

Systems biology

Microbial interactions from a new perspective: reinforcement learning reveals new insights into microbiome evolution

Parsa Ghadermazi  ¹ and Siu Hung Joshua Chan  ^{1,*}

¹Department of Chemical and Biological Engineering, Colorado State University, Fort Collins, CO 80521, United States

*Corresponding author. Department of Chemical and Biological Engineering, Colorado State University, 254 Scott Bioengineering Building, 700 Meridian Ave, Fort Collins, CO 80523, United States. E-mail: joshua.chan@colostate.edu (S.H.J.C.)

Associate Editor: Pier Luigi Martelli

Abstract

Motivation: Microbes are essential part of all ecosystems, influencing material flow and shaping their surroundings. Metabolic modeling has been a useful tool and provided tremendous insights into microbial community metabolism. However, current methods based on flux balance analysis (FBA) usually fail to predict metabolic and regulatory strategies that lead to long-term survival and stability especially in heterogenous communities.

Results: Here, we introduce a novel reinforcement learning algorithm, Self-Playing Microbes in Dynamic FBA, which treats microbial metabolism as a decision-making process, allowing individual microbial agents to evolve by learning and adapting metabolic strategies for enhanced long-term fitness. This algorithm predicts what microbial flux regulation policies will stabilize in the dynamic ecosystem of interest in the presence of other microbes with minimal reliance on predefined strategies. Throughout this article, we present several scenarios wherein our algorithm outperforms existing methods in reproducing outcomes, and we explore the biological significance of these predictions.

Availability and implementation: The source code for this article is available at: <https://github.com/chan-csu/SPAM-DFBA>.

1 Introduction

Microbes are present in almost all known biotic environments and their metabolism affects the flow of materials in their ecosystems. Microbes form intricate networks of interacting cells from various taxonomic branches with distinct functional traits which makes predicting their behavior challenging. However, determining the role of microbial life in their ecosystems can be a key to solving numerous challenges that we face today. Imbalance in human gut microbiome is consistently linked with diseases such as inflammatory bowel disease ([Segata et al. 2012](#)). At a larger scale, microbial metabolism is a major player in geochemical cycles on earth ([Rousk and Bengtson 2014](#)).

Metagenomics studies provide detailed information about the membership and biochemical functions of microbiomes. However, predicting the phenotype of microbial communities from their genotype is by nature a complex problem and has been an ongoing effort for the past few decades ([Song et al. 2014](#), [Haruta and Yamamoto 2018](#), [Kumar et al. 2019](#), [Oriano et al. 2020](#)). Trophic interactions between microbes are an important factor that significantly contributes to the evolution of microbiome composition and function in various ecosystems ([Phelan et al. 2012](#), [Amundson et al. 2022](#)) and it further complicates the prediction of emergent properties of microbial communities.

Understanding and predicting the dynamics of microbial systems has remained largely unknown despite the enormous

growth in multiomics techniques and it requires a wholistic modeling approach ([Schmidt et al. 2021](#)). Mathematical models at different abstraction levels have been developed with the goal of making predictions that can explain the experimentally observed phenotypes ([Mahadevan et al. 2002](#), [Zomorrodi and Maranas 2012](#), [Khandelwal et al. 2013](#), [Song et al. 2014](#), [Zomorrodi et al. 2014](#), [Khodayari and Maranas 2016](#), [Bauer et al. 2017](#), [Chan et al. 2017](#), [Kumar et al. 2019](#), [Cai et al. 2020](#), [Dukovski et al. 2021](#)). GEnome-scale metabolic Models (GEMs) provide a detailed view of the biochemical networks of cells that are inferred from the genome of the organism of interest. GEMs generally contain up to thousands of biochemical reactions. Predicting the emergent properties of microbial communities by merely determining flux through such biochemical reactions is one of the main challenges in systems biology that yet remains to be addressed ([Song et al. 2014](#), [Oriano et al. 2020](#)). Flux balance analysis (FBA) is a bottom-up approach that provides a scalable method for simulating cellular metabolism in the absence of reaction kinetic parameters ([Orth et al. 2010](#)). FBA converts the system of differential equations resulting from mass balance across a cell to a linear programming (LP) problem by assuming steady state condition across the cells and defining a biologically relevant objective function ([Orth et al. 2010](#)). Despite the defined objective function and the constraints on the flux values, FBA solutions are rarely unique, and feasible solutions form a large space where distinct phenotypes can coexist. Dynamic flux balance analysis (DFBA)

