

Reinforcement Learning informs optimal treatment strategies to limit antibiotic resistance.

Davis T. Weaver^{1,2}, Jeff Maltas^{2,†}, and Jacob G. Scott^{1,2,3,†}

¹Case Western Reserve University School of Medicine, Cleveland, OH, 44106, USA

²Translational Hematology Oncology Research, Cleveland Clinic, Cleveland OH, 44106, USA

³Department of Physics, Case Western Reserve University, Cleveland, OH, 44106, USA

†jeff.maltas@gmail.com, scottj10@ccf.org

ABSTRACT

Drug resistant pathogens are a wide-spread and deadly phenomenon. Antimicrobial resistance was estimated to be associated with 4.95 million deaths worldwide in 2019. If resistance continues to develop at the current rate, bacterial infections are expected to surpass cancer as the leading cause of death worldwide by 2050. Despite this troubling trend, antimicrobial drug development has all but ceased. For the few new drugs that are approved, microbes develop rapid resistance through evolution by mutation and selection. Novel approaches to designing therapy that explicitly take into account the adaptive nature of microbial cell populations are desperately needed. Approaches that can design therapies given limited information about the evolving system are particularly important due to the limitations of clinical measurement. In this study, we explore a reinforcement learning (RL) approach capable of learning effective drug cycling policies in a system defined by empirically measured fitness landscapes. Given access to a panel of 15 β -lactam antibiotics with which to treat the simulated *E. Coli* population, we demonstrate that RL agents outperform two potential treatment paradigms at minimizing the population fitness over time. We also show that RL agents approach the performance of the optimal drug cycling policy. Crucially, we show that it is possible for RL agents to learn effective drug cycling protocols using current population fitness as the only training input. Our work represents a proof-of-concept for using AI to control complex evolutionary processes.

Conflict of Interest Statement: The authors have no conflicts of interest to disclose.

1 Introduction

2 Drug resistant pathogens are a wide-spread and deadly phenomenon that were responsible for nearly 5 million deaths worldwide
3 in 2019¹. In the US alone, 3 million cases of antimicrobial resistant infections are observed each year². The increasing
4 prevalence of pan-drug resistance has prompted the CDC to declare that we have entered a “post-antibiotic era”². Despite this
5 evident public health crisis, development of novel antibiotics has all but ceased due to the poor return on investment currently
6 associated with this class of drugs³. Novel approaches to designing therapies that explicitly take into account the adaptive
7 nature of microbial cell populations while leveraging existing treatment options are desperately needed.

8 Evolutionary medicine is a rapidly growing discipline that aims to develop treatment strategies that explicitly account
9 for the capacity of pathogens and cancer to evolve^{4–10}. Such treatment strategies, termed “evolutionary therapies”, typically
10 cycle between drugs or drug doses to take advantage of predictable patterns of disease evolution. Evolutionary therapies are
11 typically developed by applying optimization methods to a mathematical or simulation-based model of the evolving system
12 under study^{11–21}. For example, in castrate-resistant prostate cancer, researchers developed an on-off drug cycling drug protocol
13 that allows drug-sensitive cancer cells to regrow following a course of treatment^{22,23}. Clinical trials have shown this therapy
14 prevents the emergence of a resistant phenotype and enables superior long-term tumor control and patient survival compared to
15 conventional strategies^{22,23}.

16 Current methods for the development of evolutionary therapies require an enormous amount of data on the evolving system.
17 For example, many researchers have optimized treatment by using genotype-phenotype maps to define evolutionary dynamics
18 and model the evolving cell population^{15,24–32}. However, most methods for optimization of these models requires a complete
19 understanding of the underlying system dynamics^{14,15,33,34}. Such detailed knowledge is currently unobtainable in the clinical
20 setting. Approaches that can approximate these optimal policies given only a fraction of the available information would fill a
21 key unmet need in evolutionary medicine.

22 We hypothesize that reinforcement learning algorithms can develop effective drug cycling policies given only experimentally
23 measurable information about the evolving pathogen. Reinforcement learning (RL) is a well-studied subfield of machine learning
24 that has been successfully used in applications ranging from board games and video games to manufacturing automation^{33,35–37}.
25 Broadly, RL methods train artificial intelligence agents to select actions that maximize a reward function. Importantly, RL
26 methods are particularly suited for optimization problems where little is known about the dynamics of the underlying system.
27 Further, RL and related optimal control methods have been previously applied for the development of clinical optimization
28 protocols in oncology and anesthesiology^{20,38–43}.

29 In this study, we develop of a novel approach to discover
30 evolving evolutionary therapies, using a well studied set of empiri-
31 cal fitness landscapes as a model system. We will explore
32 “perfect information” optimization methods such as dynamic
33 programming in addition to RL methods that can learn policies
34 given only limited information about a system. We show that
35 it is possible to learn effective drug cycling treatments given
36 extremely limited information about the evolving population,
37 even in situations where the measurements reaching the RL
38 agent is extremely noisy and the information content is low.

39 1 Methods

40 As a model system, we simulated an evolving population of *Es-*
41 *cherichia coli* (*E. coli*) using the well-studied fitness landscape
42 paradigm, where each genotype is associated with a certain
43 fitness under selection^{15,25,28}. We relied on a previously de-
44 scribed 4 allele landscape of the *E. Coli* β -lactamase gene
45 where each mutation has a measured impact on the sensitivity
46 of an *E. coli* population to one of 15 β -lactam antibiotics^{25,28}.

47 We then defined 15 different fitness regimes on the same underlying genotype space, each representing the selective effect of
48 one of 15 β -lactam antibiotics (Table 1)²⁵. We used this well-studied *E. coli* model system because it is one of the few microbial
49 cell populations for which a combinatorially complete genotype-phenotype mapping has been measured^{25,28}. By simulating an
50 evolving *E. coli* cell population using the described fitness landscape paradigm, we were able to define an optimization problem
51 on which to train RL agents (**Fig 1**).

index	drug code	drug
1	AMP	Ampicillin
2	AM	Amoxicillin
3	CEC	Cefaclor
4	CTX	Cefotaxime
5	ZOX	Ceftizoxime
6	CXM	Cefuroxime
7	CRO	Ceftriaxone
8	AMC	Amoxicillin + Clavulanic acid
9	CAZ	Ceftazidime
10	CTT	Cefotetan
11	SAM	Ampicillin + Sulbactam
12	CPR	Cefprozil
13	CPD	Cefpodoxime
14	TZP	Piperacillin + Tazobactam
15	FEP	Cefepime

Table 1. Reference codes for drugs under study

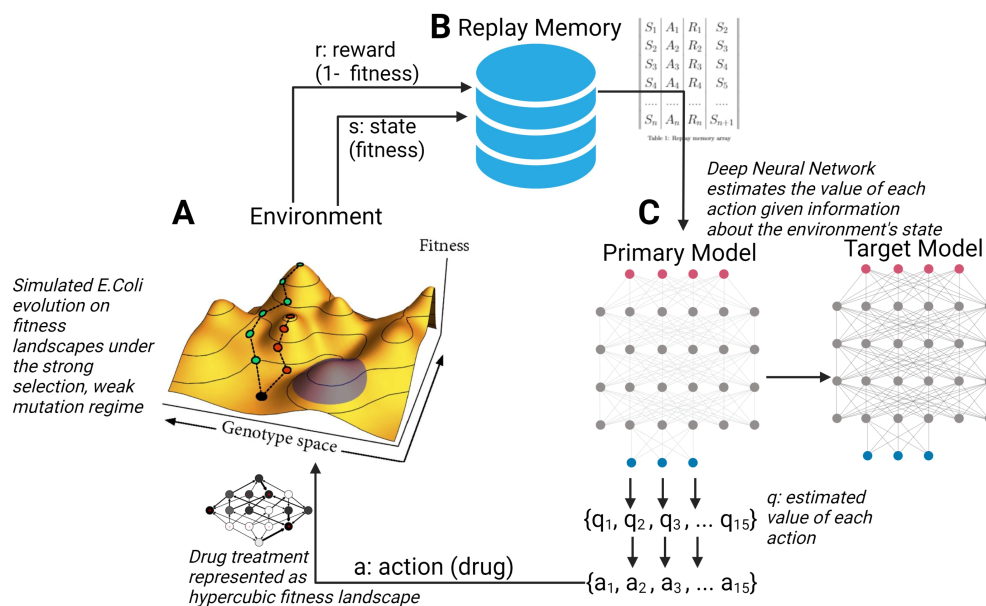


Figure 1. Schematic of artificial intelligence system for controlling evolving cell populations. **A:** *E. coli* population evolving on fitness landscapes under the strong selection, weak mutation evolutionary regime. At each time step, a reward signal r and a measure of system state s are sent to the replay memory structure. **B:** Replay memory array stores (s, a, r, s') tuples where s' is state $s+1$. These are then used to batch train the neural network. **C:** Deep Neural network estimates the value of each action given information about the environment's state. The action with the largest estimated value is then applied to the evolving cell population.

52 1.1 Simulation of Evolution Using Fitness Landscapes

53 We use a previously described fitness-landscape based model of evolution^{15,26}. In brief, we begin by modeling an evolving
 54 asexual haploid population with N mutational sites. Each site can have one of two alleles (0 or 1). We can therefore represent
 55 the genotype of a population using an N -length binary sequence, for a total of 2^N possible genotypes. We can model theoretical
 56 drug interventions by defining fitness as a function of genotype. These “drugs” can then be represented using N -dimensional
 57 hyper-cubic graphs (Fig 1A). Further, if we assume that drug evolution under drug treatment follows the strong selection and
 58 weak mutation paradigm, we can then compute the probability of mutation between adjacent genotypes and represent each
 59 landscape as a Markov chain as described by Nichol et al.¹⁵. At each time step, we sampled from the probability distribution
 60 defined by the Markov chain to simulate the evolutionary course of a single population.

61 1.2 Optimization approaches

62 We applied two related optimization approaches to identify effective drug cycling policies in this setting. First, we extended the
 63 Markov chain framework to formulate a complete Markov Decision Process (MDP). An MDP is a discrete-time framework for
 64 modeling optimal decision-making³³. Critically, the system under study must be partially under the control of the decision-
 65 making agent. MDPs can be solved using dynamic programming to generate optimal policies for the defined control problem³³.
 66 The dynamic programming algorithm requires perfect information (e.g. the complete transition matrix and instantaneous state
 67 from the MDP) in order to yield optimal policies. Next, we trained agents with imperfect information using reinforcement
 68 learning to approximate a clinical scenario where perfect information will never be available. Notably, the state set, action set,
 69 and reward assignment were shared between the perfect and imperfect information conditions. The action set corresponded to
 70 the drugs available to the optimization process. We considered this system to have a finite time horizon (20 evolutionary steps
 71 in the base case). We chose a finite time horizon rather than an infinite time horizon assumption in order to more faithfully
 72 represent clinical disease courses. For our purposes, we assume that one evolutionary time step is the equivalent of one day of
 73 evolution.

74 1.2.1 Perfect Information

75 The state set S represents all potential genotypes (16 total in our base case) that the evolving population can explore. The
 76 action set A corresponds to the 15 available β -lactam antibiotics. Finally, we define the reward set (R) and the set of transition

77 probabilities (P) as a function of the current genotype s as well as the chosen action, a (eq 1):

$$\begin{aligned} R &= 1 - f(s|a) \text{ for } s \in S \text{ and } a \in A, \text{ and,} \\ P &= f(s_{t+1}|s_t, a_t) \text{ for } s \in S \text{ and } a \in A. \end{aligned} \tag{1}$$

78 We solved the defined MDP using backwards induction, a dynamic programming approach designed to solve MDPs with finite
79 time horizons⁴⁴, to generate an optimal drug cycling policy for each evolutionary episode. Backwards induction is used to
80 estimate a value function $V(s)$ which estimates the discounted reward of being in each state s . Optimal policies $\Pi(s)$ are then
81 inferred from the value function. Throughout the remainder of the paper we will refer to this optimal drug cycling policy as the
82 “MDP” condition.

