	126.3 ± 4.2	157.4 ± 4.4
CryoRL-DQN (prioritized)	42.5 ± 4.3	86.4 ± 3.9	128.7 ± 5.1	172.0 ± 3.5
CryoRL-A2C [†]	47.0±1.3(+32.3%)	90.8±4.0(+22.7%)	$128.2 \pm 2.4 (+15.1\%)$	$163.9 \pm 4.7 (+11.5\%)$
CryoRL-C51 [†]	$47.4 \pm 0.9 (+20.3\%)$	$82.2 \pm 2.3(+7.7\%)$	$116.9 \pm 1.0(+6.7\%)$	$144.0 \pm 1.9(+2.1\%)$
$CryoRL-DQN^{\dagger}$	47.4±0.5(+13.7%)	89.0±3.1(+2.8%)	$131.8 \pm 1.8 (+0.0\%)$	$172.6 \pm 2.0 (+1.0\%)$
$CryoRL-DQN^{\dagger}$ (dueling)	$47.3 \pm 1.0 (+6.1\%)$	89.4±2.9(+0.0%)	$128.6 \pm 2.0 (+1.8\%)$	$165.4 \pm 2.5 (+5.1\%)$
CryoRL-DQN [†] (prioritized)	$47.2 \pm 1.7(+11.1\%)$	$90.6 \pm 2.8 (+4.9\%)$	$132.9 \pm 3.0 (+3.3\%)$	$174.0 \pm 3.1 (+1.2\%)$

Table 3: Perf. comparison of different CryoRL variants on Y1 († indicates action elimination (Section 4.3)). The performance gains from action elimination are highlighted by numbers in parentheses.

Test	st Training Top		p1 Acc	ol Acc. #lCT			'Fs found		
	classifier	r CryoRL	lCTF	hCTF	all	$\tau{=}120$	$\tau{=}240$	$\tau{=}360$	$\tau{=}480$
Y1	Y1	Y1	83.9	91.2	88.5	$47.4 {\pm} 0.5$	$89.0 {\pm} 3.1$	$131.8 {\pm} 1.8$	$172.6{\pm}2.0$
Y1	М	Y1	66.6	85.1	73.6	44.7 ± 2.2	70.0 ± 4.7	104.0 ± 3.6	$138.9 {\pm} 2.6$
Y2	Μ	Y2	69.5	77.4	73.5	$31.0 {\pm} 6.5$	$56.3{\pm}7.7$	87.1 ± 8.0	$125.5 {\pm} 8.3$
Y2	М	Y1	69.5	77.4	73.5	$20.9 {\pm} 4.3$	$55.8 {\pm} 3.7$	83.1 ± 3.5	$91.8{\pm}3.3$
18	able 4:	Generali	zatior	ı abili	ty of	the offli	ne classi	fier and (JrvoRL.

valid sets helps these methods learn policies more effectively. CryoRL with action elimination also achieves considerable speedups in runtime by $2\sim3$ times, as shown in Fig. 8. Since the performance differences between DQN models are minor, we focus on the vanilla DQN in the analysis below.

Comparison with Human Performance. We developed a simulation tool to benchmark human performance against the performance of cryoRL. Fifteen students from two different cryo-EM labs with various expertise levels were recruited in this human study. The users did not have any prior knowledge of this specific dataset before participating in this study. Patch images containing holes in the same dataset were shown to the user. The user had either 50 or 100 chances to select the holes to take micrographs from, corresponding to the experiment's test duration of 120 or 240 minutes. After each selection, the CTF value for the selected hole was provided to the user. The goal of the users is to select as many "good" holes as possible in 50 or 100 chances. Note that we did not penalize the users for switching to a different patch or square as we did in cryoRL. This encouraged the users to explore different patches initially and, theoretically, resulting in better performance than penalties applied. Nevertheless, we found that cryoRL outperforms the human performance in both time durations (Table 2).

Transferrability. We further evaluate the *transferability* of our proposed approach based on a new dataset \mathbf{M} , which consists of users' daily use of microscope in real-life scenarios from 2019 to 2021. Different from Y1 and Y2 data where almost each hole in the patch images was imaged, \mathbf{M} were only

Training	Test Duration				
Duration	$ \tau=120$	$\tau = 240$	$\tau {=} 360$	$\tau{=}480$	
$\tau = 120$	40.4	82.1	123.1	163.4	
$\tau = 240$	41.1	87.5	130.0	165.5	
$\tau{=}360$	45.7	90.2	125.7	163.5	

(a) Effects of time duration used in train-

ing on cryoRL performance.

