# Probing the Theoretical and Computational Limits of Dissipative Design

Shriram Chennakesavalu and Grant M. Rotskoff[*]
*Department of Chemistry, Stanford University*

Self-assembly, the process by which interacting components form well-defined and often intricate structures, is typically thought of as a spontaneous process arising from equilibrium dynamics. When a system is driven by external *nonequilibrium* forces, states statistically inaccessible to the equilibrium dynamics can arise, a process sometimes termed direct self-assembly. However, if we fix a given target state and a set of external control variables, it is not well-understood i) how to design a protocol to drive the system towards the desired state nor ii) the energetic cost of persistently perturbing the stationary distribution. Here we derive a bound that relates the proximity to the chosen target with the dissipation associated with the external drive, showing that high-dimensional external control can guide systems towards target distribution but with an inevitable entropic cost. Remarkably, the bound holds arbitrarily far from equilibrium. Secondly, we investigate the performance of deep reinforcement learning algorithms and provide evidence for the realizability of complex protocols that stabilize otherwise inaccessible states of matter.

## I. INTRODUCTION

Designing molecular materials that robustly and autonomously assemble into specific, targeted mesoscale structures remains a central challenge in a variety of fields, from materials science [1–4] to biology [5–7]. The canonical approach to this design problem is to engineer components with specific molecular interactions that stabilize a thermodynamic ground state corresponding to the target [8], an inverse approach pioneered in Refs. [9–11]. Advances in tuneable materials, such as patchy particles [12], DNA coated colloids [13], and DNA origami [14] have enabled the realization of highly-specific, directional interactions among the constituent elements and subsequent efforts to optimize these interactions [15–17]. However, this paradigm is by no means general—in fact, in many instances, self-assembly is driven not by unique and addressable [18] interactions but rather by weak and nonspecific ones. Here we explore a distinct approach, one based on *nonequilibrium* external control of a dynamically assembling system, as opposed to designing interactions among the components of a system.

Tailoring self-assembly to produce materials with exotic or desirable properties has long been a goal in the molecular sciences [8]. Early computational work on this topic focused on the inverse design of interaction potentials that produced materials with, for example, target radial distribution functions, densities, and band gaps [9]. Parameterizing flexible interaction potentials that consistently achieve the desired states is nevertheless challenging. The widespread adoption of machine learning techniques in scientific computing has led many to revisit this problem using more sophisticated numerical representations of optimizable potential energy functions. These techniques have led to significant advances, allowing both for more complex interactions [19] and new computational approaches [20], increasing the set of structures accessible to designer self-assembly.

At molecular scales, it is often impossible to alter the nature of an interaction without fundamentally changing the molecular components as well, potentially disrupting biological or chemical function. Rather than viewing self-assembly as a thermodynamic process in which the ultimate structure is determined by a minimum of the free energy, we examine the *dynamics* of assembly trajectories [8] and ask whether an external agent can perturb the assembling components so as to maximize the yield of a target structure. Because the magnitude of fluctuations at the nanoscale is comparable to the size of the system itself, the task of determining an effective external protocol for control resembles a stochastic optimal control problem.

* rotskoff@stanford.edu

Just as interaction design is inherently limited by the constituent materials, the precision of control is dictated by the external fields that couple to a given system and the spatio-temporal resolution with which we can reasonably alter these fields. Moreover, an external control approach to directed self-assembly presents new computational challenges: stochastic optimal control problems are typically formulated as high-dimensional partial differential equations, which cannot be solved either analytically or numerically for nontrivial systems. Here, we instead pose the design problem as the optimization of a Markov decision process [21], which is in turn amenable to deep reinforcement learning algorithms. Of course, these complicated high-dimensional problems have also benefited from advances in deep learning, enabling the optimization of very high-dimensional feedback protocols for essentially arbitrary physical systems.

In this paper, we investigate the theoretical and computational limits of nonequilibrium control in the context of two minimal models of molecular self-assembly. We establish theoretically a relationship between the dissipative cost of a protocol and the fidelity with which a target structure can be produced, akin to bounds that have been established for nonequilibrium growth processes [22]. We then explore the capabilities of deep reinforcement learning algorithms to control the quenched cluster size distribution of a system of particles using a feedback thermal annealing protocol as well as the steady-state cluster size distribution of nonequilibrium actively driven colloids. Taken together, our theoretical and computational results emphasize that high dimensional control can target assembly outcomes with high precision but with an inescapable dissipative cost.

## II. THE DISSIPATIVE COST OF HIGH-FIDELITY CONTROL

Consider a physical system with coordinates $\boldsymbol{x} \in \Omega \subset \mathbb{T}^d$ evolving according to overdamped Langevin equation

$$d\boldsymbol{X}_t = b(\boldsymbol{X}_t)dt + \sqrt{2D}d\boldsymbol{W}_t \tag{1}$$

where $b$ is a nonequilibrium drift, $D = k_{\mathrm{B}}T/\mu$ is the diffusion coefficient, and $\boldsymbol{W}_t$ is a Weiner process in $\mathbb{R}^d$. We assume that $\boldsymbol{X}_t$ is ergodic so that there exists a unique stationary probability density $\rho_{\mathrm{ss}} : \Omega \to \mathbb{R}$.

Our goal is to develop a feedback-guided, external control protocol $u$ that pushes the steady state distribution towards a specified target. At present, we focus on external driving that can be represented as a force, not a noise term, though we consider both regimes in the subsequent numerical experiments. We assume that the external control can be implemented as a spatially-varying external force $u$ leading to the controlled SDE

$$d\boldsymbol{X}_t^u = [b(\boldsymbol{X}_t^u) + u(\boldsymbol{X}_t^u)]dt + \sqrt{2D}d\boldsymbol{W}_t, \tag{2}$$

which in turn has an associated a steady state density $\rho_u$. While the most generic design task requires tuning $u$ to coincide with a target steady state distribution $\rho_*$, it is not clear how to specify a target density function for a large interacting particle system, as we typically characterize these systems instead by some low-dimensional observable. This more limited description requires setting some target average value of a given observable $f : \Omega \to \mathbb{R}$. Let us denote the target value of $f$ by $f_*$. The optimal controller then solves the minimization problem

$$u_* = \underset{u}{\mathrm{argmin}} |\mathbb{E}_u f - f_*| \tag{3}$$

where $\mathbb{E}_u$ denotes the expectation over the controlled process (2) and $\mathbb{E}_* f \equiv f_*$ denotes the target value of $f$, which we view as the expectation over the unknown target distribution (cf. Appendix A for a detailed discussion).

In some cases, the chosen observable $f$ might not be informative about the system. However, while we seek to carry out this minimization for a particular choice of $f$, if we instead allow $f$ to vary and solve the minimax problem to find the controller that minimizes the mean discrepancy

over all functions $g$[23], then the objective is the Kantorovich-Rubenstein dual formulation of the Wasserstein-1 distance [24]. This metric quantifies the distance between the target distribution and the steady state distribution of the controlled process,

$$\mathcal{W}_1(\rho_u, \rho_*) = \max_g \min_u |\mathbb{E}_u g - \mathbb{E}_* g|. \tag{4}$$

The Wasserstein distance is an optimal transport distance on probability distributions, measuring the cost to reallocate mass from one distribution to another. Applications of optimal transport distances have become widespread in data analysis and machine learning. See Peyré and Cuturi [25] for an applied perspective. For the fixed observable of interest, the mean discrepancy is bounded above by the Wasserstein distance, which in turn, is bounded by the Kullback-Liebler divergence or relative entropy

$$\min_u |\mathbb{E}_u f - f_*| \le \mathcal{W}_1(\rho_u, \rho_*),$$
$$\le C\sqrt{2D_{\mathrm{KL}}(\rho_u, \rho_*)}. \tag{5}$$

The first inequality follows from an application of dual formulation of the total variation distance and subsequently an application of Pinsker's inequality. Interestingly, this upper bound has a direct physical interpretation in terms of the dissipative cost of controlling the trajectory to alter the steady state distribution.

The Kullback-Leibler divergence between the stationary distribution of the controlled process and the target distribution can be interpreted as a nonequilibrium free energy difference [26] between the two distributions. Maintaining a nonequilibrium steady state distinct from the steady state of the unguided system incurs a time-extensive entropic cost. For the path measures associated with the dynamics, we can express this quantity explicitly via the Girsanov theorem, as shown in Appendix A. In particular, we see that the entropic cost of control per unit time can be written as a time-average, quantifying the cost of designing a dissipative steady state. This design cost can be measured by computing $D_{\mathrm{KL}}(\rho_u \| \rho_{\mathrm{ss}})$. For the protocol $u$, we show that the cost of reducing the Kullback-Liebler divergence in (5) to minimize (3) is

$$\sigma_{\mathrm{ex}} = \frac{1}{T} \int_0^T u(X_t^u)^2 dt \ge 0; \tag{6}$$

that is, there is an inverse relationship between the fidelity of control and the non-negative cost of control. The optimal rate is given using $u_* = b_* - b$, which is in general unknown and which vanishes when no control is needed. We can still examine the relationship indirectly by monitoring the total entropy production of the controlled system. This relation holds when the original, controlled, and target dynamics are all arbitrarily far from equilibrium.