83 1.2.2 Imperfect Information

84 In order to assess the viability of developing optimal drug therapies from potentially clinically available information, we trained
85 a Deep Q learner to interact with the evolving *E. Coli* system described above. Deep Q learning is an extremely well-studied
86 and characterized method of reinforcement learning, and is particularly suited to situations where very little *a priori* knowledge
87 about the environment is available^{33,45}. We used two different training inputs to model a gradient of information loss. In the
88 first condition, termed RL-genotype, the instantaneous genotype of the population was provided as the key training input at
89 each time point. For this condition, the neural architecture was composed of an input layer, two 1d convolutional layers, a max
90 pooling layer, a dense layer with 28 neurons, and an output layer with a linear activation function.

91 In the second condition, termed RL-fit, instantaneous population fitness of the population was provided as the key training
92 input at each time point. The neural architecture of RL-fit was composed of a neural network with an input layer, two dense
93 hidden layers with 64 and 28 neurons, and an output layer with a linear activation function. RL-fit takes population fitness at
94 time t and one-hot encoded action at time $t - 1$ as inputs and outputs Q-values. Q-values are estimates of the future value of a
95 given action. Q-value estimates are improved by minimizing the temporal difference between Q-values computed by the current
96 model and a target model, which has weights and biases that are only updated rarely. We used mean squared error (MSE) as the
97 loss function.

98 We further explored the effect of information content on learned policy effectiveness by introducing a noise parameter. With
99 noise active, fitness values $s \in S$ that were used as training inputs were first adjusted according to:

$$s_t = s_t + w \in W \sim \mathcal{N}(\mu, 0.05 \times \sigma^2). \tag{2}$$

100 For the noise experiment, μ was set to 0 such that $\sigma^2 = 0$ would introduce no noise. We then varied σ^2 (referred to as ‘noise
101 parameter’) from 0 (no noise) to 100 (profound loss of signal fidelity).

102 All code and data needed to define and implement the evolutionary simulation and reinforcement learning framework can be
103 found at https://github.com/DavisWeaver/evo_dm. The software can be installed in your local python environment using ‘pip
104 install git+https://github.com/DavisWeaver/evo_dm.git’

105 2 Results

106 In this study, we explored the viability of developing effective drug cycling policies for antibiotic treatment given less and
107 less information about the evolving system. To this end, we developed a reinforcement learning framework to design policies
108 that limit the growth of an evolving *E. Coli* population *in silico*. We evaluated this system in a well-studied *E. coli* system for
109 which empirical fitness landscapes for 15 antibiotics are available in the literature²⁵. A given RL agent could select from any of
110 these 15 drugs when designing a policy to minimize population fitness. We defined three experimental conditions. In the first,
111 we solved a Markov decision process formulation of the optimization problem under study. In doing so, we generated true
112 optimal drug cycling policies given perfect information of the underlying system (described in **Section 1.1**). In the second, RL
113 agents were trained using the current genotype of the simulated *E. Coli* population under selection (RL-genotype). Stepping
114 further down the information gradient, RL agents were trained using only observed fitness of the *E. Coli* population (RL-fit).
115 Finally, we introduced noise into these measures of observed fitness to simulate real-world conditions where only imprecise
116 proxy measures of the true underlying state may be available. Each experimental condition was evaluated based on its ability to
117 minimize the fitness of the population under study. We compared these conditions to two negative controls; a drug cycling
118 policy that selects drugs completely at random (which we will refer to as “random”), and all possible two-drug cycles (i.e
119 AMP-AM-AMP-AM-AMP). We tested 100 replicates of RL-fit and RL-genotype against each of these conditions. Each
120 replicate was trained for 500 episodes of 20 evolutionary steps (10000 total observations of system behavior). We chose 500

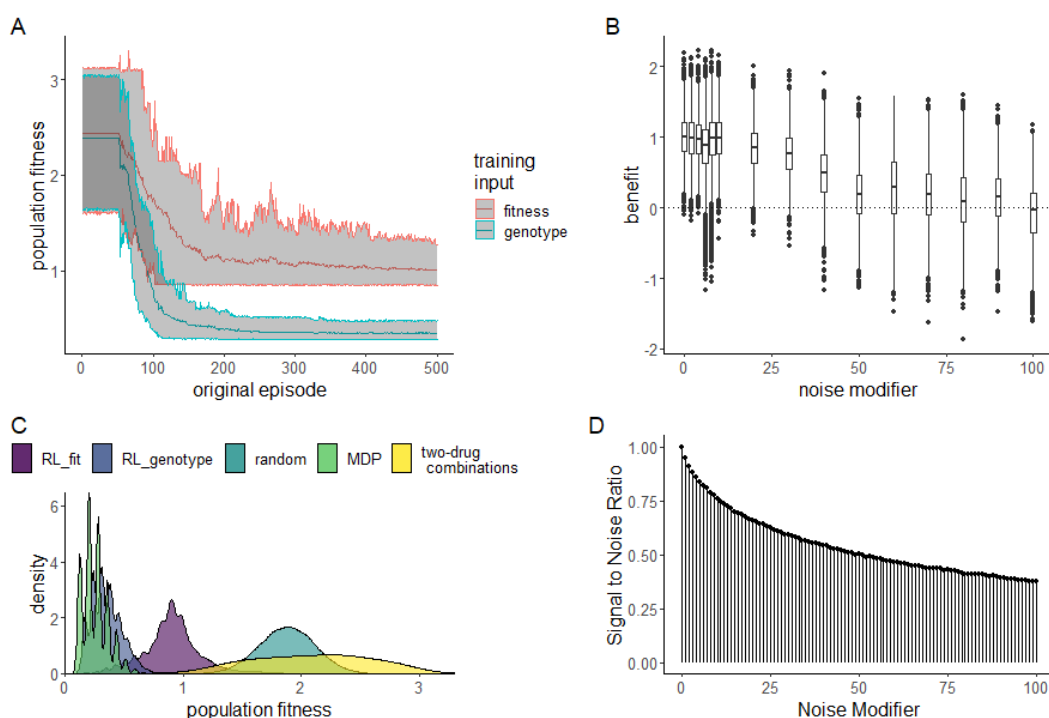


Figure 2. Performance of RL agents in a simulated *E. coli* system. **A:** Line plot showing the effectiveness of the average learned policy as training time increases on the x-axis for RL agents trained using fitness (red) or genotype (blue). **B:** Boxplot showing the effectiveness of 10 fully trained RL-fit replicates as a function of noise. Each data point corresponds to a single episode. The width of the distribution provides information about the episode by episode variability in RL-fit performance. **C:** Density plot summarizing the performance of the two experimental conditions (measured by average population fitness) relative to the three control conditions. **D:** Signal to noise ratio associated with different noise parameters. Increasing noise parameter decreases the fidelity of the signal that reaches the reinforcement learner.

121 episodes as the training time after extensive hyper-parameter tuning showed decreased or equal effectiveness with additional
 122 training.

123 Comparison of RL drug cycling policies to negative controls.

124 We found that both RL
 125 conditions dramatically re-
 126 duced fitness relative to the
 127 random policy. In both cases,
 128 the RL conditions learned ef-
 129 fective drug cycling policies
 130 after about 100 episodes of
 131 training and then fine-tuned
 132 them with minimal improve-
 133 ment through episode 500
 134 (Fig 2A). As expected, RL-
 135 genotype learned a more ef-
 136 fective drug cycling policy on average compared to RL-fit. RL-genotype had access to the instantaneous state (genotype) of
 137 the evolving population, while RL-fit was only trained using a proxy measure (population fitness). In 98/100 replicates, we
 138 observed a measurable decrease in population fitness under the learned RL-fit policy versus a random drug cycling policy
 139 (Fig S1A). Further, we found that the average RL-fit replicate outperformed all possible two-drug cycling policies (Fig 2C).
 140 RL-genotype outperformed both negative controls in all 100 replicates (Fig 2C). In some replicates, RL-genotype achieved
 141 similar performance compared to the MDP policy (Fig S1D). In addition, the distribution of performance for RL-genotype
 142 policies nearly overlapped with MDP performance (Fig 2C). Introduction of additional noise to the training process for RL-fit
 143 led to degraded performance. However, even with a large noise modifier, RL-fit still outperformed the random drug cycling
 144 condition. With a noise modifier of 40 (fitness + $\mathcal{N}(\mu = 0, \sigma^2 = 0.05 \times 40)$), RL-fit achieved an average population fitness of
 145 1.41 compared to 1.88 for the random drug cycling condition (Fig 2D).

drug sequence	replicate	condition
CTX,AMC,CTX,CPR,CTX,CPR,CTX,CPR,CTX,CPR	53	RL-fit
CTX,CPR,CPR,CPR,CTX,CPR,CPR,CPR,CTX,SAM	53	RL-genotype
CTX,AMC,CTX,AMC,CTX,AMC,CTX,AMC,CTX,CPR	23	RL-fit
CTX,AMC,CTX,AMC,CTX,AMC,CTX,AMC,CTX,AMC	23	RL-genotype
CTX,AMC,CTX,AMC,CTX,CPR,CTX,AMC,CTX,CPR	96	RL-fit
CTX,SAM,CTX,SAM,CTX,CPR,CTX,CPR,CTX,CPR	96	RL-genotype

Table 2. Example drug sequences. Here, we show the first 10 selected drugs for representative episodes of the three top-performing replicates.

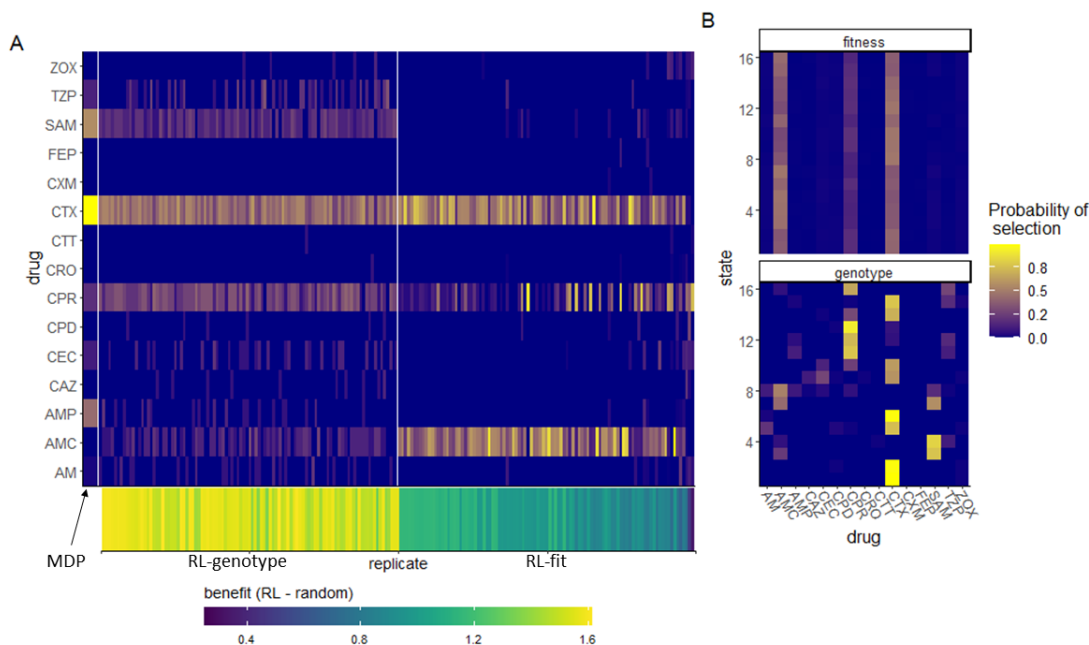


Figure 3. Drug cycling policies learned by RL-genotype and RL-fit. **A:** Heatmap depicting the learned policy for 100 replicates (on the x-axis) of the RL-genotype and 100 replicates of RL-fit. Far left column (enlarged) corresponds to the optimal policy derived from the MDP condition. The Y-axis describes the β -lactam antibiotics each RL agent could choose from while the color corresponds to the probability that the learned policy selected a given antibiotic. Bottom heatmap shows the median fitness benefit observed under the policy learned by a given replicate. **B:** Heatmap showing the average learned policy for RL-fit and RL-genotype. RL-genotype learns a more consistent mapping of state to action compared to RL-fit.