Table 5: Ablation study of CryoRL

(b) Effects of different rewards on cryoRL's performance.

41.1

43.0

41.6

41.8

Duration (minutes)

 $|\tau=120 \ \tau=240 \ \tau=360 \ \tau=480$

87.0

86.9

80.8

86.6 132.0 171.4

131.1 172.0

129.5 165.9

 $124.7 \quad 163.3$

Rewards

0.23 (default) 0.09 (default)

grid-level

0.09

 $0.09 (\times 2)$

 $0.09 (\times 2)$

square-level

 $0.23(\times 2)$

0.23

 $0.23 (\times 2)$

sparsely inspected, with a small portion of holes visited by the users. In other words, there are a lot of holes in the patch images without a CTF ground truth available. As a result, the limited coverage in M data is not sufficient for learning effective RL policies for planning microscope movement. Nevertheless, M data were collected under different realistic settings where various grid types and microscopes were used. It is much more diverse and substantially larger than Y1 data (over 100,000 holes with CTF ground truth in M vs. 4,000 in Y), making it suitable for building a foundation model for hole classification.

We split M data into training and validation sets at a ratio of 4:1 and trained a hole classifier based on Resnet50. We then applied the classifier to both test sets in Y1 and Y2, and the results are listed in Table 4. As seen from the table, the classifier achieves moderate performance on Y1 and Y2, with an accuracy of around 70% in low-CTF classification, suggesting that hole classification is still a challenging problem that needs further improvement.

We further trained RL models on Y1 and Y2 using the classification results based on the M model mentioned above. As shown in Table 4, a modest classifier (M) results in a performance drop in CryoRL (4^{th} row) as expected, but the results are still reasonably good. Additionally, we extend to test the transferability of the RL models. Specifically, we applied the RL model based on Y1 to Y2 dataset and compared the results (6^{th} row) to those from the RL model trained on Y2 itself (5^{th} row). Even though Y1 and Y2 datasets were collected with different samples and grid types, the results between these two models are still comparable, showing the good transferability of CryoRL.

5.3 Ablation Study

In this section, we conduct experiments to characterize our proposed approach. We investigate how hole time duration and rewarding affects the performance of cryoRL (i.e. the total number of low-CTF holes found ain a given amount of time). We also provide visualization of a planned path by CryoRL and the learned RL polices.

Effects of Time Duration In principle, the time duration τ used in training cryoRL controls the degree of interaction of the RL agent with the data. A small τ limits cryoRL to a few high-quality patches only, which might result in a more



Fig. 9: A trajectory of microscope movement planned by CryoRL at square level (left) and patch level (right), respectively, in a 8-hour data collection session. The blue and yellow boxes show part of the training and validation sets while the color bar represents the ground-truth CTF value. The trajectory within a specific patch (right) illustrates that cryoRL can identify patches with more good holes (CTF \leq 6.0) in a global sense and prioritize their visits first. It is also noticed that some patches with a few good holes are left untouched in the square. This is because moving to a patch in another square (not shown here) is more rewarding than staying.

conservative policy that underfits. Table 5a confirms this potential issue, showing inferior performance when a short duration of 120 minutes is used for training. **Effects of Rewarding Strategies.** In our approach, the rewards used in policy learning are empirically determined. To check the potential impact of different rewards on the performance of cryoRL, we trained more Q networks by doubling the reward for a) square switching; b) grid switching; and c) both. These changes are intended to encourage more active exploration of the data. As shown in Table 5b, the different rewarding schemes perform comparably, and increasing the reward for square switching leads to slightly better performance than the default setting. This suggests that CryoRL is not sensitive to rewards setup as long as the rewards for better performance of cryoRL is an area of improvement in future work.

Trajectory Path and RL Policy Visualization. We plot one trajectory path of the microscope movement on the atlas planned by our CryoRL at square level (left) and patch level (right), respectively, in a 8-hour data collection session. The trajectory within a specific patch (right) illustrates that cryoRL can identify patches with more good holes (CTF \leq 6.0) in a global sense and prioritize their visits first. It is also noticed that some patches with a few good holes are left untouched in the square. This is because moving to a patch in another square (not shown here) is more rewarding than staying.