applies FBA in each timepoint, and using the calculated extracellular fluxes, changes in extracellular metabolites with time are calculated. These rates of changes are used in turn to form a system of differential equations that describe the concentration profiles of different species in the system over time (Mahadevan *et al.* 2002, Uygun *et al.* 2006, Höffner *et al.* 2013, Gomez *et al.* 2014, Henson and Hanly 2014, Willemse *et al.* 2015, Zhao *et al.* 2017, Scott *et al.* 2018, Schroeder and Saha 2020, de Oliveira *et al.* 2023). However, the problem of lack of a unique solution in FBA propagates through time. Consequently, in the cases where the attempt is to model the dynamics of a heterogenous microbial community, one is faced with an extremely open solution space where different solutions can represent significantly different phenotypes while all phenotypes can satisfy FBA requirements. More importantly, DFBA relies on instantaneous biomass maximization assumption. Although in simple cases this assumption might result in realistic simulations (Mahadevan *et al.* 2002), in many other cases, it fails to predict the observed behavior of microbial systems (Cai *et al.* 2020) because depending on the environment, maximizing instantaneous growth rate can result in low fitness in future or even extinction. For example, cells that excrete extracellular amylase to breakdown starch are spending energy to do so and lower their instantaneous fitness in turn. However, secreting amylase is required for degrading starch to smaller molecules such as glucose for the cells' future use. Instantaneous biomass maximization will not allow any extracellular amylase secretion, unless previously set as a constraint on the model, while amylase secretion has been frequently observed in nature (Zhang *et al.* 2016, Song *et al.* 2019, Far *et al.* 2020).

Therefore, it is important to put the concept of Nash equilibria and evolutionary stability in this context of metabolic interactions (Cai *et al.* 2020, Schmidt *et al.* 2021). A remedy proposed by Zomorrodi and Segre (2017) was to determine Nash equilibria of the systems with several metabolic strategies of interest pre-defined. It beautifully captures the experimental results of some well-known microbial games of metabolic interactions. However, in general, it is virtually impossible to enumerate all possible metabolic strategies because of the high dimensional and continuous nature of the solution space defined by the mass balance constraints, directionality constraints, and nutrient availability. Therefore, an algorithm that can cover the entire possible solution space to a satisfactory extent when determining these stable interactions, and meanwhile does not rely on the instantaneous biomass maximization assumption, will greatly improve capability of predicting stable microbial interactions.

In this article, we aim to address these challenges by introducing a new modeling approach that integrates reinforcement learning (RL) into DFBA to model microbial metabolism in a microbiome as a decision-making process. From this perspective, microbial cells evolve by trying different metabolic strategies and learning how to improve their long-term fitness by tuning their behavior using a reinforcement learning algorithm. In this framework, each GEM is modeled as an agent capable of making decisions. The decisions in this context are flux regulations in the metabolic network and the agents make these decisions using the observable environment states. Assuming that “bad decisions” are filtered through the natural selection process, we use reinforcement learning algorithms to find the strategies that lead to the long-term optimal behavior of microbes in the system that they are interacting with. In other words, microbial models learn how to interact by trial and error