Overview of learned drug cycling policies for RL-fit and RL-genotype.

146

147 We evaluated the learned drug cycling policies of RL-fit, RL-genotype for the 15 β -lactam antibiotics under study. We
 148 compared these to the true optimal drug cycling policy as a reference. For this system, we show that the optimal drug cycling
 149 policy relies heavily on Cefotaxime, Ampicillin + Sulbactam, and Ampicillin (Fig 3A). Cefotaxime was used as treatment
 150 in more than 50% of time-steps, with Ampicillin + Sulbactam and Ampicillin used next most frequently. The optimal drug
 151 cycling policy used Cefprozil, Piperacillin + Tazobactam, and Cefaclor infrequently. The remaining drugs were not used at
 152 all. The different RL-fit replicates largely converged on a similar policy. They relied heavily on Cefotaxime and Amoxicillin
 153 + Clavulanic acid. However, they relied infrequently on Cefprozil. RL-genotype replicates also converged on a relatively
 154 conserved policy. Further, RL-genotype replicates showed a much more consistent mapping of state to action compared to
 155 RL-fit (Fig 3B). All optimization paradigms identified complex drug cycles that use 3 or more drugs to treat the evolving
 156 cell population. None of the tested two-drug combinations compete with policies learned by RL-genotype, and are generally
 157 out-performed by RL-fit. We show that policies that do not rely on Cefotaxime are suboptimal in this system. The three
 158 replicates that showed the least benefit compared to the random drug cycling case did not use Cefotaxime at all (Fig 3B). The
 159 importance of Cefotaxime is likely explained by the topography of the CTX drug landscape (Fig S6). More than half of the
 160 available genotypes in the CTX landscape lie in fitness valleys, providing ample opportunities to combine CTX with other
 161 drugs and "trap" the evolving population in low-fitness genotypes.

Evolutionary trajectories observed under RL-Genotype, RL-fit, and MDP drug policies.

162

163 Next, we compared the evolutionary paths taken by the simulated *E. coli* population under the MDP, RL-fit, RL-genotype,
 164 and random policy paradigms. The edge weights (corresponding to the probability of observed state transitions) of the
 165 RL-genotype and MDP landscapes show a 0.96 pearson correlation (Fig 4). In contrast, the edge weights of the RL-fit and MDP
 166 landscapes show a 0.82 pearson correlation (Fig 4). During the course of training the MDP condition, the backwards induction
 167 algorithm generated a value function $V(s, a)$ for all $s \in S$ and $a \in A$. In Figure 4D, we use this value function to show that
 168 certain genotypes (namely 1,5,6, and 13) were more advantageous to the evolving population than to the learner. These states
 169 were frequented much more often under the random drug cycling condition compared to any of the experimental conditions
 170 (Fig 4D). We also show that other genotypes (namely 12 and 11) were particularly advantageous for the learner compared to
 171 the evolving population. These states were frequented much more often under the experimental conditions compared to the
 172 random drug cycling condition (Fig 4D).

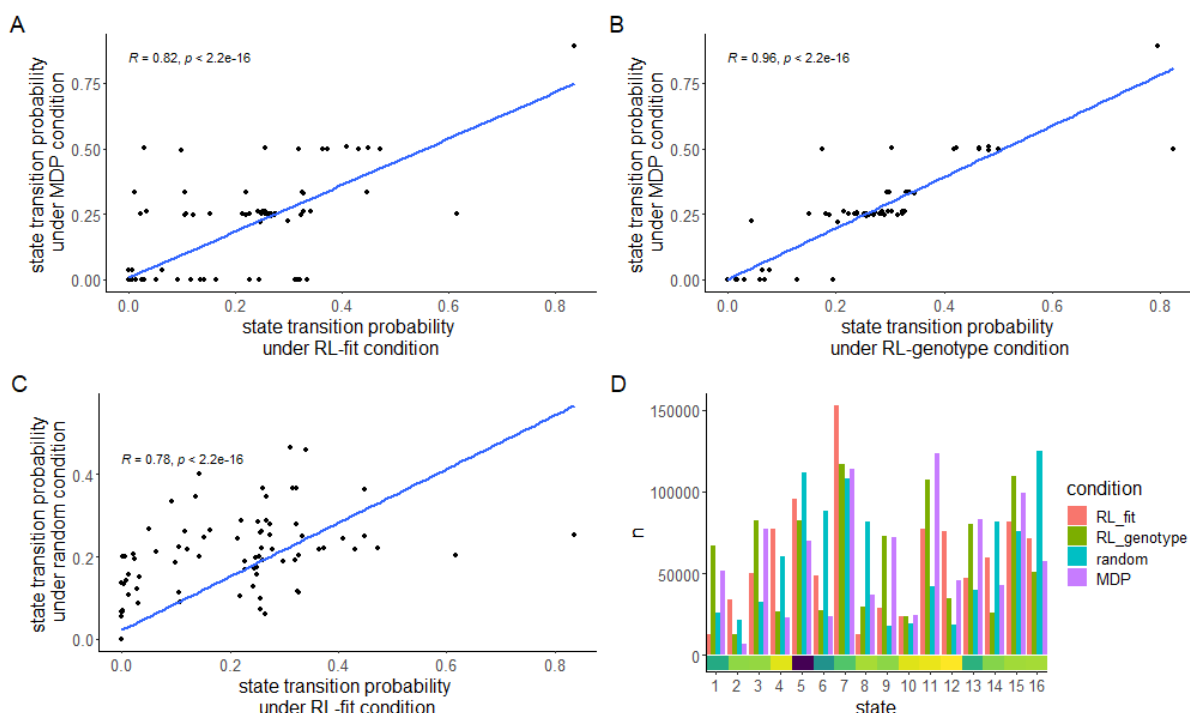


Figure 4. Comparison of evolutionary trajectories seen under different regimes A-C: Selected Pairwise comparisons of state transition frequency under different experimental conditions. State transition frequency is nearly identical for the RL-genotype and MDP conditions ($R=0.97$). In contrast, state-transition-frequency for the RL-fit and MDP conditions are related but less strongly correlated ($R=0.75$). As expected, state transition frequency is not very similar between the RL-fit and random conditions ($R=0.62$). **D:** Bar chart comparing the frequency that states are observed under different experimental conditions. The value of each state (to the learner) is highlighted for each state by the bottom heatmap. High value states are observed more frequently in RL-fit, RL-genotype, and MDP conditions compared to the random condition.

173 We also show that certain state transitions occur more frequently than others, independent of experimental conditions.
 174 For example, the population nearly always transitioned from genotype 5 to genotype 7 (**Fig 5**). This transition highlights the
 175 way these learned policies use drug landscapes to guide evolution. Genotype 5 (0100) is a fitness peak in most of the drug
 176 landscapes used in the learned policies, and is therefore a very disadvantageous state for the controlling agent. CTX, the most
 177 commonly used drug in all effective policies, has a slightly higher peak at genotype 7 (0110), which forces the population away
 178 from genotype 5 (**Fig S4**). As another example, the evolving population very rarely transitioned from state 1 to state 9 in the
 179 RL-fit condition. This state transition occurred commonly in the MDP and RL-genotype conditions (**Fig 5**). This difference is
 180 explained by the policies shown in **Figure 3B**. Under the RL-genotype policy, CTX was selected every time the population was
 181 in state 1 (the initial condition). The CTX landscape topography allows transition to 3 of the 4 single mutants, including state 9
 182 (1000) (**Fig S6**). Under the RL-fit policy, CTX and AMC were used in about equal proportion when the population is in state 1.
 183 Unlike the CTX landscape, the AMC landscape topography does not permit evolution from state 1 to state 9 (**Fig S6**).

184 Characteristics of selected drug policies

185 To better understand why certain drugs were used so frequently by RL-genotype, RL-fit, and the MDP policies, we
 186 developed the concept of an “opportunity landscape”. An opportunity landscape is an optimistic summation of n fitness
 187 landscapes. We computed each opportunity landscape by taking the minimum fitness value for each genotype from a given
 188 set of fitness landscapes. This simplified framework gives a sense of a potential best case scenario if the drugs in a given
 189 combination are used optimally. For example, the MDP policy relied heavily on CTX, CPR, AMP, SAM, and TZP to control
 190 the simulated *E. Coli* population. The resultant opportunity landscape (**Fig 6A**) contains only a single fitness peak, with 15/16
 191 of the genotypes in or near fitness valleys. In **Figure 6B**, we show the actual state transitions observed during evolution under
 192 the MDP policy. We also color the nodes based on the value function estimated by solving the MDP. As expected, the value
 193 function estimated by the MDP aligns closely with the topography of the opportunity landscape. There is only one genotype
 194 that the value function scores as being very poor for the learner, corresponding to the single peak in the opportunity fitness
 195 landscape (**Fig 6**). Interestingly, the opportunity landscape predicted that the population would evolve to the single fitness peak
 196 and fix. In contrast, the observed state transitions suggest that the MDP policy was able to guide the population away from that
 197 single fitness peak. A more detailed discussion of opportunity landscapes can be found in the supplemental materials.