The relation (5) helps quantify the dissipative cost of controlling a system, even when the state of the system is only partially observed (i.e., through an observable). Importantly, the relation that we derive is only an upper bound—the discrepancy could be small for some uninformative choice of $f$ without any significant perturbation. Similarly, many dissipative processes will not necessarily reduce the discrepancy between the mean and its target value—that is, the control could be imprecise. Nevertheless, when this bound is tight, the KL divergence quantifies the physical cost of driving the system towards the externally specified target.

## III. REINFORCEMENT LEARNING FOR HIGH-DIMENSIONAL CONTROL PROTOCOLS

*A priori*, finding the optimal $u$ requires both detailed information about the dynamics of the system and the ability to implement complex interactions. Fortunately, the minimization problem (3) is also amenable to model-free reinforcement learning, an approach we explore here. Moreover, the
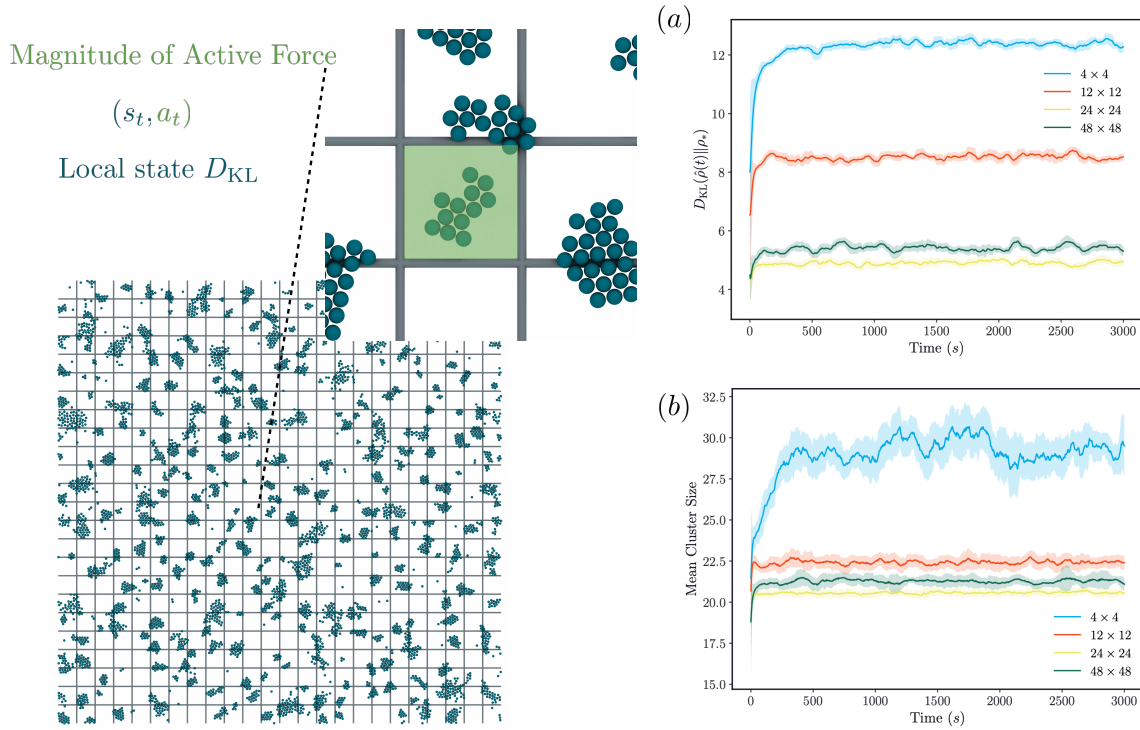
FIG. 1. Activity based control of transient motility induced phase separated clusters. Left: the state-action pair of the Markov decision process is depicted graphically. The activity is modulated externally with spatial resolution indicated by the grid. The state is mean cluster size within a region. ($a$) The KL divergence from instantaneous cluster size distribution to the target distribution as a function of time for optimized protocols with increasing control resolution: $4 \times 4$, $12 \times 12$, $24 \times 24$, and $48 \times 48$. ($b$) The mean value (target 20) plotted as a function of time for the optimized protocols.

infimum over protocols $\boldsymbol{u} : \mathbb{R}^d \to \mathbb{R}^d$ requires minimizing with respect to arbitrary many-body force functions which are unlikely to be realizable in experimental settings. We pursue a more pragmatic approach by building the constraints of control directly into our protocol $\boldsymbol{u}$; these constraints are imposed by restricting to a fixed class of experimentally accessible protocols.

We drive the system so that the observable of interest, a cluster size distribution throughout this paper, matches an externally specified target. We defined the observable by first introducing a map $h : \mathbb{R}^d \to \mathbb{N}^n$ which counts the number of clusters of size $k$ and stores that number in the $k$th entry of an $n$-dimensional vector where $n$ is the total number of particles in the system and hence the maximum cluster size. We then define the normalized histogram of cluster sizes for the configuration $\hat{\rho}(h(\boldsymbol{x}))$ and measure the discrepancy between this histogram and the target using the Kullback-Leibler divergence,

$$C(\boldsymbol{x}) = D_{\mathrm{KL}}(\hat{\rho}(h(\boldsymbol{x})) \| \rho_*). \tag{7}$$

This empirical distribution implicitly depends on the external controller $\boldsymbol{u}$ through the sampled state $\boldsymbol{x}$. Importantly, this cost functional makes evaluation of the loss function independent of the particle dynamics, which in turn allows for model free reinforcement learning. Using a cluster size distribution rather than the average cluster size ensures that the observable robustly describes the target state even when the region of control is large.

Because the objective (7) does not depend explicitly on an unspecified path measure, it is a tractable target for optimization. We consider control functions $\boldsymbol{u}$, depicted schematically in Fig. 1,

in which the external control drives a system locally with a fixed spatial patterning. While the protocol is not an arbitrary many-body force, in our case it remains high-dimensional and $u$ may be a complicated function of the configuration $x$. Despite this complexity, neural networks offer a robust, high-dimensional function representation, which we exploit in our representation of $u$. The steady state $\rho_u$ distribution depends on the dynamics of the system, so direct, gradient-based optimization of (7) is challenging. In our setup, the duration of the period between protocol updates is sufficiently long that the gradients of the control parameters become too small to meaningfully optimize the objective by backpropagating through the dynamics. There is also a conceptual reason for employing a model-free optimization algorithm—in experimental settings, we cannot require precise knowledge of the microscopic dynamics of the system to design an external protocol.

Because the optimal $u$ is time-independent by construction, we use a time-local representation of the joint dynamics of the system and the controller, also known as a Markov decision process. Optimization problems of this type are the basic framework for reinforcement learning and have been studied extensively in the machine learning and control literature [21]. Within this framework, maintaining the optimal steady state distribution requires incorporating information about the expected future divergence from the target distribution, as measured by $C$. For a fixed protocol $u$, the cumulative expected future divergence from the target distribution is an expectation over the dynamics of the system, starting from a given state $x_t$, is

$$\bar{C}(x_t, u) = \mathbb{E}_u^{x_t} \sum_{k=0}^{\infty} \gamma^k C(x_{t+\tau(k+1)}) \tag{8}$$

where $\gamma < 1$ is a so-called discount factor that ensures that the sum is convergent and gives additional weight to temporally proximate states. Some reinforcement learning algorithms, e.g. policy gradient [21], use (8) as a direct target for optimization, but the lack of time locality makes gradient based optimization challenging when the dynamics occurs over long time scales [27].