in their environment through self-play mechanism (Laterre *et al.* 2018), without the need to pre-define metabolic and regulatory strategies. Reinforcement learning has shown great promise in solving very complex problems in the past decade (Mnih *et al.* 2013, Silver *et al.* 2018, Brown *et al.* 2020) and have been used with success in different fields of science and engineering (Mousavi *et al.* 2017, Treloar *et al.* 2020, Jebellat *et al.* 2021, Kargar *et al.* 2022, Kiran *et al.* 2022, Lotfi *et al.* 2022). Although still relying on FBA, this approach is fundamentally different from biomass maximizing agents assumed commonly in traditional FBA and DFBA as the long-term consequences of actions are also considered in a dynamic context to find strategies that are also performing well in future rather than only an instance of time. In several cases, discussed shortly, greedily optimizing for biomass production will lead to early community extinction. Rationally, such strategies should be eliminated by the natural selection process. The strategies taken by RL agents after training can be useful to understand why certain types of behaviors are observed in real microbial systems.

2 Materials and methods

2.1 Reinforcement learning

In a reinforcement learning problem, one or more agents interact with their environment. An agent is an entity that is learning by interacting with the environment and is capable of making decisions in a given state. An environment can be defined as the collection of all entities that are surrounding the agent. Depending on the problem of interest, the agents can observe the entire or a part of their environment and based on a mapping called policy decide what actions to take in each state. Through this interaction with the environment the agents receive a reward and go to the next state according to the environment dynamics until it reaches the terminal state, that is when one episode is completed. The final goal of an agent is to maximize return during an episode. Return at time t is defined as the discounted sum of the collected reward after time t until the end of the episode (Sutton and Barto 2018):

$$G_t = \sum_{k=t+1}^t \gamma^{k-t-1} R_k. \quad (1)$$

Here γ is the discounting factor which determines the importance of future rewards and is the reward at time t . Policy is a function that describes the behavior of an agent in each state. Policy function is a mapping that outputs the probability of taking action a_t when the agent is in state s_t (Sutton and Barto 2018):

$$\pi : A \times S \rightarrow [0, 1]. \quad (2)$$

$$\pi(s, a) = P(a_t, s_t). \quad (3)$$

Another important definition is the **value function** under policy which is defined as Sutton and Barto (2018):

$$v_\pi(s) = E_\pi[G_t | S = s]. \quad (4)$$

Value for state s under policy π is the expected return of being in state s and following policy for the upcoming states

(Sutton and Barto 2018). In other word, the value function tells us how valuable it is to be in state s when following a specific policy function.

Depending on the type of problem at hand, there are families of RL algorithms that can be used (Arulkumaran *et al.* 2017, Sutton and Barto 2018). In the proposed framework (described in detail in Section 2.3), the states are observable extracellular metabolite concentrations which are continuous variables. On the other hand, actions are the reaction fluxes which are also continuous variables. This means that both the action and state space are continuous and multidimensional. For this type of problem, *policy gradient* family of algorithms is a good choice (Sutton and Barto 2018). Among them, proximal policy optimization (PPO) (Schulman *et al.* 2017) is a policy gradient algorithm which has been used with success in many Artificial Intelligence problems recently (Holubar and Wiering 2020, Han and Liang 2022, Wu *et al.* 2022, Yu *et al.* n.d.) and is appropriate for modeling microbial communities.

2.2 Proximal policy optimization

In policy gradient algorithms, the policy can be defined as a parameterized function, here a neural network, which maps a state to a real number, and the parameters of this function are tuned in a way so that the agent takes actions that lead to higher return because of the policy gradient theorem (Sutton and Barto 2018). One issue with using reinforcement learning algorithms is that the underlying mathematical operation in our algorithm is a linear programming problem. The suggested actions by the policy function can easily form a solution that is not in the feasible region of the LP problem and infeasible region can be in the proximity of areas where the return is maximum. This complicates the training process and we observed that many of the RL algorithms that we tested, failed on simple test cases. PPO addresses this issue by avoiding abrupt changes in the policy space. In the PPO algorithm, the agent tries to maximize the following surrogate objective function:

$$L^{\text{CLIP}}(\theta) = E_t \left[\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right], \quad (5)$$

where $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$ and $\hat{A}_t = R_t - v(s_t)$.