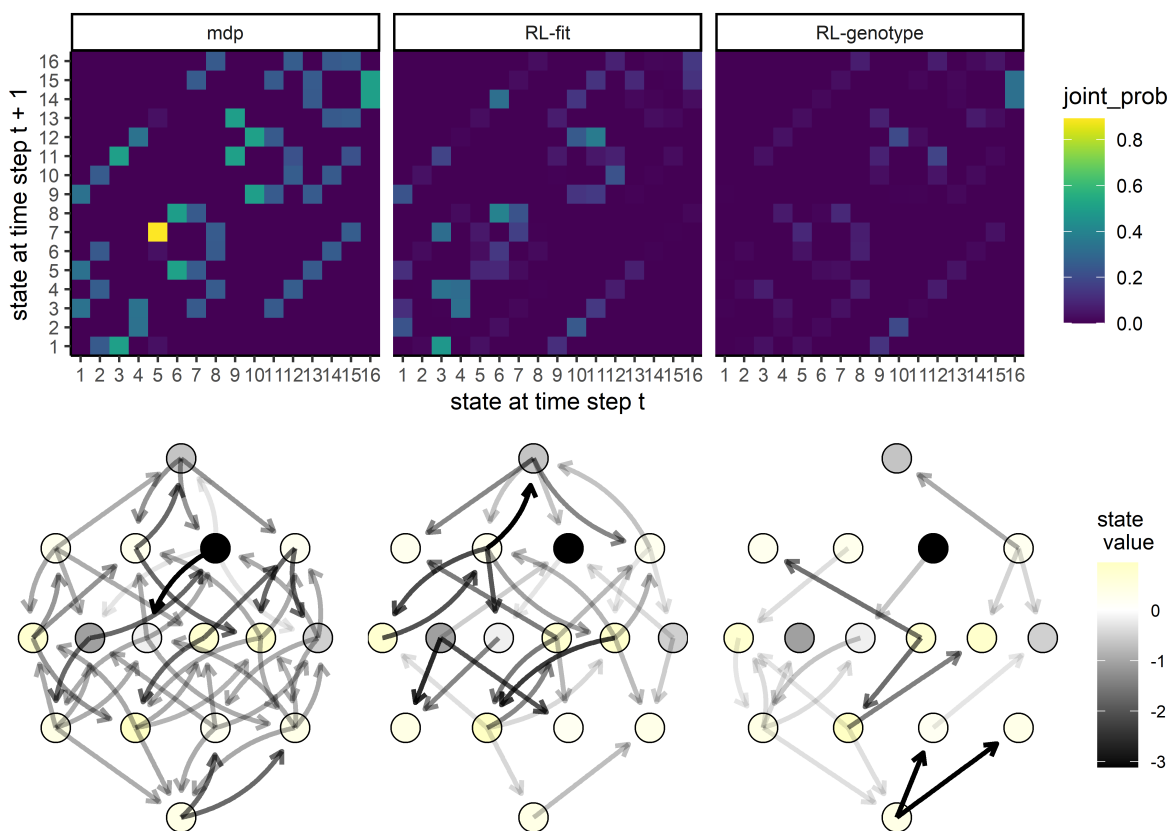


Figure 5. Movement of simulated *E. Coli* population through the genomic landscape. **Top row:** Heatmap depicting the joint probability distribution for each state transition under the different experimental conditions. The second two show the difference in state transition probability compared to the MDP condition. **Bottom row:** Graph depicting the fitness landscape, beginning with the wild type (bottom) all the way to the quadruple mutant (top). Size of arrow depicts the frequency with which a state transition was observed under the labeled experimental condition. The color of each node corresponds with the expected value (to the learner) of being in that state. As above, the second two arrows correspond to the observed difference between RL-Fit or RL-genotype and the MDP condition.

198 We also show that both the MDP and RL-genotype conditions select the drug with the lowest fitness for most genotypes (**Fig**
 199 **6C**). There are a few notable exceptions to this rule, which highlight RL-genotype’s capacity for treatment planning. A greedy
 200 policy that selects the lowest drug-fitness combination for every genotype would select Amoxicillin (AM) when the population
 201 is identified as being in genotype 5. The AM drug landscape then strongly favors transition back to the wild-type genotype
 202 (state 1). From state 1, most available drugs encourage evolution back to the genotype 5 fitness peak. As we see in **Figure 6B**,
 203 state 5 is by far the least advantageous for the learner. The greedy policy therefore creates an extremely disadvantageous cycle of
 204 evolution. In fact, none of the tested policies rely heavily on AM in state 5 (**Fig 3B**), instead taking a fitness penalty to select
 205 Cefotaxime (CTX). The CTX drug landscape encourages evolution to the double mutant, which has access to the highest value
 206 areas of the landscape. Finally, we rank drug landscapes based on the number of genotypes with a fitness value < 1 (**Fig 6D**).
 207 Based on the defined reward function, these genotypes would be considered advantageous to the learner. We show that drugs
 208 identified as useful by the optimal policy or RL-genotype tend to have more advantageous genotypes in their drug landscape.
 209 The only two highly permissive landscapes (CPD, CPR) that aren’t used have extremely similar topography to CTX, which
 210 most policies were built around.

211 3 Discussion

212 The evolution of widespread microbial drug resistance is driving a growing public health crisis around the world. In this study,
 213 we show a proof of concept for how existing drugs could be leveraged to control microbial populations without increasing drug
 214 resistance. To that end, we tested optimization approaches given decreasing amounts of information about an evolving system
 215 of *E. Coli*, and showed that it is possible to learn highly effective drug cycling policies given only empirically measurable
 216 information. To accomplish this, we developed a novel reinforcement learning approach to control an evolving population of *E.*
 217 *Coli in silico*. We focused on 15 empirically measured fitness landscapes pertaining to different clinically available β -lactam

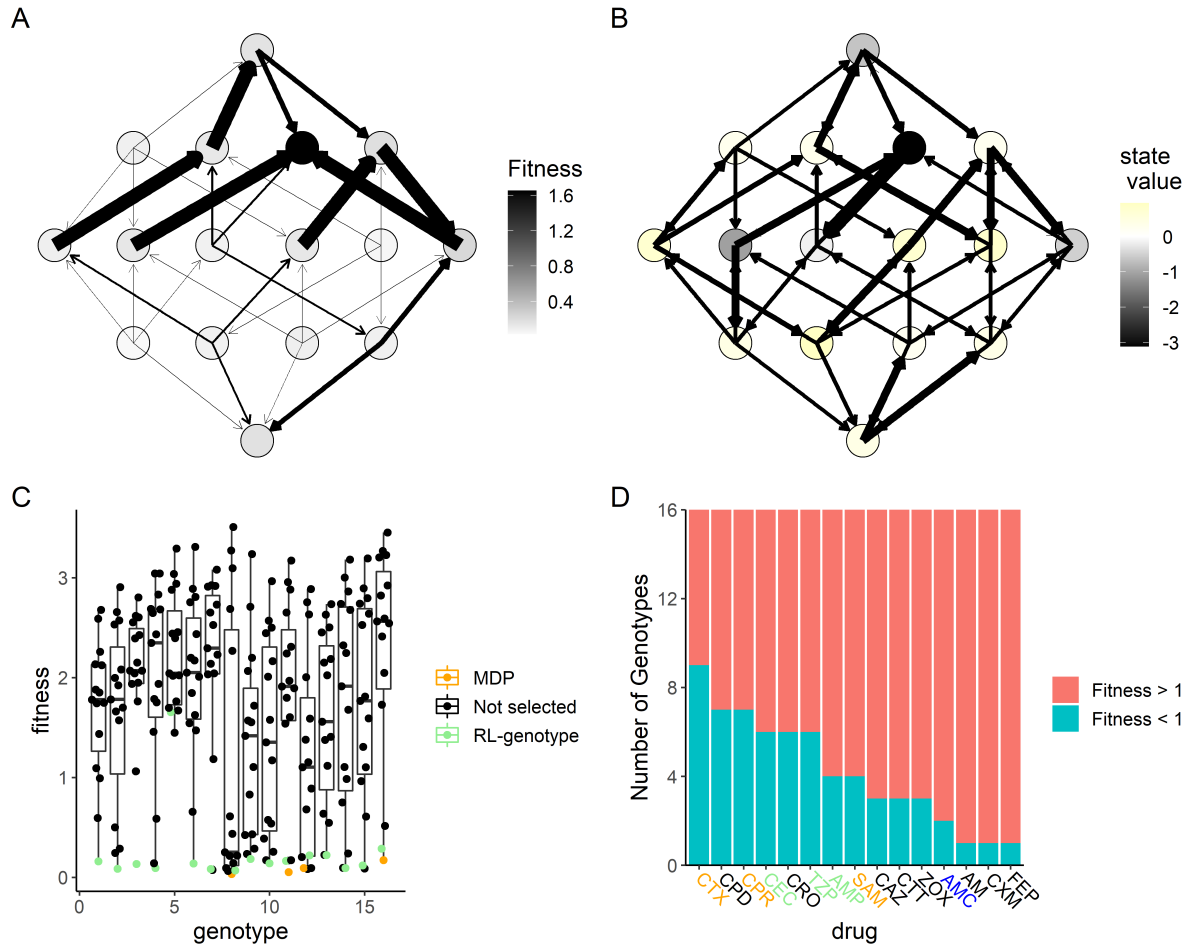


Figure 6. MDP value function closely matches opportunity landscape for drugs commonly used under MDP policy. Panels A and B show the 16-genotype fitness landscape under study, starting with the wild type at the top, progressing through the single mutants, double mutants, triple mutants, and finally the quadruple mutant at the bottom. **A:** Opportunity landscape for the 5 drugs most commonly used under the MDP policy (CTX, CPR, AMP, SAM, and TZP). **B:** Observed state transitions under the MDP policy. The node color corresponds to the value function estimated by solving the MDP. Lower values correspond to states the MDP policy attempts to avoid while higher values correspond to states the MDP policy attempts to steer the population. **C:** Scatter Plot showing the distribution of fitness with respect to genotype for the 15 β -lactam antibiotics under study. The drug selected by RL-genotype in a given genotype is highlighted in light blue. In cases where the MDP selected a different drug than RL-genotype, that drug is highlighted in orange. **D:** Number of genotypes with fitness above or below 1 for each drug under study. Drugs that are used by both the MDP and RL-genotype are highlighted in orange. Drugs that are used by only the MDP are highlighted in green. Drugs that are used by only RL-genotype are highlighted in blue.

218 antibiotics (**Table 1**). In this setting, RL agents selected treatments that, on average, controlled population fitness much more
219 effectively than either of the two negative controls. We showed that RL agents with access to the instantaneous genotype of the
220 population over time approach the MDP-derived optimal policy for these landscapes. Critically, we showed that RL agents were
221 capable of developing effective drug cycling protocols even when the measures of fitness used for training were first adjusted
222 by a noise parameter. This suggests that even imperfect measurements of an imperfect measure of population state (the kind of
223 measurements we are able to make in clinical settings) may be sufficient to develop effective control policies. We also show
224 that RL or MDP-derived policies consistently outperform simple alternating drug cycling policies. Finally, we introduced the
225 concept of the "Opportunity Landscape" which can provide powerful intuition into the viability of various drug combinations.

226 Our work expands a rich literature on the subject of evolutionary control through formal optimization approaches.
227 Our group and others have developed and optimized perfect information systems to generate effective drug cycling poli-
228 cies^{11, 12, 14, 16, 17}. Further, a limited number of studies have used RL-based methods for the development of clinical optimization
229 protocols^{20, 38–41, 43}. These studies have been limited so far to contrived simulated systems, including a recent study that
230 introduced Cellulose, a RL framework capable of controlling evolving bacterial populations in a stochastic simulated system⁴².

231 Much like the studies noted above, we show that AI or MDP-based policies for drug selection or drug dosing dramatically
232 outperform sensible controls in the treatment of an evolving cell population. We also extend this literature in two key ways. To
233 our knowledge, ours is the first optimization protocol capable of learning effective drug cycling policies using only observed
234 population fitness (a clinically tractable measure) as the key training input. Second, we grounded our work with empirically
235 measured fitness landscapes which will facilitate more natural extension to the bench.

236 There are several limitations to this work which bear mention. We assume that selection under drug therapy represents
237 a strong-selection and weak mutation regime in order to compute transition matrices for our models. While this is likely
238 true in most cases, it is possible that other selection regimes emerge in cases of real world pharmacokinetics where the drug
239 concentration fluctuates dramatically. In addition, we chose to keep drug concentration constant throughout are analysis, largely
240 owing to the lack of robust empirical data linking genotype to phenotype under dose varying conditions (sometimes called a
241 fitness seascape)⁴⁶. As more empirical fitness seascape data becomes available, a natural extension would be to explore the
242 efficacy of the RL system in controlling a population by varying both drug and dose.