Deep reinforcement learning algorithms based on $Q$-learning lift the expected cost $\bar{C}$ so that it depends on a given state-action pair $(x_t, a_t)$ and protocol $u$. We assume that the system evolves with $u(x_s) = a_t$ constant for a fixed duration $s \in [t, t + \tau]$ so that $x_{t+\tau}$ is obtained by solving (2) with the initial condition $x_t$; this means that the feedback protocol has a time lag and in practice we choose $\tau$ to be sufficiently short that changes in the cost were minimal. The lifted cost functional, conventionally called $Q$, quantifies the future cost assuming that the action $a_t$ is taken at time $t$, explicitly,

$$Q^u(x_t, a_t) = \mathbb{E}_u^{x_t} \left[ C(x_{t+\tau}; a_t) + \sum_{k=1}^{\infty} \gamma^k C(x_{t+\tau(k+1)}) \right]. \tag{9}$$

In (9), the expectation is carried out over a trajectory initialized at $x_t$ and subject to the action $a_t$ until time $\tau$ and subsequently using the protocol $u$. The optimal next action for a given protocol $u$ is then simply $\mathrm{argmin}_a Q^u(x_t, a)$. Employing $Q$-learning requires estimating $Q$ usually with a value iteration algorithm [28]; traditionally this has been carried out using a tabular representation of $Q$, meaning that every state-action pair must be visited by the dynamics in order to provide an accurate representation. Of course, if the state or action space is high-dimensional, this is infeasible. Recently, alternative strategies that represent $Q$ as a deep neural network have shown promise in a variety of contexts. Importantly, when $Q$ is represented as a neural network, high-dimensional state and action spaces become tractable. We employ a variant of deep $Q$-learning [29, 30] that uses two deep neural networks to approximate the $Q$ function, called double $Q$ learning which helps avoid minimization bias in the estimate of the minimizer; we discuss the details of this approach in Appendix B.
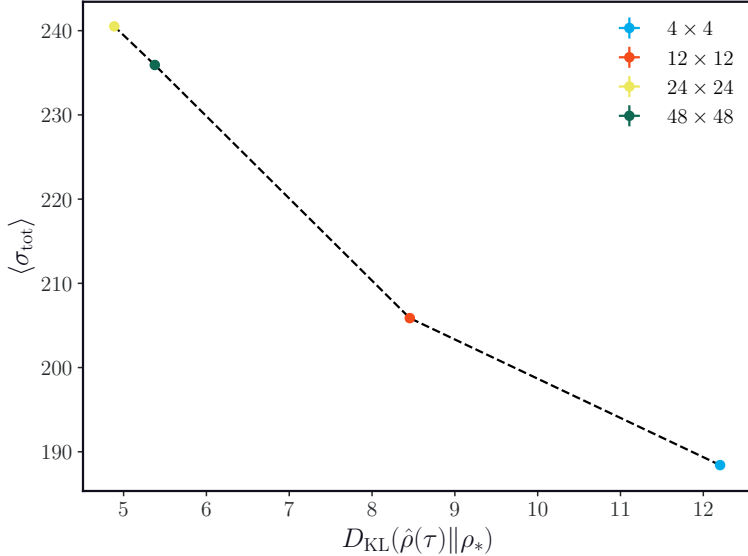
FIG. 2. The total entropy production of the controlled system plotted against the fidelity of control as measured by the KL divergence from the target.

## IV.   CONTROLLING CLUSTER SIZES IN ACTIVE COLLOIDS

To test this reinforcement learning approach to dissipative design, we first considered a model colloidal particle system actively driven by externally controlled light sources. Models of self-propelled or active matter evince rich phase behavior and have rapidly become canonical models for pattern formation out of equilibrium. In these systems, when the Péclet number is sufficiently large, the mobility depends strongly on the local density, which results in nonequilibrium phase separation [31]. This phenomenon, called motility induced phase separation (MIPS), requires energy consumption and is largely independent of the inter-particle interactions [31–33].

Because the activity can be externally modulated with simple controls, for example by selectively illuminating a portion of the system with light of variable intensity, it offers a natural example for control. Reinforcement learning has shown some success in controlling active systems: recently, Falk *et al.* [34] examined enhancing transport properties using low-dimensional external protocols optimized with an actor-critic model. The externally modulated activity leads to clustered states, but when this activity is turned off, the particles diffuse apart and the clusters disintegrate. The limitations of this control will necessarily limit the steady state distributions that can be accessed. In turn, this means that excess dissipation associated with the feedback protocol may enhance control, but it is possible that some of the energy expended goes to waste.

To explore the limits of an activity-inducing dissipative external control, we sought to maintain a steady state distribution consisting of clusters of particles much smaller than the macroscopic aggregate that forms when there is constant activity. We specified a target distribution $\rho_*$ of cluster sizes as in (7) using three distinct target distributions, all with identical mean and variance. The distribution of cluster sizes is discrete, so we first tested a binomial distribution with $n = 20, p = 3/4$. Because this distribution decay rapidly in the tails, we also tested a Gamma distribution $\Gamma(k, \theta)$ with shape parameter $k = 25$ and scale parameter $\theta = 3/5$ and computed the corresponding probability mass function by integrating the probability density over the bins. Finally, to remove any effects due to asymmetry of the target distribution, we used a Gaussian target distribution with $\mu = 20, \sigma^2 = 5$.

We used deep $Q$-learning to optimize an external protocol which controlled the intensity of the activity over a spatial grid, as depicted in Fig. 1. Both the cost function and the corresponding action were evaluated locally; that is, the cluster size distribution was computed over a given region and then the activity was chosen to minimize the estimated $Q$ function given the observed state. We tested this approach with grids of increasing resolution, corresponding to increasingly fine-tuned spatial control. The optimization was carried out for at least 20 "episodes" of 350 decisions, until the relative entropy between the instantaneous distribution and the target had converged.[35]. Once the protocol had converged, we computed the relative entropy and the mean cluster size over a collection of long test trajectories, shown in Fig. 1 (*a*) and (*b*). When the number of control regions is small (e.g, $4 \times 4$) perturbations to the activity are not localized enough to prevent the formation of large clusters. As a result, the mean cluster size is substantially larger than the target. On the other hand, the benefits of high-resolution control diminish as the number of regions becomes very large ($48 \times 48$). In effect, once a single region can accommodate only a few clusters of the target size, protocols of increasing resolution achieve the same outcome. Furthermore, these high-dimensional protocols come with an additional computational or practical burden.

Protocols with sufficient spatial resolution to execute local control perform well (Fig. 1), leading to mean values for the average cluster size that are close to the target. The variance of the empirical distribution of the steady state under the optimized protocol is larger than the target, emphasizing that our relatively coarse external control is still limited. The timescale over which a cluster diffusively disintegrates in the absence of activity exceeds the time required to form a cluster. Because the decision period has a fixed duration and activity will favor the formation of large clusters, there is a relative abundance of clusters with a size that exceeds the mean. These clusters also contribute to variance in the left tail of the distribution, because dissolving them creates a large number of free particles. Nevertheless, as shown in the schematic of Fig. 1, the typical clusters are in line with the target distribution and the formation of macroscopic clusters is always avoided.

We examine the generic relationship between dissipative cost and fidelity of control by computing the total entropy production for increasing resolution of control. In general we found that more control regions led to better control as measured by the cost function (7). Fig. 2 plots the total dissipation $\sigma_{\text{tot}}$ as a function of the KL-divergence from the target; each point is averaged over 100 realizations of the optimized feedback protocol. As the number of control regions increases from $4 \times 4$ up to $24 \times 24$, the fidelity increases but with a clear increase in dissipation, providing evidence of the utility of the bound (5). At the highest resolutions, diminishing returns become evident because the typical cluster size become comparable to the region itself—in this regime the cost function is essentially exactly the mean discrepancy in (3). For other target cluster size distributions (see Appendix C), we observed a similar trend. We emphasize that the generic trade-off between dissipation and fidelity is not necessarily a monotonic relationship, as observed for different cost functions in Appendix C because not all dissipative dynamics will lead to improvement of the cost.

## V. FEEDBACK GUIDED ANNEALING

While thermal annealing has a long history for macroscopic systems, repeated annealing cycles is an important part of the preparation of a wide variety of nanoscopic materials, from thin films [36] to DNA origami [37]. Annealing is limited as a mechanism for control because there is essentially only one parameter that can be tuned, the rate at which the temperature is decreased. This process will eventually find a global free energy minimum, meaning that the ultimate structure is determined entirely by the thermodynamic properties of the system. Indeed, there has been significant focus on designing interaction potentials that lead to specific structural motifs in a variety of contexts [9, 19, 20].

We examine an alternative paradigm that exercises more localized control with measurement-guided feedback to design an annealing schedule. Rather than globally tuning a temperature, we locally update the temperature on a grid, see Fig. 3. The specific temperature that the protocol
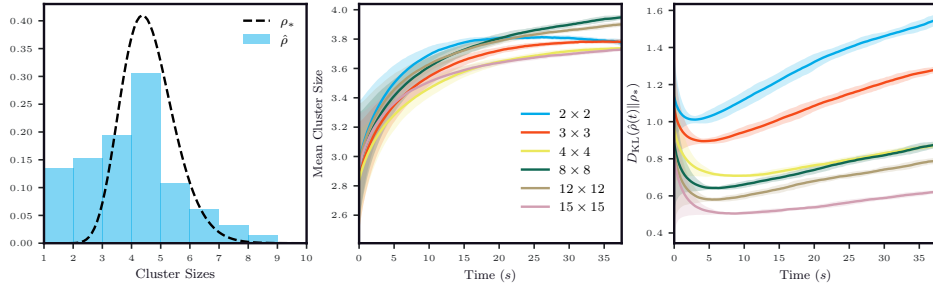
FIG. 3. (*a*) The target distribution of cluster sizes $\rho_*$ and the empirical distribution obtained using the optimized annealing protocol denoted $\hat{\rho}$. While there is an over-representation of isolated particles, as in the case of the active particle system, the model of the distribution and the tail are well-matched by using the protocol. (*b*) The evolution of the mean cluster size during the feedback annealing process for increasing resolution of temperature control. (*c*) The evolution of the discrepancy between the instantaneous distribution and the target as a function of time for increasing resolution of temperature control.

prescribes depends on the local configuration. While the annealing process itself is substantially more complicated in this framework, the approach could be used with arbitrary materials and does not require the realization of highly specific interactions, which can be enormously challenging to engineer in nanoscale systems.