We further compare and visualize the policies learned by our approach as well as the strategies used by human users. Specifically, we count how often the microscope visits a pair of hole-level images (i.e patches) in the 50 trials of our results and illustrate such information by an undirected graph. A node of the graph represents a patch and a blue edge between two patches indicates the



Fig. 10: Illustration of data collection policies from cryoRL and human subjects. Here a graph node denotes a patch in our data and the size of the node indicates the quality of the patch (i.e the number of low-CTF holes). Patches from the same grid are grouped by color and linked by light grey edges. A blue edge between a pair of patches shows how often the two patches are visited by the microscope. Intuitively, an effective policy should demonstrate strong connections between large-sized nodes, which is the case for the learned policy by our approach. As opposed to the RL policy, the human users presents random behaviors (b)).

frequency of them being visited by the microscope. Note that the node size here denotes the quality of a patch determined by the number of good holes in the patch, and the node color indicates the grid the patch belongs to. Intuitively, a good policy should show strong connections between large-sized nodes. As observed in Fig 10a), our learned RL policy favors larger-size nodes, clearly demonstrating that CryoRL enables efficient data collection. Oppositely, the behavior of human users is random, with a lot of more patches being explored. This is because that the users were not penalized for switching different patches in the human study, and may also be due to the large variance in the user expertise.

6 Conclusion

To summarize, by combining supervised classification and deep RL, cryoRL provides a new framework for cryo-EM data collection. It can not only return the quality predictions for lower magnified hole level images but can also plan the trajectory for data acquisition. We have shown that cryoRL combined with an offline hole classifier achieves better performance than average human users. Nevertheless, cryoRL needs squares to be pre-selected and all their corresponding patch-level images to be pre-captured. Future work will be needed to further optimize the RL system to consider more of this hierarchical process of cryo-EM data collection. The specific hyper-parameters, especially the penalties in the reward function, can also be improved for a more practical application.

References

- Satyajeet Singh Ahuja Yuandong Tian Ying Zhang Xin Jin ang Zhu, Varun Gupta. Network Planning with Deep Reinforcement Learning. In SIGCOMM. ACM, 2021.
- Philip R Baldwin, Yong Zi Tan, Edward T Eng, William J Rice, Alex J Noble, Carl J Negro, Michael A Cianfrocco, Clinton S Potter, and Bridget Carragher. Big data in cryoem: automated collection, processing and accessibility of em data. *Current opinion in microbiology*, 43:1–8, 2018.
- Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *International Conference on Machine Learning*, pages 449–458. PMLR, 2017.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. arXiv preprint arXiv:1912.06680, 2019.
- 5. Anders Brahme. Comprehensive biomedical physics. Newnes, 2014.
- Jack PK Bravo, Tyler L Dangerfield, David W Taylor, and Kenneth A Johnson. Remdesivir is a delayed translocation inhibitor of sars-cov-2 replication. *Molecular cell*, 81(7):1548–1552, 2021.
- Anchi Cheng, Edward T Eng, Lambertus Alink, William J Rice, Kelsey D Jordan, Laura Y Kim, Clinton S Potter, and Bridget Carragher. High resolution single particle cryo-electron microscopy using beam-image shift. *Journal of structural biology*, 204(2):270–275, 2018.
- 8. Ahmed Fawzy Gad. Pygad: An intuitive genetic algorithm python library, 2021.
- 9. Eric Hand. 'we need a people's cryo-em.' scientists hope to bring revolutionary microscope to the masses. *Science*, 2020.
- Anssi Kanervisto, Christian Scheller, and Ville Hautamäki. Action space shaping in deep reinforcement learning. In 2020 IEEE Conference on Games (CoG), pages 479–486. IEEE, 2020.
- Zunlong Ke, Joaquin Oton, Kun Qu, Mirko Cortese, Vojtech Zila, Lesley McKeane, Takanori Nakane, Jasenko Zivanov, Christopher J Neufeldt, Berati Cerikan, et al. Structures and distributions of sars-cov-2 spike proteins on intact virions. *Nature*, 588(7838):498–502, 2020.
- 12. Scott Kirkpatrick. Optimization by simulated annealing: Quantitative studies. Journal of statistical physics, 34(5):975–986, 1984.
- Goran Kokic, Hauke S Hillen, Dimitry Tegunov, Christian Dienemann, Florian Seitz, Jana Schmitzova, Lucas Farnung, Aaron Siewert, Claudia Höbartner, and Patrick Cramer. Mechanism of sars-cov-2 polymerase stalling by remdesivir. *Nature* communications, 12(1):1–7, 2021.
- Thomas Kollar and Nicholas Roy. Trajectory optimization using reinforcement learning for map exploration. The International Journal of Robotics Research, 27(2):175–196, 2008.
- Alexander Krull, Peter Hirsch, Carsten Rother, Augustin Schiffrin, and C Krull. Artificial-intelligence-driven scanning probe microscopy. *Communications Physics*, 3(1):1–8, 2020.
- Gabriel C Lander, Scott M Stagg, Neil R Voss, Anchi Cheng, Denis Fellmann, James Pulokas, Craig Yoshioka, Christopher Irving, Anke Mulder, Pick-Wei Lau, et al. Appion: an integrated, database-driven pipeline to facilitate em image processing. *Journal of structural biology*, 166(1):95–102, 2009.