243 While we present the most extensive genotype-phenotype modeling work to date on this subject, we still only modeled the
244 effect of mutations at 4 genotypic positions. The real *E. Coli* genome is approximately 5×10^6 base pairs⁴⁷. The evolutionary
245 landscape for living organisms is staggeringly large, and not tractable to model *in silico*. It is possible that empirical measures
246 of fitness like growth rate or cell count may not provide a robust enough signal of the underlying evolutionary state on real
247 genomes. *In vitro* implementations of reinforcement learning-based drug cycle optimization systems are needed to address this
248 potential shortcoming. Another potential alternative would be to use the comparatively low-dimensional phenotype landscape
249 of drug resistance⁴⁸.

250 In this work, we present a novel reinforcement-learning framework capable of controlling an evolving population of *E. Coli*
251 *in silico*. We show that RL agents stably learn multi-drug combinations that were state specific and reliably out-performed a
252 random drug cycling policy as well as all possible two-drug cycling policies. We also highlight key features of the types of drug
253 landscapes that are useful for the design of evolutionary control policies. Our work represents an important proof-of-concept
254 for AI-based evolutionary control, an emerging field with the potential to revolutionize clinical medicine.

255 Acknowledgements

256 This work was made possible by the National Institute of Health (5R37CA244613-03, 5T32GM007250-46, and T32CA094186)
257 and American Cancer Society (RSG-20-096-01). Figure 1 was created with BioRender.com.

References

- 258 **1.** Murray, C. J. *et al.* Global burden of bacterial antimicrobial resistance in 2019: a systematic analysis. *The Lancet* **399**,
259 629–655, DOI: [10.1016/S0140-6736\(21\)02724-0](https://doi.org/10.1016/S0140-6736(21)02724-0) (2022). Publisher: Elsevier.
- 260 **2.** Centers for Disease Control and Prevention (U.S.). Antibiotic resistance threats in the United States, 2019. Tech. Rep.,
261 Centers for Disease Control and Prevention (U.S.) (2019). DOI: [10.15620/cdc:82532](https://doi.org/10.15620/cdc:82532).
- 262 **3.** Plackett, B. Why big pharma has abandoned antibiotics. *Nature* **586**, DOI: [10.1038/d41586-020-02884-3](https://doi.org/10.1038/d41586-020-02884-3) (2020).
- 263 **4.** Stearns, S. C. Evolutionary medicine: Its scope, interest and potential. *Proc. Royal Soc. B: Biol. Sci.* **279**, 4305–4321,
264 DOI: [10.1098/rspb.2012.1326](https://doi.org/10.1098/rspb.2012.1326) (2012).
- 265 **5.** Grunspan, D. Z., Nesse, R. M., Barnes, M. E. & Brownell, S. E. Core principles of evolutionary medicine. *Evol. Medicine,*
266 *Public Heal.* **2018**, 13–23, DOI: [10.1093/emph/eox025](https://doi.org/10.1093/emph/eox025) (2018).
- 267 **6.** Perry, G. H. Evolutionary medicine. *eLife* **10**, e69398, DOI: [10.7554/eLife.69398](https://doi.org/10.7554/eLife.69398) (2021).
- 268 **7.** Andersson, D. I. *et al.* Antibiotic resistance: turning evolutionary principles into clinical reality. *FEMS Microbiol. Rev.* **44**,
269 171–188 (2020). Publisher: Oxford University Press.
- 270 **8.** Manrubia, S. *et al.* From genotypes to organisms: State-of-the-art and perspectives of a cornerstone in evolutionary
271 dynamics. *Phys. Life Rev.* **38**, 55–106, DOI: [10.1016/j.plrev.2021.03.004](https://doi.org/10.1016/j.plrev.2021.03.004) (2021).
- 272 **9.** Stracy, M. *et al.* Minimizing treatment-induced emergence of antibiotic resistance in bacterial infections. *Science* **375**,
273 889–894, DOI: [10.1126/science.abg9868](https://doi.org/10.1126/science.abg9868) (2022). Publisher: American Association for the Advancement of Science.
- 274 **10.** Baym, M., Stone, L. K. & Kishony, R. Multidrug evolutionary strategies to reverse antibiotic resistance. *Science* **351**,
275 aad3292 (2016). Publisher: American Association for the Advancement of Science.
- 276 **11.** Yoon, N., Vander Velde, R., Marusyk, A. & Scott, J. G. Optimal Therapy Scheduling Based on a Pair of Collaterally
277 Sensitive Drugs. *Bull. Math. Biol.* **80**, 1776–1809, DOI: [10.1007/s11538-018-0434-2](https://doi.org/10.1007/s11538-018-0434-2) (2018).
- 278 **12.** Maltas, J. & Wood, K. B. Pervasive and diverse collateral sensitivity profiles inform optimal strategies to limit antibiotic
279 resistance. *PLoS biology* **17**, e3000515, DOI: [10.1371/journal.pbio.3000515](https://doi.org/10.1371/journal.pbio.3000515) (2019).
- 280 **13.** Maltas, J. & Wood, K. B. Dynamic collateral sensitivity profiles highlight challenges and opportunities for optimizing
281 antibiotic sequences. *bioRxiv* (2021).
- 282 **14.** Gluzman, M., Scott, J. G. & Vladimirovsky, A. Optimizing adaptive cancer therapy: dynamic programming and evolutionary
283 game theory. *Proc. Royal Soc. B: Biol. Sci.* **287**, 20192454, DOI: [10.1098/rspb.2019.2454](https://doi.org/10.1098/rspb.2019.2454) (2020). Publisher: Royal
284 Society.
- 285 **15.** Nichol, D. *et al.* Steering Evolution with Sequential Therapy to Prevent the Emergence of Bacterial Antibiotic Resistance.
286 *PLOS Comput. Biol.* **11**, e1004493, DOI: [10.1371/journal.pcbi.1004493](https://doi.org/10.1371/journal.pcbi.1004493) (2015). Publisher: Public Library of Science.
- 287 **16.** Yoon, N., Krishnan, N. & Scott, J. Theoretical modeling of collaterally sensitive drug cycles: shaping heterogeneity to
288 allow adaptive therapy. *J. Math. Biol.* **83**, 47, DOI: [10.1007/s00285-021-01671-6](https://doi.org/10.1007/s00285-021-01671-6) (2021).
- 289 **17.** Iram, S. *et al.* Controlling the speed and trajectory of evolution with counterdiabatic driving. *Nat. Phys.* **17**, 135–
290 142, DOI: [10.1038/s41567-020-0989-3](https://doi.org/10.1038/s41567-020-0989-3) (2021). Bandiera_abtest: a Cg_type: Nature Research Journals Number: 1
291 Primary_atype: Research Publisher: Nature Publishing Group Subject_term: Biophysics;Statistical physics;Theoretical
292 physics Subject_term_id: biophysics;statistical-physics;theoretical-physics.
- 293 **18.** Maltas, J., Singleton, K. R., Wood, K. C. & Wood, K. B. Drug dependence in cancer is exploitable by optimally constructed
294 treatment holidays. *bioRxiv* (2022).
- 295 **19.** Chakrabarti, S. & Michor, F. Pharmacokinetics and drug interactions determine optimum combination strategies in
296 computational models of cancer evolution. *Cancer Res.* **77**, 3908–3921, DOI: [10.1158/0008-5472.CAN-16-2871](https://doi.org/10.1158/0008-5472.CAN-16-2871) (2017).
- 297 **20.** Newton, P. K. & Ma, Y. Nonlinear adaptive control of competitive release and chemotherapeutic resistance. *Phys. Rev. E*
298 **99**, 022404, DOI: [10.1103/PhysRevE.99.022404](https://doi.org/10.1103/PhysRevE.99.022404) (2019).
- 299 **21.** Kim, S., Lieberman, T. D. & Kishony, R. Alternating antibiotic treatments constrain evolutionary paths to multidrug
300 resistance. *Proc. Natl. Acad. Sci.* **111**, 14494–14499, DOI: [10.1073/pnas.1409800111](https://doi.org/10.1073/pnas.1409800111) (2014). Publisher: Proceedings of
301 the National Academy of Sciences.
- 302 **22.** Zhang, J., Cunningham, J. J., Brown, J. S. & Gatenby, R. A. Integrating evolutionary dynamics into treatment of metastatic
303 castrate-resistant prostate cancer. *Nat. Commun.* **8**, DOI: [10.1038/s41467-017-01968-5](https://doi.org/10.1038/s41467-017-01968-5) (2017).
- 304