To assess the prospects of this feedback guided annealing procedure, we studied a minimal example of cluster formation using a 2D Lennard-Jones system. At low temperatures, the free energy minimum corresponds to a single cluster, but at low densities, an instantaneous temperature quench from a high-temperature state will yield a kinetically trapped configuration that consists of small clusters.

We sought to control the distribution of these intermediate clusters by optimizing a thermal annealing function using reinforcement learning. We fixed a target cluster size distribution, chosen to be a Gamma distribution with variance $\sigma^2 = 1$ and a mean $\mu = 4$, shown as a dashed line in Fig. 3. We trained an external controller that used the local cluster size distribution to determine the subsequent temperature within each region of control. With only coarse control, the annealing reliably produced distributions of cluster sizes with a mean value close to the target (Fig. 3 (c)). However, high-resolution control was required to yield a distribution close—as measured by the KL-divergence—to the target distribution. The empirical cluster size distribution, averaged over 1000 annealing trajectories, is shown in Fig. 3 (a) for a $15 \times 15$ grid. While there is an over-representation of isolated particles by small amount, the distribution is remarkably close to the target. For fewer regions, the KL-divergence (Fig. 3 (d)) is substantially larger and the distributions (cf. Appendix D) differ markedly from the target.

Because the external control goes beyond the regime that we treat theoretically in Sec. II. While our analysis extends to systems in which the external control is represented as a drift function $u$, this system relies on changes to the diffusion tensor, which requires a substantially different mathematical treatment, a topic we plan to explore in future work.

## VI. CONCLUSION

Experimental advances enabling high-resolution external control create new opportunities to produce materials with exotic properties. In this work, we seeking to address several fundamental questions about the "realizability" of high-fidelity control protocols. In doing so, we derive a general bound that establishes a trade-off between the fidelity of control and the dissipative cost of

implementing it.

In many applications, the external fields that could be modulated with high-resolution will be *imprecise.* That is, we will not necessarily be able to tune a field directly conjugate to the observable of interest. We do not a priori know how to choose observables that robustly approximate the optimal transport distance, which is a significant topic for future work. Despite the arguably "coarse" control we have—at least compared with the case of interaction design—we nevertheless find that the optimization is successful at producing states with well-defined target mean cluster sizes. That is, our results demonstrate with appropriate objective functions, reinforcement learning algorithms can identify protocols that closely match target structures even without highly specific interactions. What is more, though we do not directly design the protocols to show the relation between dissipation and fidelity, we observe a trend consistent with the prediction throughout.

The approach we have pursued raises many additional questions. Foremost among these perhaps is the tightness of the bound (5) for a given observable, a question we hope to examine in model systems in future work. Many other reinforcement learning strategies could be deployed on the systems we studied, including approaches based on policy gradient. Some of these approaches may enable model-free control for more complicated systems, leading to insights about driven self-assembly in complex environments.

[1] Y. Yin and A. P. Alivisatos, Nature **437**, 664 (2005).
[2] Y. Ma and A. L. Ferguson, Soft Matter **15**, 8808 (2019).
[3] K. R. Gadelrab, A. F. Hannon, C. A. Ross, and A. Alexander-Katz, Molecular Systems Design & Engineering **2**, 539 (2017).
[4] H. Ronellenfitsch, N. Stoop, J. Yu, A. Forrow, and J. Dunkel, Physical Review Materials **3**, 095201 (2019).
[5] J. C. Cameron, S. C. Wilson, S. L. Bernstein, and C. A. Kerfeld, Cell **155**, 1131 (2013).
[6] C. Sigl, E. M. Willner, W. Engelen, J. A. Kretzmann, K. Sachenbacher, A. Liedl, F. Kolbe, F. Wilsch, S. A. Aghvami, U. Protzer, M. F. Hagan, S. Fraden, and H. Dietz, Nature Materials , 1 (2021).
[7] G. M. Rotskoff and P. L. Geissler, Proceedings of the National Academy of Sciences **112**, 201802499 (2018).
[8] E. M. Furst, Soft Matter **9**, 9039 (2013).
[9] M. C. Rechtsman, F. H. Stillinger, and S. Torquato, Physical Review Letters **95**, 228301 (2005).
[10] M. Rechtsman, F. Stillinger, and S. Torquato, Physical Review E **73**, 011406 (2006), number of pages: 12 Publisher: American Physical Society.
[11] M. C. Rechtsman, F. H. Stillinger, and S. Torquato, Physical Review E **74**, 021404 (2006), number of pages: 7 Publisher: American Physical Society.
[12] G.-R. Yi, D. J. Pine, and S. Sacanna, J. Phys.: Condens. Matter **25**, 193101 (2013).
[13] Y. Wang, Y. Wang, X. Zheng, É. Ducrot, J. S. Yodh, M. Weck, and D. J. Pine, Nat. Commun. **6**, 7253 (2015).
[14] Y. Ke, L. L. Ong, W. M. Shih, and P. Yin, Science **338**, 1177 (2012).
[15] D. Chen, G. Zhang, and S. Torquato, The Journal of Physical Chemistry B **122**, 8462 (2018), tex.eprint: https://doi.org/10.1021/acs.jpcb.8b05627.
[16] Z. Ma, E. Lomba, and S. Torquato, Physical Review Letters **125**, 068002 (2020), number of pages: 6 Publisher: American Physical Society.
[17] F. Romano, J. Russo, L. Kroc, and P. Šulc, Physical Review Letters **125**, 118003 (2020).
[18] L. O. Hedges, R. V. Mannige, and S. Whitelam, Soft Matter **10**, 6404 (2014).
[19] A. Das and D. T. Limmer, The Journal of Chemical Physics **154**, 014107 (2021).
[20] C. P. Goodrich, E. M. King, S. S. Schoenholz, E. D. Cubuk, and M. P. Brenner, Proceedings of the National Academy of Sciences **118**, e2024083118 (2021).
[21] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, second edition ed., Adaptive Computation and Machine Learning Series (The MIT Press, Cambridge, Massachusetts, 2018).
[22] M. Nguyen and S. Vaikuntanathan, Proceedings of the National Academy of Sciences **113**, 14231 (2016).
[23] The set of functions $g$ satisfies technical assumptions detailed in Appendix A.
[24] C. Villani, *Topics in Optimal Transportation*, Graduate Studies in Mathematics No. 58 (American

Mathematical Society, Providence, Rhode Island, 2003).

[25] G. Peyré and M. Cuturi, arXiv:1803.00567 [stat] (2018), arXiv:1803.00567 [stat].

[26] D. A. Sivak and G. E. Crooks, Physical Review Letters **108**, 150601 (2012).

[27] B. Recht, Annual Review of Control, Robotics, and Autonomous Systems **2**, 253 (2019).

[28] C. J. C. H. Watkins and P. Dayan, in *Reinforcement Learning*, edited by R. S. Sutton (Springer US, Boston, MA, 1992) pp. 55–68.

[29] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, in *Deep Learning Workshop: Advances in Neural Information Processing Systems*, Vol. 26, edited by C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger (Curran Associates, Inc., 2013).

[30] S. Fujimoto, H. van Hoof, and D. Meger, in *Proceedings of the 35th International Conference on Machine Learning*, Proceedings of Machine Learning Research, Vol. 80, edited by J. Dy and A. Krause (PMLR, 2018) pp. 1587–1596.

[31] M. E. Cates and J. Tailleur, Annual Review of Condensed Matter Physics **6**, 219 (2015).

[32] M. F. Hagan, A. Baskaran, and G. S. Redner, Physical Review Letters **110**, 055701 (2013).

[33] J. Palacci, S. Sacanna, A. P. Steinberg, D. J. Pine, and P. M. Chaikin, Science **339**, 936 (2013).

[34] M. J. Falk, V. Alizadehyazdi, H. Jaeger, and A. Murugan, arXiv:2105.04641 [cond-mat] (2021), arXiv:2105.04641 [cond-mat].

[35] It is difficult to directly compare the duration of training because the replay buffer grows more quickly for systems with more regions of control.

[36] G. Makrides, B. Zinsser, A. Phinikarides, M. Schubert, and G. E. Georghiou, Renewable Energy **43**, 407 (2012).

[37] S. Dey, C. Fan, K. V. Gothelf, J. Li, C. Lin, L. Liu, N. Liu, M. A. D. Nijenhuis, B. Saccà, F. C. Simmel, H. Yan, and P. Zhan, Nature Reviews Methods Primers **1**, 1 (2021).