Setting this objective function (L^{CLIP}) causes the policy function to increase the probability of actions that lead to higher return during an episode. In Equation (5), $r_t(\theta)$ represents the ratio of the probability of taking action a_t at state s_t under the improved and the old policy. \hat{A}_t is called the advantage function which is the difference between reward collected at state s and the expected return of state s . The function does not allow $r_t(\theta)$ to go beyond $1 - \epsilon$ and $1 + \epsilon$. This objective definition has two implications on the policy function:

- 1) Probability of actions that result in more positive advantage will increase and the reverse will happen in the case of low advantage values.
- 2) The parameters of the policy function will only change if the change in the policy space is in a limited range.

Using this technique, PPO effectively stabilizes the training process. ϵ is a hyper parameter that depends on the problem of interest. In all of our simulations we used $\epsilon = 0.1$.

The value function and the policy function (called critic and actor networks, respectively) in all of our experiments are

neural networks with 10 linear layers followed by hyperbolic tangent activation function (Ramachandran *et al.* n.d.) with Adam optimizer as the network optimizer in both functions. In all simulation cases, we used 0.001 and 0.0001 as learning rates for critic and actor network, respectively.

2.3 Self-PIAying Microbes in Dynamic Flux Balance Analysis

Self-PIAying Microbes in Dynamic Flux Balance Analysis (SPAM-DFBA) is a dynamic algorithm. Every DFBA simulation happens over a defined period called episodes. When the final timepoint of the episode is reached, the agent go to the terminal state and a new episode is started. In the beginning of the algorithm, the agents are defined by assigning them a metabolic model. At each time point, or “state,” the agents observe a part of their environment, by sensing the concentrations of some of the extracellular metabolites. According to their policy function, a neural network in our implementation, these agents regulate a subset of fluxes through their metabolic network in each state by posing constraints on the reaction flux bounds (Fig. 1).

Decisions on uptake fluxes are limited by the user-provided kinetic expressions and parameters. If the flux through an exchange reaction is higher than what is allowable by the corresponding kinetic rates, the flux values are clipped to the highest value that is allowable by the kinetic rates for that compound. The constraints imposed by the policy network are coming from feeding the state to the actor function. The main difference between SPAM-DFBA and the standard DFBA is that during FBA for an agent at each timepoint, additional constraints returned by the agent’s policy function are imposed. Next, FBA is performed, and the rewards, actions, and states are recorded in an array until the end of the episode (Supplementary Fig. S1). Depending on the resources available, parallel episodes are simulated for the current policy at the same time. After one batch of episodes is finished, the critic network, i.e. the value function, gets updated based on the states and reward values. This step in reinforcement learning training process is called policy evaluation (Sutton and Barto 2018). Next, according to the PPO objective function

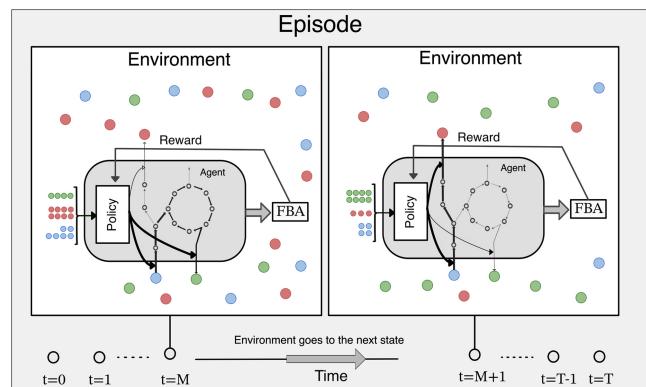


Figure 1. A schema of dynamic microbiome modeling viewed as a reinforcement learning problem. At each time step, the “agents” take an action according to their policy function without having a model of the environment dynamics. The actions in this sense are the constraints on the flux values. Based on the calculated flux with FBA, a reward is given to the agent. Changes in the environment is calculated using the exchange fluxes and the environment goes to the next state. An episode contains all the time points from $t=0$ to $t=7$. After an episode is finished, the agents improve their policies to maximize the reward signals.