- 18 Q. Fan et al.
- Alessandro Lazaric, Marcello Restelli, and Andrea Bonarini. Reinforcement learning in continuous action spaces through sequential monte carlo methods. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2008.
- Ngan Le, Vidhiwar Singh Rathour, Kashu Yamazaki, Khoa Luu, and Marios Savvides. Deep reinforcement learning in computer vision: A comprehensive survey. *CoRR*, abs/2108.11510, 2021.
- 19. Kyowoon Lee, Sol-A Kim, Jaesik Choi, and Seong-Whan Lee. Deep reinforcement learning in continuous action spaces: a case study in the game of simulated curling. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2937–2946. PMLR, 10–15 Jul 2018.
- Florian A Lempp, Leah B Soriaga, Martin Montiel-Ruiz, Fabio Benigni, Julia Noack, Young-Jun Park, Siro Bianchi, Alexandra C Walls, John E Bowen, Jiayi Zhou, et al. Lectins enhance sars-cov-2 infection and influence neutralizing antibodies. *Nature*, 598(7880):342–347, 2021.
- Yilai Li, Jennifer N Cash, John JG Tesmer, and Michael A Cianfrocco. Highthroughput cryo-em enabled by user-free preprocessing routines. *Structure*, 28(7):858–869, 2020.
- 22. Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning, 2013.
- Mohammadreza Nazari, Afshin Oroojlooy, Lawrence Snyder, and Martin Takác. Reinforcement learning for solving the vehicle routing problem. Advances in neural information processing systems, 31, 2018.
- O. Klimov A. Nichol M. Plappert A. Radford J. Schulman S. Sidor Y. Wu P. Dhariwal, C. Hesse and P. Zhokhov. Openai baselines. 2017.
- 25. Matthew Perry. Simanneal: https://github.com/perrygeo/simanneal, 2020.
- 26. Alexis Rohou and Nikolaus Grigorieff. Ctffind4: Fast and accurate defocus estimation from electron micrographs. *Journal of structural biology*, 192(2):216–221, 2015.
- Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. arXiv preprint arXiv:1511.05952, 2015.
- Johannes F Scheid, Christopher O Barnes, Basak Eraslan, Andrew Hudak, Jennifer R Keeffe, Lisa A Cosimi, Eric M Brown, Frauke Muecksch, Yiska Weisblum, Shuting Zhang, et al. B cell genomics behind cross-neutralization of sars-cov-2 variants and sars-cov. *Cell*, 2021.
- 29. David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of Go with deep neural networks and tree search. Nature, 529(7587):484–489, jan 2016.
- Markus Stabrin, Fabian Schoenfeld, Thorsten Wagner, Sabrina Pospich, Christos Gatsogiannis, and Stefan Raunser. Transphire: automated and feedback-optimized on-the-fly processing for cryo-em. *Nature communications*, 11(1):1–14, 2020.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018.
- 32. Dimitry Tegunov and Patrick Cramer. Real-time cryo-electron microscopy data preprocessing with warp. *Nature methods*, 16(11):1146–1152, 2019.