[38] A. L. Gibbs and F. E. Su, International Statistical Review **70**, 419 (2002).

[39] Y. Pantazis and M. A. Katsoulakis, The Journal of Chemical Physics **138**, 054115 (2013).

**Appendix A: Relation between coarse-grained control and dissipation**

Let $\boldsymbol{x} \in \Omega \subset \mathbb{T}^d$ denote a configuration on the system, taking coordinates on the torus due to periodic boundary conditions. We set

$$C = \sup_{\boldsymbol{x},\boldsymbol{y} \in \Omega} \|\boldsymbol{x} - \boldsymbol{y}\|, \tag{A1}$$

the diameter on this compact space. Because it is impractical to directly specify the target distribution $\rho_* : \Omega \to \mathbb{R}$ for a high-dimensional interacting particle system, we instead specify the target value of a given observable $f : \mathbb{R}^d \to \mathbb{R}$ and denote this target value by $f_*$. For technical reasons described below, we assume that $f$ is Lipschitz continuous (essentially meaning that its derivative remains bounded) with Lipschitz constant $K = 1$; any value of $K$ could be used and consequently the bound derived below will have a prefactor of $CK$. We assume that the dynamics of system is governed by an overdamped Langevin equation controlled by an external force $u$; that is,

$$d\boldsymbol{X}_t^u = [b(\boldsymbol{X}_t^u) + u(\boldsymbol{X}_t^u)]dt + \sqrt{2D}d\boldsymbol{W}_t, \tag{A2}$$

as discussed in the main text. We also assume that the resulting dynamics is ergodic so that the process (A2) relaxes to a unique stationary distribution $\rho_u$.

Naturally, the objective function for the controller seeks to match the time-averaged value of the observable $f$ with the target value $f_*$. That is, we want to find the optimal controller $u_*$ that solves

$$u_* = \underset{u}{\operatorname{argmin}} |\mathbb{E}_u f - f_*| \tag{A3}$$

over all control functions $u$. Here the notation $\mathbb{E}_u$ denotes an expectation

$$\mathbb{E}_u f = \int_\Omega f(\boldsymbol{x})\rho_{\mathrm{ss}}^u(\boldsymbol{x})d\boldsymbol{x} = \lim_{T \to \infty} \frac{1}{T} \int_0^T f(\boldsymbol{X}_t^u)dt, \tag{A4}$$

where $\boldsymbol{X}_t^u$ solves (A2). Though we do not know functional form of the target steady state density, we assume that it is "implementable" in the sense that

$$f_* = \mathbb{E}_* f = \lim_{T \to \infty} \frac{1}{T} \int_0^T f(\boldsymbol{X}_t^*)dt, \tag{A5}$$

where

$$d\boldsymbol{X}_t^* = b_*(\boldsymbol{X}_t^*)dt + \sqrt{2D}d\boldsymbol{W}_t, \tag{A6}$$

for some unknown $b_*$.

In the approach we take, the optimal controller depends on the choice of observable. Some choices of $f$ may not be particularly informative about the system—for example, the bulk density could remain fixed for steady state distributions with differing microscopic structures. However, issues with uninformative observables could be overcome by allowing any 1-Lipschitz function $g : \Omega \to \mathbb{R}$ and carrying out the "adversarial" optimization

$$\max_g \min_u |\mathbb{E}_u g - \mathbb{E}_* g| \tag{A7}$$

so that the optimal controller must minimize the mean discrepancy for *every* observable. This stringent requirement actually amounts to minimizing a measure of the distance between the probability distribution with density $\rho_u$ and the unspecified target distribution $\rho_*$. In fact, the maximization coincides with the Kantorovich-Rubenstein dual formulation of the Wasserstein-1 optimal transport distance [24],

$$\mathcal{W}_1(\rho_u, \rho_*) = \inf_{\pi \in \Pi(\rho_u, \rho_*)} \int_{\Omega \times \Omega} |\boldsymbol{x} - \boldsymbol{y}|d\pi(\boldsymbol{x}, \boldsymbol{y}) \tag{A8}$$

where $\Pi(\rho_u, \rho_*)$ denotes all distributions with marginals $\rho_u$ and $\rho_*$. This so-called integral probability metric is related to other measures of "distance" between probability distribution [38] and is upper bounded by the total variation distance

$$\mathcal{W}_1(\rho_u, \rho_*) \leq C\|\rho_u - \rho_*\|_{\mathrm{TV}} \tag{A9}$$

which is proved by using the coupling definition of the total variation distance and noting that the expected distance on the coupling that realizes the infimum in (A8) must be less than the maximum distance between points in $\Omega$ multiplied by the expected fraction of unequal points. Via Pinsker's inequality, this bound can be extended to the KL divergence between the two distributions:

$$\mathcal{W}_1(\rho_u, \rho_*) \leq C\|\rho_u - \rho_*\|_{\mathrm{TV}}, \leq C\sqrt{D_{\mathrm{KL}}(\rho_u\|\rho_*)}. \tag{A10}$$

The KL divergence, or relative entropy, is a well-studied object in nonequilibrium statistical mechanics because it provides a physical measure of entropic distance between distributions [26]. The steady state distributions appearing in (5) differ with an entropic cost that can be bounded by considering the time average of the relative entropy of the path measures associated with the corresponding nonequilibrium trajectories of duration $T$ [39],

$$
\begin{aligned}
D_{\mathrm{KL}}(\rho_u\|\rho_*) &\leq \frac{1}{T} D_{\mathrm{KL}}(\mathbb{P}_u\|\mathbb{P}_*) \\
&= \frac{1}{T}\mathbb{E}_{\mathbb{P}_u} \log \frac{d\mathbb{P}_u}{d\mathbb{P}_*}, \\
&= \mathbb{E}_{\mathbb{P}_u} \frac{1}{\sqrt{2DT}} \int_0^T \delta u(X_t) dW_t + \frac{1}{2DT}\int_0^T \|\delta u(X_t)\|^2 dt.
\end{aligned}
\tag{A11}
$$

In the expressions above $\delta u = b_* - (b + u)$. The stochastic integral vanishes in the limit $T \to \infty$ and we see that the obvious minimizer can be deduced:

$$\mathbb{E}_{\mathbb{P}_u} \frac{1}{2DT} \int_0^T \|b_*(X_t) - b(X_t) - u(X_t)\|^2 dt \implies u_*(X_t) = b_*(X_t) - b(X_t) \tag{A12}$$

where $u_* = \mathrm{argmin}_u D_{\mathrm{KL}}(\mathbb{P}_u\|\mathbb{P}_*)$, which matches the target drift function exactly. However, there is an entropic cost associated with this choice compared to the path measure $\mathbb{P}$ of the unperturbed dynamics;

$$D_{\mathrm{KL}}(\mathbb{P}_u\|\mathbb{P}) = \frac{1}{2DT}\mathbb{E}_{\mathbb{P}_u} \int_0^T \|u(X_t^u)\|^2 dt \equiv \sigma_{\mathrm{ex}} \geq 0 \tag{A13}$$

is non-negative excess entropy associated with the control. This quantity vanishes when $u \equiv 0$ meaning that no control is required to produce the steady state. In the long time limit, it could also vanish if $u = \nabla V$, for some potential energy $V$, is a conservative force, the "interaction design" paradigm discussed in the introduction—in this case any excess dissipation would only be transient. We conclude that for the protocol $u$ that results from optimization (via reinforcement learning or some other technique) will dissipate at least $\sigma_{\mathrm{ex}}$ in order match the true target distribution.

### Appendix B: Computational details of for Deep Q-Learning

*Code Availability:* All simulations were conducted using OpenMM 7.5.1 and PyTorch 1.9.0. Our code is available at https://github.com/rotskoff-group/dissipative-design.

$Q$-learning is a model-free, off-policy Reinforcement Learning algorithm used to estimate the optimal state-action value function $Q^\star$. Because $Q$-learning is model-free we can avoid incorporating the dynamics of the system into the optimization framework, which may be unknown in experimental systems.