- 305 **23.** Cunningham, J. J., Brown, J. S., Gatenby, R. A. & Staňková, K. Optimal control to develop therapeutic strategies for
306 metastatic castrate resistant prostate cancer. *J. Theor. Biol.* **459**, 67–78, DOI: [10.1016/j.jtbi.2018.09.022](https://doi.org/10.1016/j.jtbi.2018.09.022) (2018).
- 307 **24.** Weinreich, D. M., Watson, R. A. & Chao, L. Perspective: Sign Epistasis and Genetic Constraint on Evo-
308 lutionary Trajectories. *Evolution* **59**, 1165–1174, DOI: [10.1111/j.0014-3820.2005.tb01768.x](https://doi.org/10.1111/j.0014-3820.2005.tb01768.x) (2005). [_eprint:
309 https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.0014-3820.2005.tb01768.x](https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.0014-3820.2005.tb01768.x).
- 310 **25.** Mira, P. M. *et al.* Rational Design of Antibiotic Treatment Plans: A Treatment Strategy for Managing Evolution and
311 Reversing Resistance. *PLOS ONE* **10**, e0122283, DOI: [10.1371/journal.pone.0122283](https://doi.org/10.1371/journal.pone.0122283) (2015). Publisher: Public Library
312 of Science.
- 313 **26.** Maltas, J., McNally, D. M. & Wood, K. B. Evolution in alternating environments with tunable inter-landscape correlations.
314 *Evol. international journal organic evolution* **75**, 10–24, DOI: [10.1111/evo.14121](https://doi.org/10.1111/evo.14121) (2021).
- 315 **27.** de Visser, J. A. G. M. & Krug, J. Empirical fitness landscapes and the predictability of evolution. *Nat. Rev. Genet.* **15**,
316 480–490, DOI: [10.1038/nrg3744](https://doi.org/10.1038/nrg3744) (2014). Bandiera_abtest: a Cg_type: Nature Research Journals Number: 7 Primary_atype:
317 Reviews Publisher: Nature Publishing Group Subject_term: Epistasis;Evolutionary genetics;Experimental evolution
318 Subject_term_id: epistasis;evolutionary-genetics;experimental-evolution.
- 319 **28.** Weinreich, D. M., Delaney, N. F., Depristo, M. A. & Hartl, D. L. Darwinian evolution can follow only very few mutational
320 paths to fitter proteins. *Sci. (New York, N.Y.)* **312**, 111–114, DOI: [10.1126/science.1123539](https://doi.org/10.1126/science.1123539) (2006).
- 321 **29.** Ogbunugafor, C. B., Wylie, C. S., Diakite, I., Weinreich, D. M. & Hartl, D. L. Adaptive Landscape by Environment
322 Interactions Dictate Evolutionary Dynamics in Models of Drug Resistance. *PLoS Comput. Biol.* **12**, 1–20, DOI: [10.1371/
323 journal.pcbi.1004710](https://doi.org/10.1371/journal.pcbi.1004710) (2016).
- 324 **30.** Toprak, E. *et al.* Evolutionary paths to antibiotic resistance under dynamically sustained drug selection. *Nat. genetics* **44**,
325 101–105 (2012). Publisher: Nature Publishing Group.
- 326 **31.** Greenbury, S. F., Louis, A. A. & Ahnert, S. E. The structure of genotype-phenotype maps makes fitness landscapes
327 navigable. *Nat. Ecol. & Evol.* 1–11, DOI: [10.1038/s41559-022-01867-z](https://doi.org/10.1038/s41559-022-01867-z) (2022). Publisher: Nature Publishing Group.
- 328 **32.** Baym, M. *et al.* Spatiotemporal microbial evolution on antibiotic landscapes. *Sci. (New York, N.Y.)* **353**, 1147–1151, DOI:
329 [10.1126/science.aag0822](https://doi.org/10.1126/science.aag0822) (2016).
- 330 **33.** Sutton, R. S. & Barto, A. G. *Reinforcement learning: An introduction* (MIT press, 2018).
- 331 **34.** Du Plessis, L., Leventhal, G. E. & Bonhoeffer, S. How good are statistical models at approximating complex fitness
332 landscapes? *Mol. biology evolution* **33**, 2454–2468 (2016). Publisher: Oxford University Press.
- 333 **35.** Silver, D. *et al.* Mastering the game of Go with deep neural networks and tree search. *Nature* **529**, DOI: [10.1038/nature16961](https://doi.org/10.1038/nature16961)
334 (2016).
- 335 **36.** Mnih, V. *et al.* Human-level control through deep reinforcement learning. *Nature* **518**, 529–533, DOI: [10.1038/nature14236](https://doi.org/10.1038/nature14236)
336 (2015). Bandiera_abtest: a Cg_type: Nature Research Journals Number: 7540 Primary_atype: Research Publisher: Nature
337 Publishing Group Subject_term: Computer science Subject_term_id: computer-science.
- 338 **37.** Vinyals, O. *et al.* Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature* **575**, 350–354,
339 DOI: [10.1038/s41586-019-1724-z](https://doi.org/10.1038/s41586-019-1724-z) (2019). Bandiera_abtest: a Cg_type: Nature Research Journals Number: 7782
340 Primary_atype: Research Publisher: Nature Publishing Group Subject_term: Computer science;Statistics Subject_term_id:
341 computer-science;statistics.
- 342 **38.** Moore, B. L. *et al.* Reinforcement Learning for Closed-Loop Propofol Anesthesia: A Study in Human Volunteers. *J. Mach.*
343 *Learn. Res.* **15**, 655–696 (2014).
- 344 **39.** Petersen, B. K. *et al.* Deep Reinforcement Learning and Simulation as a Path Toward Precision Medicine. *J. Comput. Biol.*
345 *A J. Comput. Mol. Cell Biol.* **26**, 597–604, DOI: [10.1089/cmb.2018.0168](https://doi.org/10.1089/cmb.2018.0168) (2019).
- 346 **40.** Padmanabhan, R., Meskin, N. & Haddad, W. M. Reinforcement learning-based control of drug dosing for cancer
347 chemotherapy treatment. *Math. Biosci.* **293**, 11–20, DOI: [10.1016/j.mbs.2017.08.004](https://doi.org/10.1016/j.mbs.2017.08.004) (2017).
- 348 **41.** Ahn, I. & Park, J. Drug scheduling of cancer chemotherapy based on natural actor-critic approach. *Bio Syst.* **106**, 121–129,
349 DOI: [10.1016/j.biosystems.2011.07.005](https://doi.org/10.1016/j.biosystems.2011.07.005) (2011).
- 350 **42.** Engelhardt, D. Dynamic Control of Stochastic Evolution: A Deep Reinforcement Learning Approach to Adaptively
351 Targeting Emergent Drug Resistance. *J. Mach. Learn. Res.* **21**, 1–30 (2020).
- 352 **43.** Martin, R. B. Optimal control drug scheduling of cancer chemotherapy. *Automatica* **28**, 1113–1123, DOI: [10.1016/
353 0005-1098\(92\)90054-J](https://doi.org/10.1016/0005-1098(92)90054-J) (1992).

- 354 **44.** Kallenberg, L. *Lecture Notes Markov Decision Problems - version 2020* (2020).
- 355 **45.** Arulkumaran, K., Deisenroth, M. P., Brundage, M. & Bharath, A. A. Deep reinforcement learning: A brief survey. *IEEE*
356 *Signal Process. Mag.* **34**, 26–38, DOI: [10.1109/MSP.2017.2743240](https://doi.org/10.1109/MSP.2017.2743240) (2017).
- 357 **46.** King, E. S. *et al.* Fitness seascapes facilitate the prediction of therapy resistance under time-varying selection, DOI:
358 [10.1101/2022.06.10.495696](https://doi.org/10.1101/2022.06.10.495696) (2022). Pages: 2022.06.10.495696 Section: New Results.
- 359 **47.** Rode, C. K., Melkerson-Watson, L. J., Johnson, A. T. & Bloch, C. A. Type-Specific Contributions to Chromosome Size
360 Differences in *Escherichia coli*. *Infect. Immun.* **67**, 230–236 (1999).
- 361 **48.** Iwasawa, J. *et al.* Analysis of the evolution of resistance to multiple antibiotics enables prediction of the *Escherichia coli*
362 phenotype-based fitness landscape. *PLOS Biol.* **20**, e3001920, DOI: [10.1371/journal.pbio.3001920](https://doi.org/10.1371/journal.pbio.3001920) (2022). Publisher:
363 Public Library of Science.

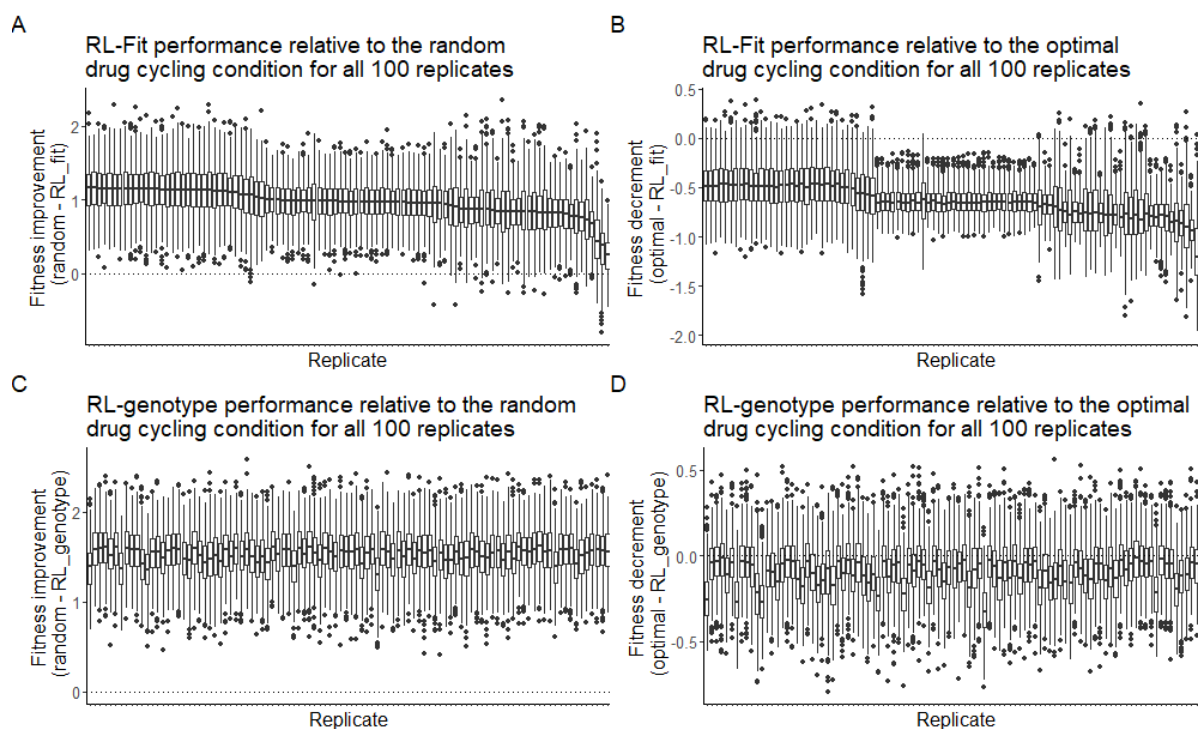


Figure S1. Performance of RL-fit and RL-genotype for each replicate. **A:** Fitness observed under RL-fit policy compared to random drug cycling condition. **B:** Fitness observed under RL-fit policy compared to fitness observed under optimal policy. **C:** Fitness observed under RL-genotype policy compared to random drug cycling condition. **D:** Fitness observed under RL-genotype policy compared to fitness observed under optimal policy.

Supplemental Materials

1.1 Hyperparameter tuning

We varied key hyperparameters one at a time in order to identify optimal values to promote learning in this setting. Parameter ranges and the selected value are shown in **Table S1**. Due to the long run-times of the training process, we were unable to make use of more formal hyper-parameter optimization approaches. Future work will increase the efficiency of training reinforcement learners in this setting, opening up a number of interesting follow-on studies.

1.2 Additional performance data for RL agents

As mentioned in the main text, we tested both the RL-fit and RL-genotype conditions 100 times each. In **Figure S1**, we show the performance of all 100 RL-fit and RL-genotype replicates. In 98/100 replicates, RL-fit outperformed the random drug cycling case (**Fig S1A**). The very best RL-fit replicates still fell short of the MDP-derived optimal policy (**Fig S1B**). In all 100 replicates, RL-genotype outperformed the random drug cycling case (**Fig S1C**). RL-genotype performance approached the performance of the optimal policy (**Fig S1D**).

Table S1. Key Hyperparameters for reinforcement learner

Parameter	Value	Range
gamma	0.99	0-1
learning rate	0.0001	0.000001-0.1
minibatch size	60	20-500
update target model frequency	310	100-1000

376 1.3 Policy Clustering

377 We performed PCA on the policies for all 100 replicates from the RL-genotype and RL-fit conditions. We then used the
378 silhouette method, implemented in the factoextra package, to estimate the appropriate number of clusters. We found that either
379 2 or 5 clusters would be optimal. As two clusters would only recapitulate our original RL-genotype and RL-fit conditions, we
380 performed kmeans clustering with 5 centers. Results are shown in **Fig S2**.