CryoRL: Reinforcement Learning Enables Efficient Cryo-EM Data Collection

- 33. Peng Tong, Juan Liu, Xijun Wang, Bo Bai, and Huaiyu Dai. Deep reinforcement learning for efficient data collection in uav-aided internet of things. In 2020 IEEE International Conference on Communications Workshops (ICC Workshops), pages 1-6. IEEE, 2020.
- Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In *Interna*tional conference on machine learning, pages 1995–2003. PMLR, 2016.
- Felix Weis and Wim JH Hagen. Combining high throughput and high quality for cryo-electron microscopy data collection. Acta Crystallographica Section D: Structural Biology, 76(8):724–728, 2020.
- Thomas Weise. Global optimization algorithms-theory and application. Self-Published Thomas Weise, 361, 2009.
- Jiayi Weng, Huayu Chen, Dong Yan, Kaichao You, Alexis Duburcq, Minghao Zhang, Hang Su, and Jun Zhu. Tianshou: A highly modularized deep reinforcement learning library. arXiv preprint arXiv:2107.14171, 2021.
- Daniel Wrapp, Nianshuang Wang, Kizzmekia S Corbett, Jory A Goldsmith, Ching-Lin Hsieh, Olubukola Abiona, Barney S Graham, and Jason S McLellan. Cryoem structure of the 2019-ncov spike in the prefusion conformation. *Science*, 367(6483):1260–1263, 2020.
- Hangping Yao, Yutong Song, Yong Chen, Nanping Wu, Jialu Xu, Chujie Sun, Jiaxing Zhang, Tianhao Weng, Zheyuan Zhang, Zhigang Wu, et al. Molecular architecture of the sars-cov-2 virus. *Cell*, 183(3):730–738, 2020.
- 40. Deheng Ye, Zhao Liu, Mingfei Sun, Bei Shi, Peilin Zhao, Hao Wu, Hongsheng Yu, Shaojie Yang, Xipeng Wu, Qingwei Guo, et al. Mastering complex control in moba games with deep reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 6672–6679, 2020.
- 41. Wanchao Yin, Chunyou Mao, Xiaodong Luan, Dan-Dan Shen, Qingya Shen, Haixia Su, Xiaoxi Wang, Fulai Zhou, Wenfeng Zhao, Minqi Gao, et al. Structural basis for inhibition of the rna-dependent rna polymerase from sars-cov-2 by remdesivir. *Science*, 368(6498):1499–1504, 2020.
- 42. Tom Zahavy, Matan Haroush, Nadav Merlis, Daniel J Mankowitz, and Shie Mannor. Learn what not to learn: Action elimination with deep reinforcement learning. Advances in Neural Information Processing Systems, 31, 2018.
- 43. Yu Zhang, Zhiyu Mou, Feifei Gao, Ling Xing, Jing Jiang, and Zhu Han. Hierarchical deep reinforcement learning for backscattering data collection with multiple uavs. *IEEE Internet of Things Journal*, 8(5):3786–3800, 2020.
- 44. Xiangyu Zhao, Liang Zhang, Zhuoye Ding, Long Xia, Jiliang Tang, and Dawei Yin. Recommendations with Negative Feedback via Pairwise Deep Reinforcement Learning. In Proceedings of the 24th International Conference on Knowledge Discovery and Data Mining, pages 1040–1048. ACM, 2018.

A Appendix

Methods	classifier	$\tau = 120$	$\tau{=}240$	$\tau {=} 360$	$\tau = 480$
CryoRL-A2C	1	37.0 ± 6.7	71.9 ± 9.5	104.1 ± 9.4	$144.8 {\pm} 9.3$
CryoRL-C51		37.2 ± 4.2	70.6 ± 5.0	98.1 ± 5.0	128.0 ± 3.6
CryoRL-DQN	Resnet18	42.9 ± 3.6	80.8 ± 3.0	123.3 ± 5.9	168.5 ± 2.0
CryoRL-DQN (dueling)		42.9 ± 4.2	86.9 ± 5.2	125.2 ± 5.3	159.5 ± 6.9
CryoRL-DQN (prioritized)		42.3 ± 4.1	$86.0 {\pm} 4.3$	128.3 ± 3.6	174.1 ± 5.5
$CryoRL-A2C^{\dagger}$	1	46.0±2.6(+24.3%)	86.4±1.2(+20.2%)	124.4±2.2(+19.5%)	158.8±4.5(+9.7%)
CryoRL-C51 [†]		$46.5 \pm 0.8 (+25.0\%)$	$78.1 \pm 1.3 (+10.6\%)$	$116.7 \pm 1.0 (+18.9\%)$	$138.2 \pm 2.8 (+8.0\%)$
$CryoRL-DQN^{\dagger}$	Resnet18	47.4±2.0(+10.5%)	91.0±2.5(+12.6%)	$132.8 \pm 2.1 (+7.7\%)$	$176.5 \pm 3.5 (+4.7\%)$
$CryoRL-DQN^{\dagger}$ (dueling)		$47.2 \pm 1.0 (+10.0\%)$	$89.1 \pm 2.7 (+10.0\%)$	$129.2 \pm 1.8(+3.2\%)$	$166.2 \pm 5.0 (+4.2\%)$
CryoRL-DQN ^{\dagger} (prioritized)		47.1±2.3(+11.3%)	$90.4 \pm 2.3 (+2.5\%)$	133.0±3.0(+3.7%)	177.4±4.1(+1.9%)
CryoRL-DQN	1	41.7 ± 3.1	86.6±3.0	132.0 ± 2.3	171.4 ± 2.0
$CryoRL-DQN^{\dagger}$	Resnet50	47.4±0.5(+13.7%)	89.0±3.1(+5.1%)	131.8±1.8(+0.0%)	172.6±2.0(+1.0%)