In many RL algorithms, we assume that the environment is a Markov decision process; in our case, this means we consider a dynamics specified by states $\boldsymbol{x} \in \mathbb{R}^d$, actions of the external protocol $\boldsymbol{a} \in [0, a_{\max}]^{n_{\text{regions}}}$, and a cost function $C$. In deep $Q$-learning[29], we represent the state-action value function, $Q$, using a deep neural network, and optimize this network using the Bellman dynamic programming principle. For a deterministic policy $u$, the Bellman equation reads

$$Q^u(\boldsymbol{x}_t, \boldsymbol{a}_t) = \mathbb{E}^{\boldsymbol{x}_t}_u \left[ C(\boldsymbol{x}_{t+\tau}; \boldsymbol{a}_t) + \sum_{k=1}^{\infty} \gamma^k C(\boldsymbol{x}_{t+\tau(k+1)}) \right]. \tag{B1}$$

In this expression, the cost is evaluated for a state $\boldsymbol{x}_{t+\tau}$ after taking action $\boldsymbol{a}_t$ and subsequently following the protocol $u$. In practice, we use a second target neural $Q'$ network to estimate the predicted value of $Q^u(\boldsymbol{x}_{t+\tau}, u(\boldsymbol{x}_{t+\tau}))$ to ensure more stable updates. We update the weights of the target network to match the weights of the $Q$ network at a rate $\tau$. Finally, we use experience replay, a technique in which we store past experiences in a replay buffer $\mathcal{R}$ and sample from the buffer during update steps. Each experience here is a tuple of $(\boldsymbol{x}_t, \boldsymbol{a}_t, c_t, \boldsymbol{x}_{t+1})$ that was observed within each "grid" (i.e. a spatial region) of control.

Our state space $\mathcal{S}$ represents the normalized empirical cluster-size distribution over some sample space $\Omega$, which was specified differently for each of our target distributions, as described in the main text. We found that this representation worked the best compared to other possible state representations, such as images of the system configurations. Each state $\boldsymbol{x}^i_t$ was the normalized cluster-size distribution within a grid of the system $\boldsymbol{X}_t$. This allowed us to utilize local information about the system to spatially control the system.

$Q$-Learning methods are known to suffer from underestimation bias [30]. This bias arises during training because we are using $\min_{\boldsymbol{a}} Q^u(\boldsymbol{x}, \boldsymbol{a})$ as an estimate for $\min_{\boldsymbol{a}} Q^{\star}(\boldsymbol{x}, \boldsymbol{a})$. By Jensen's Inequality, we have

$$\mathbb{E} \min_{\boldsymbol{a}} Q(\boldsymbol{x}, \boldsymbol{a}) \leq \min_{\boldsymbol{a}} \mathbb{E} Q(\boldsymbol{x}, \boldsymbol{a}) = \min_{\boldsymbol{a}} Q^{\star}(\boldsymbol{x}, \boldsymbol{a}) \tag{B2}$$

Underestimation bias results in suboptimal actions, which will have an artificially lower state-action value, being selected as the optimal action. One approach to mitigating this bias is to use clipped double $Q$-learning, where we maintain two estimates of $Q^{\star}$: $Q_1$ and $Q_2$ [30]. We then use the maximum of the two estimates provided by these networks as an upper bound of the estimate of the state-action value. With this estimate, we update both $Q_1$ and $Q_2$. Finally, we update the weights of the target networks $Q'_1$ and $Q'_2$ towards the weights of $Q_1$ and $Q_2$ at a rate $\tau$.

Reinforcement Learning algorithms are constrained by a trade-off between exploration and exploitation [21]. In order to determine an optimal policy $\pi$, it is necessary to *exploit* our current estimates of the state-action value. However, it is necessary to sufficiently *explore* the state-action space in order to improve our current estimates of the state-action value. One common approach to this dilemma is to use an $\epsilon$-greedy search, where a random action is selected with probability $\epsilon$ during training to promote exploration. In our approach, we instead always randomly select an action during the first $e_{\text{explore}}$ episodes of training. After the first $e_{\text{explore}}$ episodes, we decided to train our state-action value function using a greedy approach (i.e. selecting the optimal action that minimizes the cost). We used this approach because of the small size of our action space $\mathcal{A}$ and because the decisions made in each grid were considered to be independent experiences. This allowed us to more thoroughly explore the action-space during the first $e_{\text{explore}}$ episodes, especially when using higher-resolutions of control.

When using higher-resolutions of control, it is possible for a grid to not contain any particles. During training, if a grid did not contain any particles at time $t$ or $t+1$, it was not included in the replay buffer $\mathcal{R}$. While training with the greedy approach (i.e. if the episode $e \geq e_{\text{explore}}$), we set the action to be $a_{\max}$ if the grid did not contain any particles at some time $t$. For our Lennard-Jones system, $a_{\max}$ represented the highest possible temperature, in order to ensure that any free particle that entered the region would continue to diffuse until it reached a region with more particles. For our active matter system $a_{\max}$ represented the highest possible activity, to similarly ensure that any

| Hyper-Parameters | | |
|---|---|---|
| Hyperparameter | Active Colloids | LJ System |
| $Q$ learning rate | $3 \times 10^{-4}$ | $3 \times 10^{-4}$ |
| Optimizer | Adam | Adam |
| Target Update Rate ($\tau$) | $5 \times 10^{-3}$ | $5 \times 10^{-3}$ |
| Batch Size | 32 | 32 |
| Discount Factor ($\gamma$) | 0.9 | 0.95 |
| Number of Hidden Layers | 2 | 1 |
| $e_{\text{explore}}$ | 5 | 25 |
| $T$ | 350 | 150 |
| $a_{high}$ | 1.5 | 1.0 |
| $\mathcal{A}$ | $[0, 0.25, 0.5, 0.75, 1.0, 1.25, 1.5]$ | $[0.01, 0.25, 1.0]$ |

TABLE I. Relevant parameters for $Q$-learning training

free particle that entered the region would continue self-propulsion until it was able to form clusters. Finally, when a cluster was located in multiple grids, each of these grids was considered to contain the entire cluster. This effectively allows neighboring grids that shared a cluster to "communicate," which was especially important during the high-resolution control cases.

---

**Algorithm 1** Clipped Double Q Learning Training

---

1: **Initialize Replay Buffer $\mathcal{R}$ to capacity $\mathcal{N}$**
2: **Randomly initialize $Q_1$ and $Q_2$ networks with weights $\theta^{Q_1}$ and $\theta^{Q_2}$**
3: **Initialize target networks $Q_1'$ and $Q_2'$ networks with weights $\theta^{Q_1'} \leftarrow \theta^{Q_1}$ and $\theta^{Q_2'} \leftarrow \theta^{Q_2}$**
4: **for** e $= 0 \ldots \boldsymbol{M}$ **do**
5:     Initialize state $\boldsymbol{X}_0$ for episode e
6:     **for** t $= 0 \ldots \boldsymbol{T}$ **do**
7:         **for each $\boldsymbol{x}_t$ in $\boldsymbol{X}_t$ do**
8:             **if $\boldsymbol{x}_t$ is empty and e $\geq e_{\text{explore}}$ then**
9:                 Select $a_t = a_{\max}$ and execute action $a_t$
10:                 **continue**
11:             **end if**
12:             **if e $< e_{\text{explore}}$ then**
13:                 Select a random action $a_t$ from $\mathcal{A}$
14:             **else**
15:                 Select $a_t = \text{argmin}_{\boldsymbol{a}} Q_1(\boldsymbol{x}_t, \boldsymbol{a})$
16:             **end if**
17:             Execute action $a_t$ and observe next state $x_{t+1}$ and cost $c_t \equiv C$
18:             **if $\boldsymbol{x}_t$ or $\boldsymbol{x}_{t+1}$ is empty then**
19:                 **continue**
20:             **end if**
21:             Store transition $(\boldsymbol{x}_t, \boldsymbol{a}_t, c_t, \boldsymbol{x}_{t+1})$ in $\mathcal{R}$
22:         **end for**
23:         Sample random batch of transitions $(\boldsymbol{x}_j, \boldsymbol{a}_j, c_j, \boldsymbol{x}_{j+1})$ from $\mathcal{R}$
24:         Set $y_j = c_j + \gamma \max_{1,2} \min_{\boldsymbol{a}'} Q_{1,2}'(\boldsymbol{x}_{t+1}, \boldsymbol{a})$
25:         Update $Q_{1,2}$ by minimizing $L_{1,2} = (y_j - Q_{1,2}(\boldsymbol{x}_j, \boldsymbol{a}_j))^2$
26:         Update the target networks $\theta^{Q_{1,2}'} \leftarrow (1 - \tau)\theta^{Q_{1,2}'} + \tau\theta^{Q_{1,2}}$
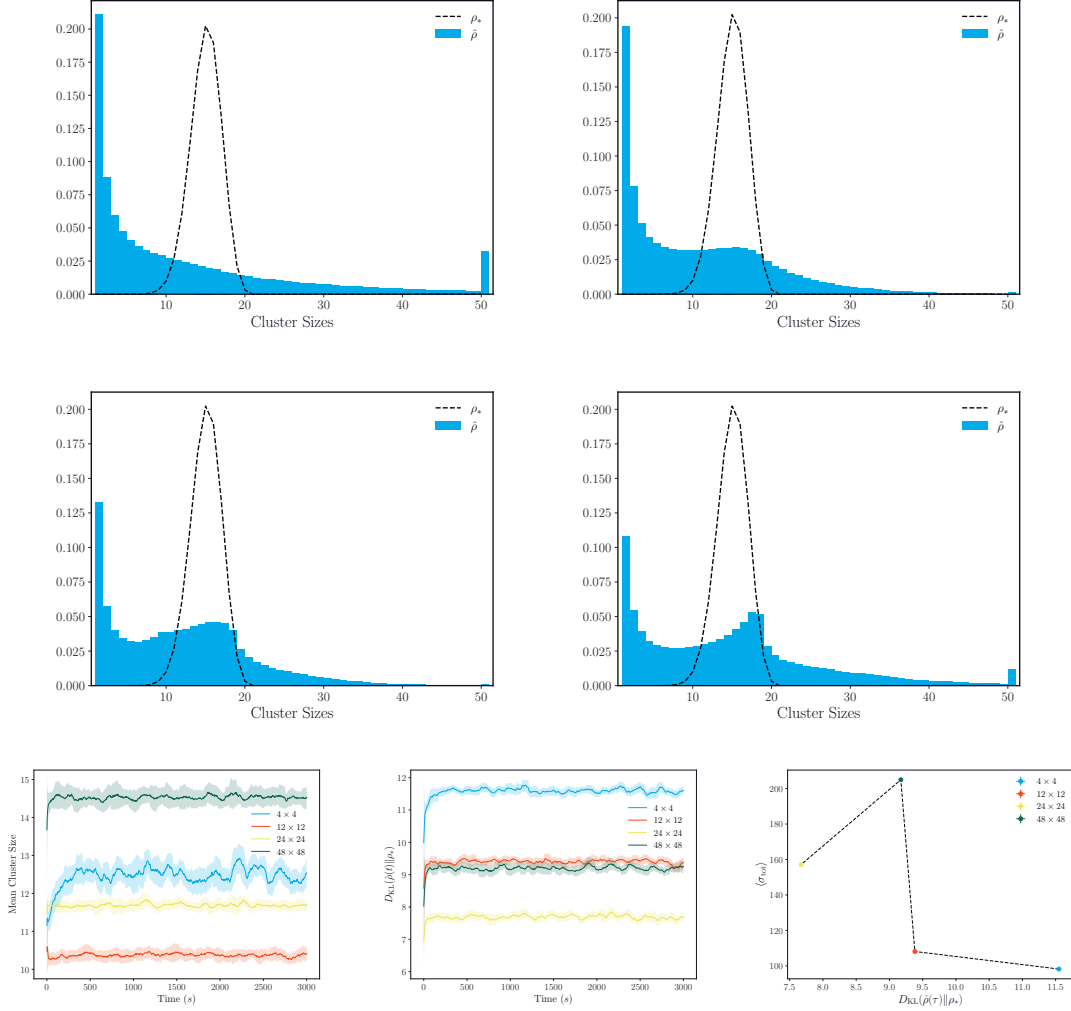27:     **end for**
28: **end for**