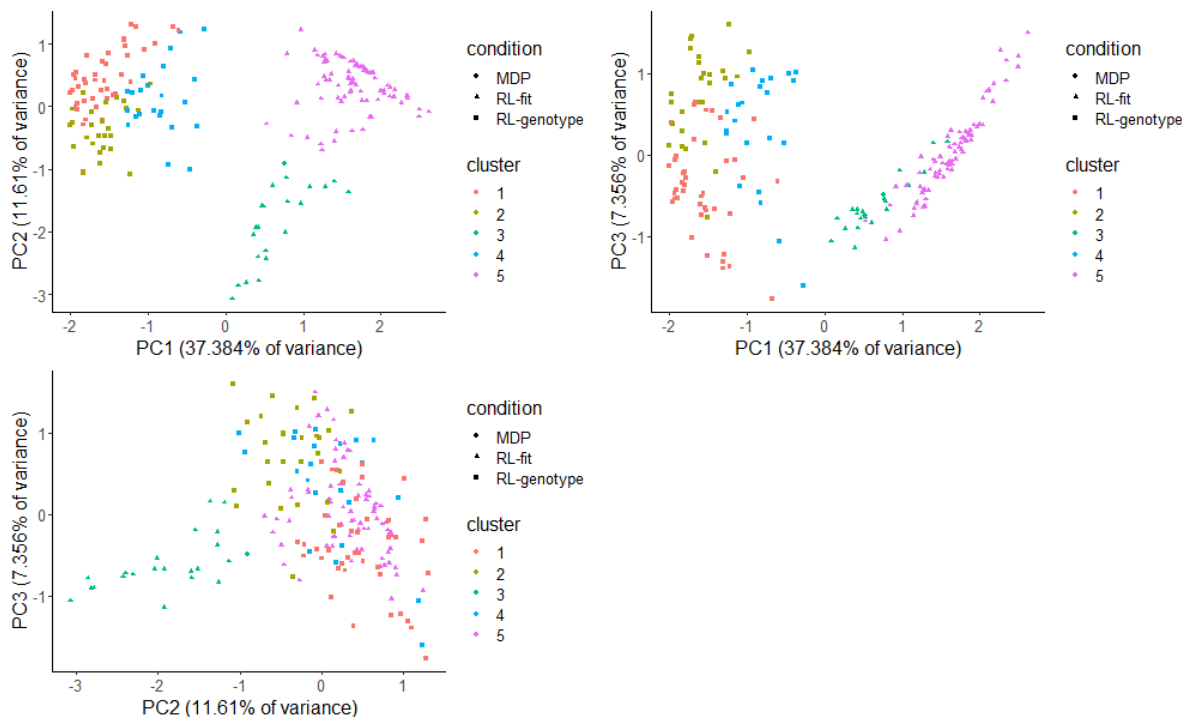


Figure S2. Identification of policy groups using PCA and kmeans clustering. RL-fit replicates were separated into 2 distinct replicates. RL-genotype replicates were separated into 3 distinct groups but it is unclear if these groups are meaningfully different.

381 We plotted the top 3 pairwise comparisons of principal components, which together account for about 56.4% of the variance
382 in this dataset. The RL-fit and RL-genotype policies were clearly separated by this method, with RL-genotype being split
383 between clusters 1,2, and 4. RL-fit was split between clusters 3 and 5 (**Fig S2**). Clusters 3 and 5 represent meaningfully
384 different policy motifs that RL-fit found frequently over the course of 100 replicates. Cluster 3 replicates were Cefprozil
385 dominant, and used Cefotaxime and Amoxicillin + Clavulanic acid infrequently. Cluster 5 replicates were Cefotaxime and
386 Amoxicillin + Clavulanic acid dominant, and used Cefprozil infrequently (**Fig S3**). Notably, cluster 3 replicates tended to have
387 worse performance compared to cluster 5 replicates (**Fig 3**).

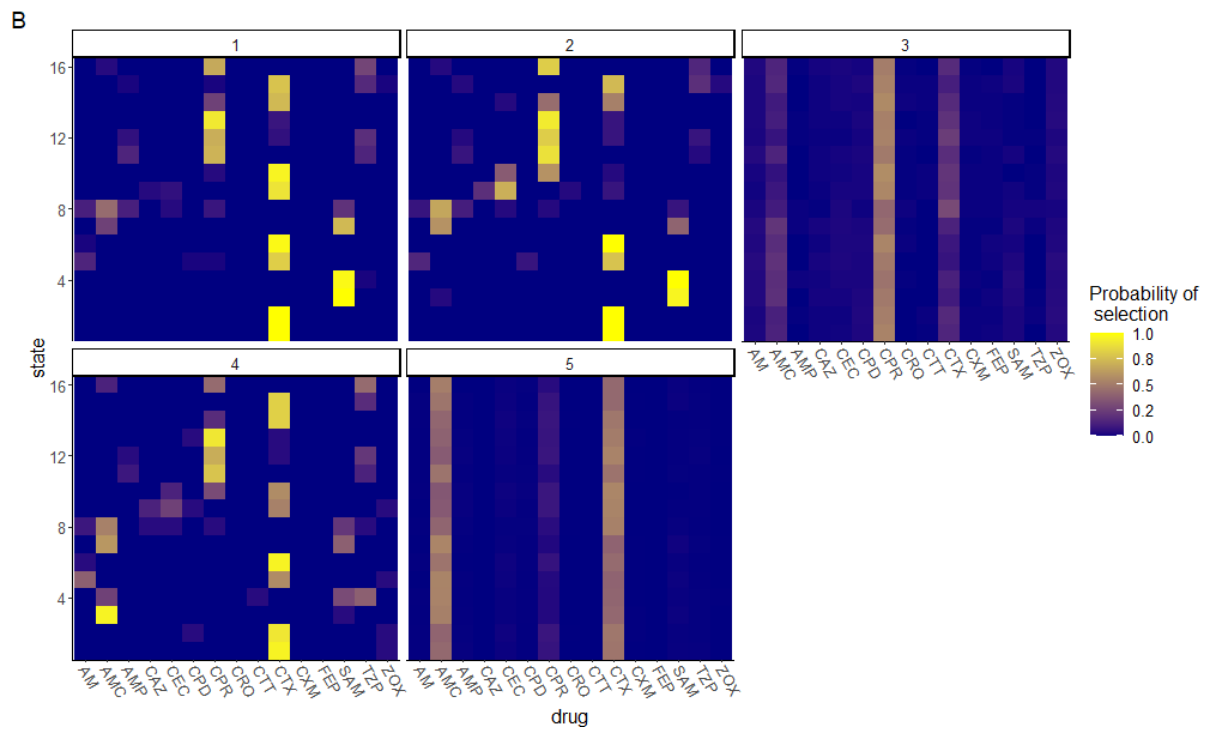


Figure S3. policy heatmaps for groups identified using PCA and kmeans clustering. Groups 1,2, and 4 correspond to RL-genotype policies. Groups 3 and 5 correspond to RL-fit policies. color gradient represents probability that a given drug (x-axis) will be selected when population is in a given state (y-axis).

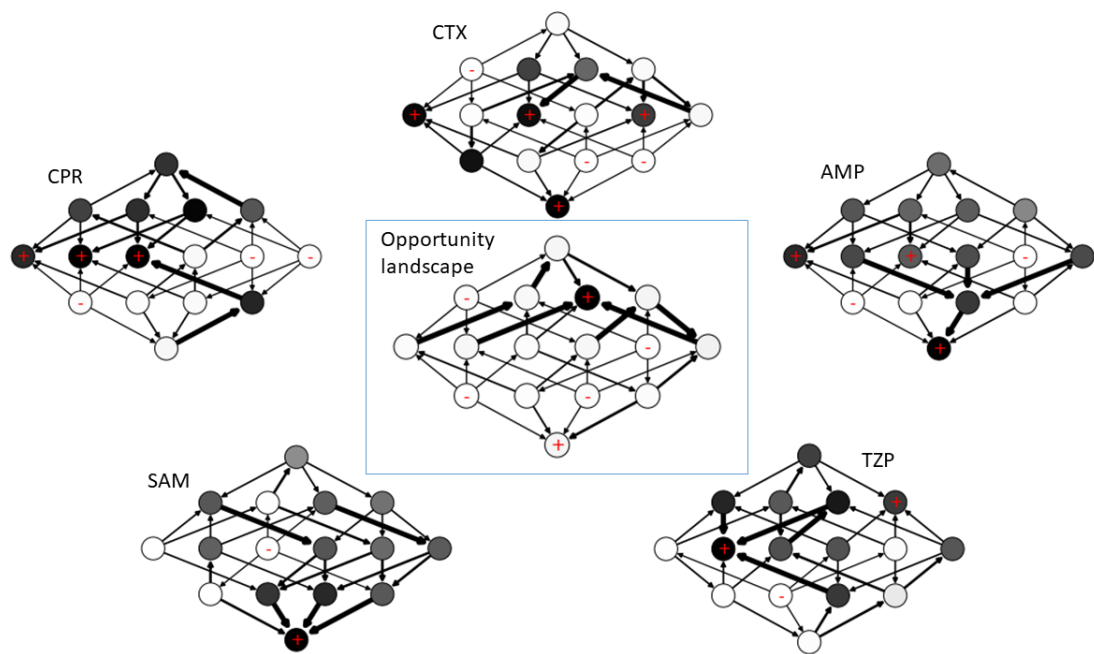


Figure S4. Opportunity Landscape for MDP-derived policy. Opportunity landscape is an optimistic combination of 5 empirically measured drug landscapes. Just 1/16 genotypes is near a fitness peak on the opportunity landscape, helping to explain the extremely low fitness observed in the simulated E.Coli population when the MDP-derived policy is applied.

1.4 Opportunity Landscapes

388
389
390
391
392
393
394
395

We define an opportunity landscape to be the most optimistic combination of n landscapes, formed by taking the minimum possible fitness at each genotypic position. This construct can help us better understand how the learner uses different combinations of drugs to maintain the evolving population at extremely low fitness values. Figure S4 describes the opportunity landscape discovered by the MDP condition. As noted in the main text, the MDP primarily uses 5 drugs (CTX, CPR, AMP, SAM, and TZP) in combination to trap the evolving population of E. Coli at extremely low fitness genotypes. In the combined opportunity landscape, just one genotype (0100) had a high fitness in all 5 drugs. As expected, the opportunity landscape closely matches the value function estimated by the MDP (Fig 5)

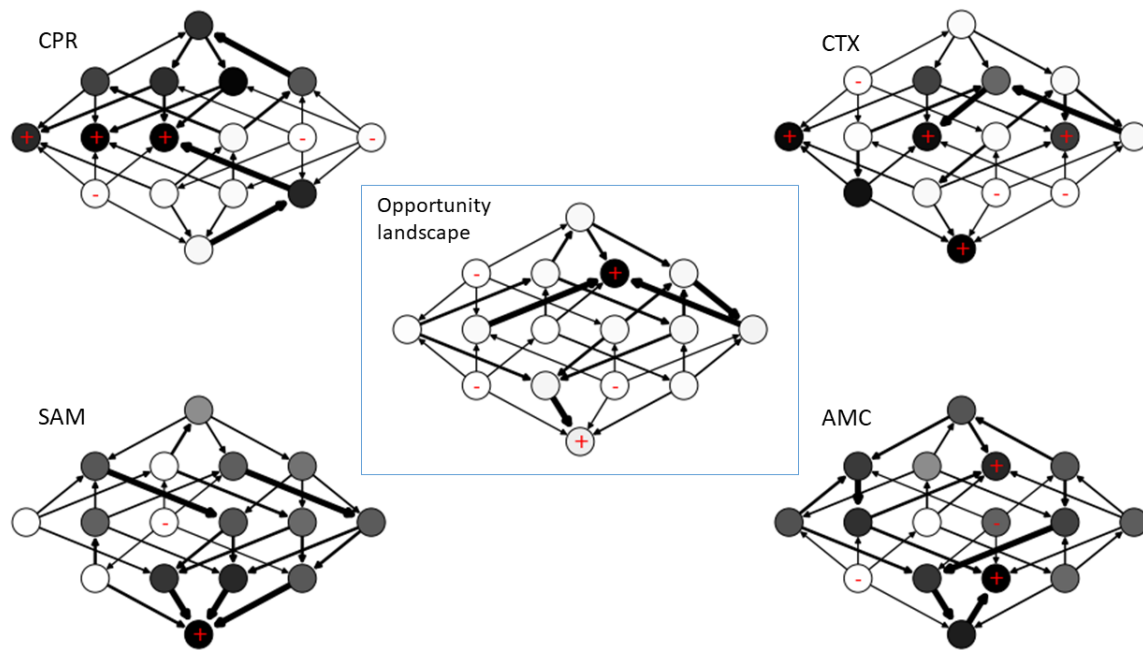


Figure S5. Opportunity landscape for most common policy identified in the RL-genotype condition. As in the MDP-derived policy, just 1/16 genotypes is near a fitness peak in the opportunity landscape.