Table 6: Performance of different CryoRL variants on the Y1 dataset using Resnet18 as the offline hole classifier († indicates action elimination.). The performance gains from action elimination are highlighted by numbers in parentheses. The numbers in bold mark the best performance achieved by CryoRL under different time durations using Resnet18 as the classifier.

Resnet18 Results. We adopted Resnet18 as the offline classifier for CryoRL, which achieves better low-CTF classification accuracy than Resnet50 (91.0% v.s 83.9%), but lower high-CTF classification accuracy (87.5% v.s 91.2%). This suggests that Resnet18 yield more falsely classified good holes. As a result, CryoRL based on Resnet18 underperforms its counterpart based on Resnet50 (Table 6). However, when action elimination is applied, the performance of Resnet18 is significantly boosted and even gets slightly better than that of Resnet50. Additionally, action elimination greatly improves A2C and C51, similar to what's shown in the main paper.

Require: States S, Actions A, Rewards R1: procedure Action_ELIM (P, L, C, β, τ) **Require:** Learning Rate α , Discounting factor γ , 2: $N_{max} \leftarrow \beta * max_lCTF(P,C,\tau) \triangleright max-$ imum lCTFs found assuming that all holes Elimination coefficient beta**Require:** Switching costs C, Duration τ are good 1: procedure QLEARNING $AE(S, A, R, C, \alpha, \beta, \gamma, 3)$ $n \leftarrow 0$ $A' \leftarrow \{\}$ 4: 2. $P \leftarrow [p_0, p_1, \cdots, p_n]$ \triangleright Patches for p_i in P do 5: 3: $L \leftarrow [l_0, l_1, \cdots, l_n] \triangleright \#$ of predicted lCTFs in $n \leftarrow n + l_i$ $A' \leftarrow A' \bigcup \{h_j \in p_i | j = 1 \cdots m_i\}$ if $n \ge N_{max}$ then here h 6: each patch 7: 4: $A^{\overline{\prime}}$ \leftarrow Action $Elim(P, L, C, \tau)$ 8: \triangleright standard $\overset{\circlearrowright}{9}$: 5: $Q \leftarrow QLearning(S, A', R, \alpha, \gamma)$ breakQ learning 10: end if return O 11: end for 6: 7: 8: end procedure return A'12: end procedure

Alg. 7: Fast CryoRL with Action Elimination

Algorithm for Action Elimination. The psudo code for action elimination is illustrated in Alg. 7. In the algorithm, Action Elim returns a list of valid actions, which are provided to the standard QLearning procedure or other policy learners for policy learning. The procedure max lCTF finds an upper limit of the number of low-CTF holes within a time duration τ under the assumption that all holes are in good quality. The elimination coefficient β controls the size of the valid action set. During training, β should be set large to ensure sufficient training data with diversity. However, in test, β can be set smaller to eliminate bad microscope movements while making action execution efficient.

Experimental Setup for Genetic Algorithm (GA) and Simulated Annealing (SA) As mentioned in the main paper (Section 5.2), the solutions of both GA and SA are assessed based on the same objective function used for RL, i.e Eq. 1 in the main paper. We implemented CryoRL-GA based on pyGAD [8] and Cryo-SA base on SimAnneal [25]. For CryoRL-GA, we set the number of generations to 40 and the solutions per population to 10. We use single-point crossover and and random mutation. For CryoRL-SA, the minimum and maximum temperatures are chosen as 1e - 8 and \sqrt{N} , respectively, where N is the total number of training samples. The temperature reduction rate is set to 0.995.