---

FIG. 4. Summary of results for Binomial target distribution. Histograms of cluster sizes for $4 \times 4$ (top left) $12 \times 12$ (top right) $24 \times 24$ (center left) and $48 \times 48$ (center right). Mean cluster sizes (bottom left), KL cost function (bottom center), and total entropy production as a function $D_{\mathrm{KL}}$ (bottom right).

## Appendix C: Active Colloids

We modeled our active colloids based on the experimental system in [33]. In this system, colloidal activity can be modulated by blue light. In addition to self-propulsion, these particles exhibit an attractive interaction due to phoretic and osmotic effects when activated by blue light.

We work in normalized units of the particle diameter and the self-propulsion velocity. The model consists of a purely repulsive interaction

$$U_{\mathrm{WCA}}(r) = \begin{cases} 0 & r > r_{\mathrm{cut}} \\ 4\epsilon \left[ \left( \frac{\sigma}{r} \right)^{12} - \left( \frac{\sigma}{r} \right)^{6} \right] + \epsilon & r \leq r_{\mathrm{cut}} \end{cases} \tag{C1}$$

with $r_{\mathrm{cut}} = 2^{1/6}\sigma$. In addition to the hard sphere repulsion, there is an attractive force induced by
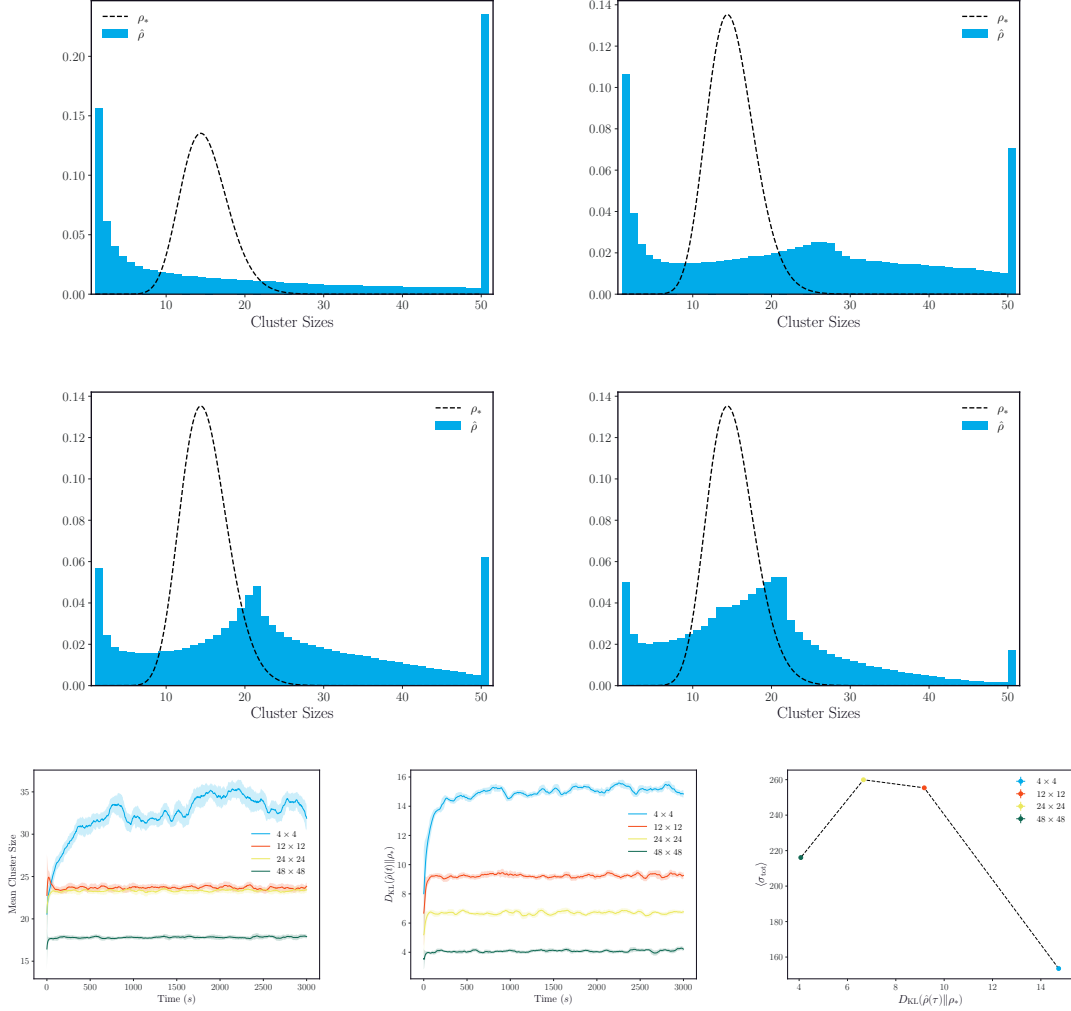
FIG. 5. Summary of results for Gamma target distribution. Histograms of cluster sizes for $4 \times 4$ (top left) $12 \times 12$ (top right) $24 \times 24$ (center left) and $48 \times 48$ (center right). Mean cluster sizes (bottom left), KL cost function (bottom center), and total entropy production as a function $D_{\mathrm{KL}}$ (bottom right).

the activity which models hydrodynamic effects,

$$U_{\mathrm{attract}}(\boldsymbol{x}_i, \boldsymbol{x}_j) = \frac{\sqrt{A_i A_j}}{\|\boldsymbol{x}_i - \boldsymbol{x}_j\|^2} \tag{C2}$$

where the coefficient $A_i$ is determined by the instantaneous value of the activity

$$A_i = \alpha(\boldsymbol{x}_i)^2 A_0 \tag{C3}$$

and $A_0$ is a constant. Here, $\alpha(\boldsymbol{x}_i)$ is the predicted "action" by the RL algorithm.