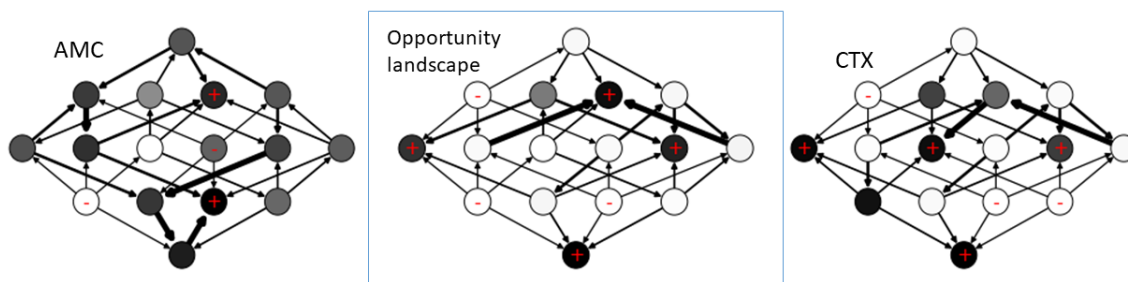


Figure S6. Opportunity landscape for the most common policy identified in the RL-fit condition. The most common RL-fit policy relies on AMC and CTX to control the E. Coli population. Assuming the most optimistic combination of these two drug landscapes, 4/16 genotypes are near a fitness peak.

396 The opportunity landscape for the RL-genotype is almost identical to the opportunity landscape observed for the MDP
397 policy (Fig S5). Interestingly, RL-genotype only uses 3 of the 5 drugs in the MDP policy; CPR, CTX, and SAM. RL-fit
398 discovered policies that typically only used two drugs. The most effective RL-fit policies relied heavily on AMC and CTX. We
399 present the resulting opportunity landscape in Figure S6. As expected, there are more genotypes with high fitness values under
400 this two-drug paradigm compared to the 4 or 5 drug policies discovered by RL-genotype and the MDP, respectively.

401 1.5 MDP policy

402 As mentioned in the main text, we computed the MDP policy by formulating a markov decision process of the strong selection,
403 weak mutation model of evolution under study. We then solved the MDP using backward induction, an algorithm designed
404 to identify an optimal policy for a finite time discrete MDP. The identified policy is a function of current state and current
405 time step, making it even more specific than the policies identified by the reinforcement learning conditions. We show the
406 time and state-specific MDP policy in Figure S7. Near the end of an episode (steps 19 and 20), we see a switch to a greedy
407 policy that simply selects the drug with the minimum fitness for a given genotype. We also varied the discount rate (gamma),
408 between 0 and 1 during the hyperparameter tuning process. In Figure S8, we show the effect of gamma on the average fitness
409 achieved by the MDP policy. While gamma didn't have a large effect, likely due to the relatively short length (20 time steps) of
410 each episode, we show that increasing gamma led to increased performance of the computed MDP policy. We also show that
411 increasing gamma led to increased use of CTX (drug 4).

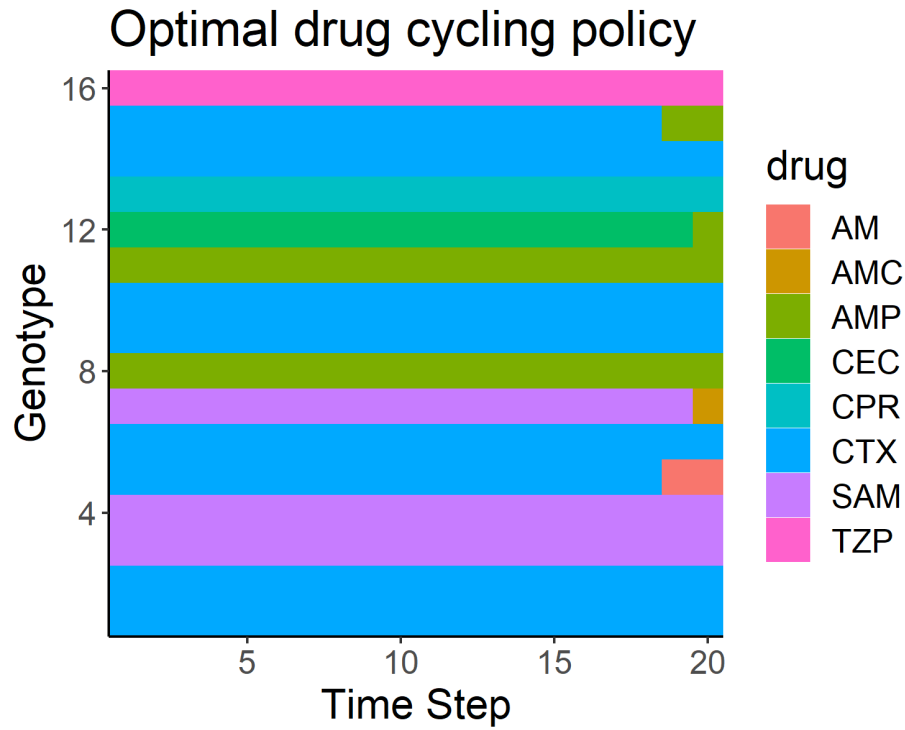


Figure S7. MDP-derived optimal policy.

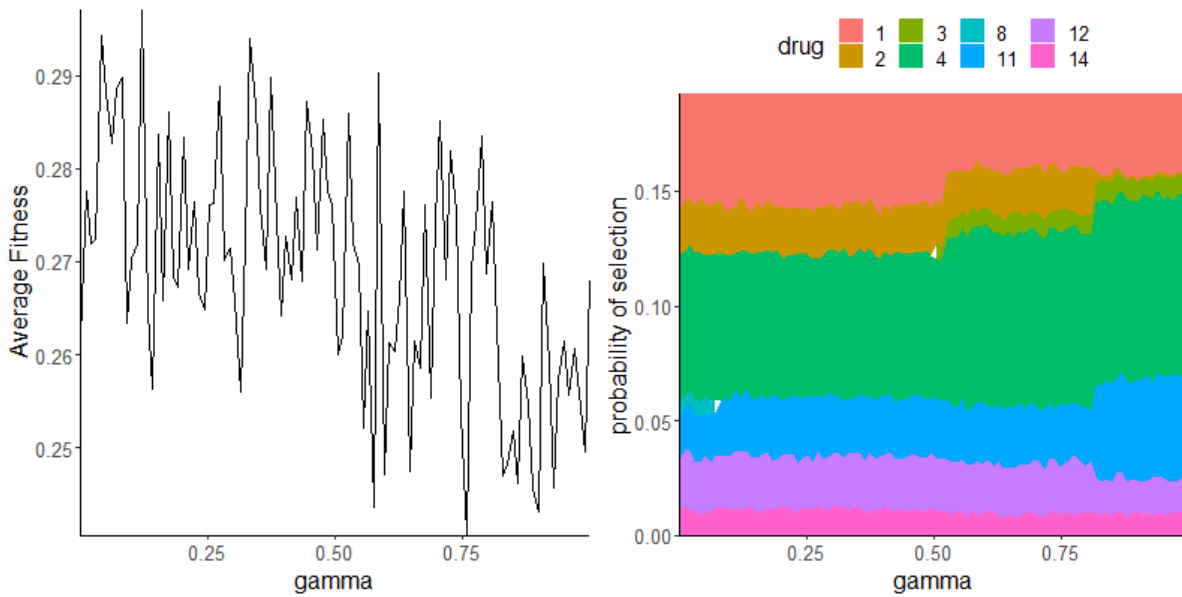


Figure S8. Effect of variation in gamma on optimal policy performance and composition

412 1.6 Additional Analyses

413 As noted in **Figure 2C** in the main text, we evaluated the performance of all A-B-A-B two-drug cycles to use as a comparison
414 group for RL-fit and RL-genotype. In **Figure S9A**, we examine these combinations in greater depth. We also show the
415 landscape correlation between the two drugs in every combination. We show that anti-correlated landscapes tend to make
416 more effective combinations, likely due to collateral sensitivity. Highly correlated landscapes tend to make ineffective drug
417 combinations, likely due to collateral resistance.

418 Finally, we evaluated the effect of starting population genotype on the performance of each two-drug combination. We
419 found that the starting genotype of the population had no effect on the overall distribution of performance for these two-drug
420 combinations (**Fig S9**).

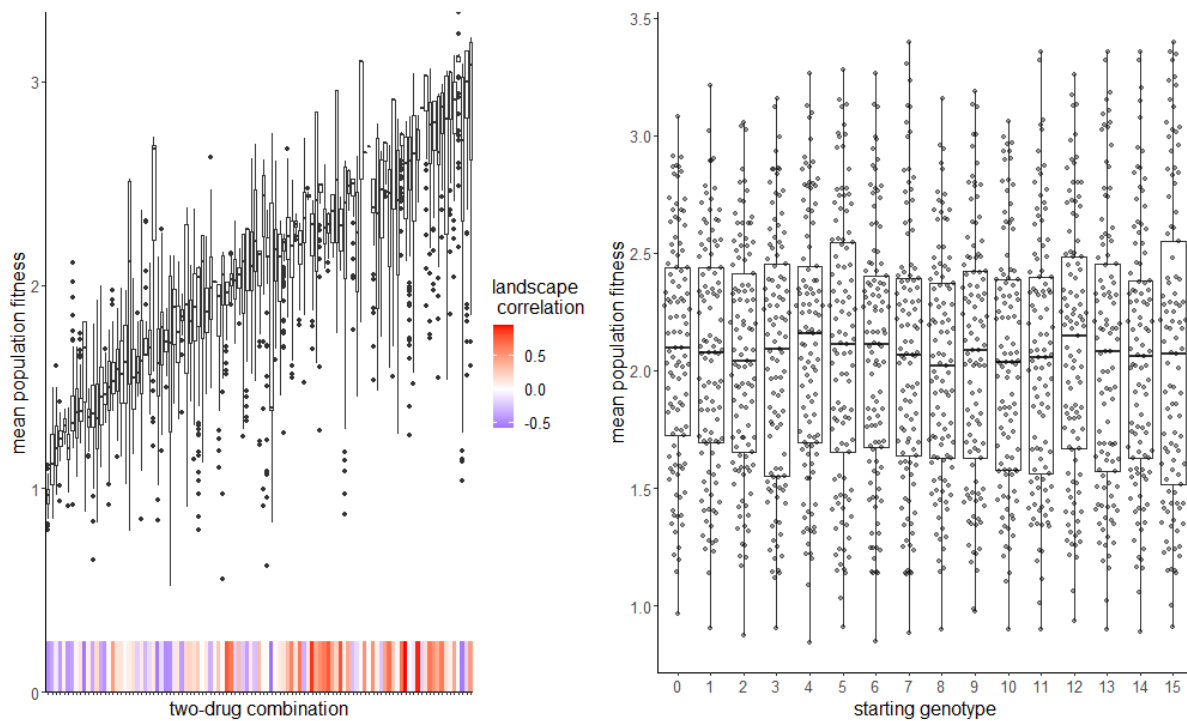


Figure S9. Two-drug cycling policies.