In addition to the conservative force that arises from these potential terms, there is an active force

$$F_{\mathrm{active}}(\boldsymbol{x}) = (\alpha(\boldsymbol{x})\cos(\boldsymbol{\theta}), \alpha(\boldsymbol{x})\sin(\boldsymbol{\theta}), 0) \tag{C4}$$
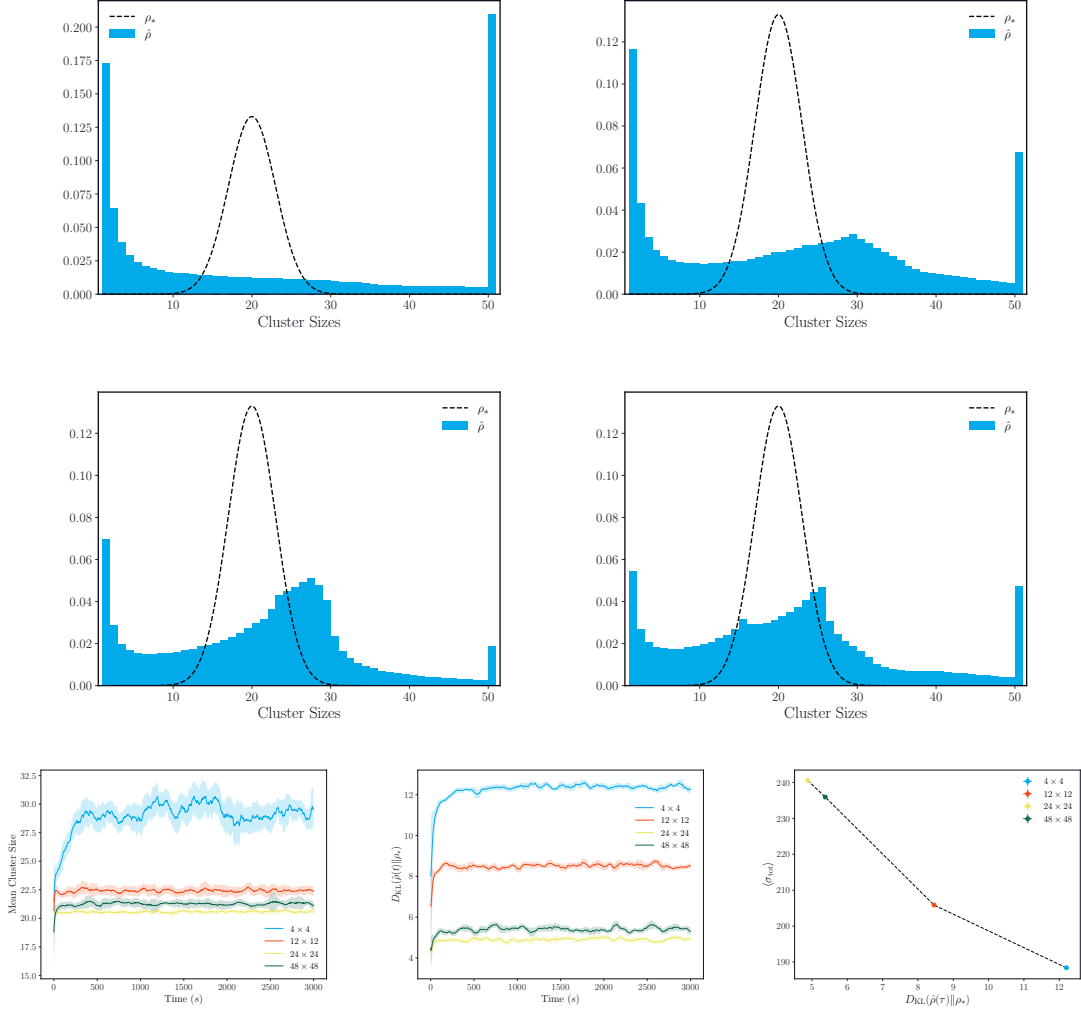
FIG. 6. Summary of results for Gaussian target distribution. Histograms of cluster sizes for $4 \times 4$ (top left) $12 \times 12$ (top right) $24 \times 24$ (center left) and $48 \times 48$ (center right). Mean cluster sizes (bottom left), KL cost function (bottom center), and total entropy production as a function $D_{\mathrm{KL}}$ (bottom right).
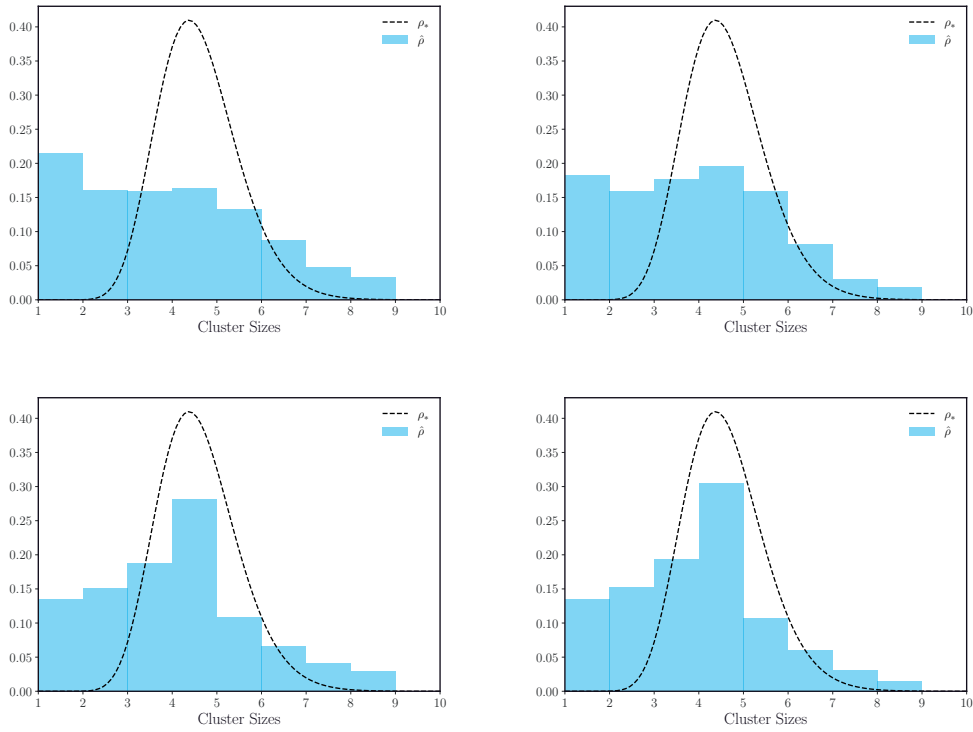
where $\boldsymbol{\theta}$ is a vector of particle directions which itself has a purely diffusive dynamics

$$d\boldsymbol{\theta}_t = \sqrt{2D_r}d\boldsymbol{W}_t. \tag{C5}$$

$$dX_t = -[\nabla U_{\mathrm{attract}}(X_t) - \nabla U_{\mathrm{WCA}}(X_t) + F_{active}(X_t, \boldsymbol{\theta}_t)]dt + \sqrt{2D_t}d\boldsymbol{W}_t \tag{C6}$$

| Simulation Parameters | |
|---|---|
| Parameter | Value (Normalized Units) |
| $D_{\mathrm{r}}$ | 0.125 |
| $D_{\mathrm{t}}$ | 0.041667 |
| $\epsilon$ | 0.5 |
| $\sigma$ | 1 |
| $A_0$ | 0.87 |
| $\Delta t$ | 0.00005 |

TABLE II. Parameters for active colloid system.



FIG. 7. Histograms of cluster sizes for $3 \times 3$ (top left) $4 \times 4$ (top right) $12 \times 12$ (bottom left) and $15 \times 15$ (bottom right).

## Appendix D: Lennard-Jones System

To investigate feedback guided thermal annealing, we modeled a system of colloidal particles that interact via Lennard-Jones interactions.

$$U_{\mathrm{LJ}}(\boldsymbol{x}_i, \boldsymbol{x}_j) = 4\epsilon \left[ \left( \frac{\sigma}{\|\boldsymbol{x}_i - \boldsymbol{x}_j\|} \right)^{12} - \left( \frac{\sigma}{\|\boldsymbol{x}_i - \boldsymbol{x}_j\|} \right)^{6} \right] \tag{D1}$$

Equation of Motion:

$$dX_t = -\nabla U_{\mathrm{LJ}}(X_t)dt + \sqrt{2\beta_t^{-1}(X_t)}dW_t \tag{D2}$$

As described, we update the temperature of each grid of the system based on our Reinforcement Learning approach. The $\beta_t^{-1}$ for a particle depends on where the particle is at the beginning of a decision and is not changed during the duration of a decision. Because our decision length is only 0.25 seconds, our particle will not, on average, diffuse between grids within a decision even for our highest resolution of control. At the beginning of the next decision, we instantaneously update the temperature of each grid and subsequently update the $\beta_{t+1}^{-1}$ for each particle depending on the temperature of the grid in which it is located.

| Simulation Parameters | |
|---|---|
| Parameter | Value (Normalized Units) |
| $\epsilon$ | 0.5 |
| $\sigma$ | 1 |
| $\Delta t$ | 0.0002 |

TABLE III. Parameters for the thermal annealing simulations.