(Equation 5) and using the information collected during the batch of episodes, the actor function is updated using gradient ascent to improve the policy (Supplementary Fig. S1). This step is called policy improvement. As a result of this iterative process, the agents become better at taking actions that maximize their own return by trial and error in the feasible solution space.

In our implementation, the reward function has two components: a negative reward as penalty for infeasible flux distributions and a positive reward for growth/biomass production. When the policy network generates flux values that are outside of the feasibility range, the agent will get a negative reward to learn to stay away from infeasibility. For the positive reward for the growth rate determined by FBA at each time point, we emphasize that it is different from the maximization of biomass production for the immediate time point in those future rewards also affect the decision made by the agent in any timepoint. As a result, an action with low immediate reward but high future reward might be favored over immediate biomass optimization strategy.

2.4 Implementation

SPAM-DFBA is implemented in Python. All the case studies are simulated with Python v3.10, and COBRApy v0.25.0 (Ebrahim et al. 2013). GLPK solver v0.4.7 was used to perform FBA in the toy communities and Gurobi optimizer (Gurobi Optimization 2023) with academic license was used to perform FBA on the GEMs using COBRApy interface. PyTorch v1.12.0 (Paszke et al. 2019) was used for building and training the neural networks. Same network structure and hyperparameters were used for all three cases to illustrate the robustness of this method with respect to the hyperparameters. Ray library was used to perform parallel computing (Moritz et al. 2017). Plotly v5.9.0 (Plotly Technologies Inc. 2015) was used to generate all the plots. Jupyter Notebooks are provided for reproducing all simulations. SPAM-DFBA is available as a package in PyPI with a detailed documentation website that facilitates using SPAM-DFBA for future studies. For more information, please refer to: <https://github.com/chan-csu/SPAM-DFBA>

3 Results

We created multiple toy microbial communities that exemplify the weakness of FBA or DFBA which are inspired by NECom (Cai et al. 2020) to demonstrate the advantage of SPAM-DFBA. The following subsections will provide a detailed description of these biologically relevant toy communities.

3.1 Amylase secretion without mass transfer considerations

This group of toy communities were designed to emulate a case that microbial cells are grown on a mixture of starch and glucose in a well-mixed chemostat system (Fig. 2). The cells are capable of secreting amylase to degrade the available starch. However, producing amylase is an energy-consuming step in the organism's metabolism and it requires ATP and precursors that would otherwise be used in biomass production. This poses a challenge on modeling the dynamics of such systems using DFBA because instant maximization of biomass would not allow any amylase secretion unless amylase production is set as a constraint on the underlying LP problem.

Additionally, the amount of amylase secretion is also an important consideration as too much amylase production impedes growth in the environment. Exoenzyme production in microorganisms has been a system of interest in studying microbial games because of the cheater-producer coexistence problem (Gore et al. 2009). Cheaters that do not secrete exoenzyme might evolve from an exoenzyme producer population because they benefit from the oligo-/mono-mers released by exoenzyme secreted by producers. In SPAM-DFBA, the agents take future awards into consideration in regulating their metabolic flux. As a result, they can learn to control amylase production in a manner that promotes their long-term survival by exploring various strategies in the environment and refining their policies depending on the rewards received.

3.1.1 Convergence and amylase secretion predicted in single-agent simulations

We tested what strategies the intelligent agents in SPAM-DFBA will learn in terms of exoenzyme production when only one homogenous population exists (single-agent simulations, Fig. 2A) vs. multiple phenotypes are allowed (two- and five-agent simulations, Fig. 2B and C). In all cases, the agents converge to a stable policy which cannot be further improved (Fig. 2D–F). There are significant growth and starch degradation in all cases (Fig. 2G–I), suggesting the capability of the algorithm to identify feasible and biologically relevant solutions.

3.1.2 Agents learn to minimize exploitability in multiagent environments

Simple DFBA cannot predict any amylase production by the metabolic network as it contradicts the biomass maximization assumption. However, SPAM-DFBA not only predicts amylase secretion on starch but also reveals an interesting pattern when comparing the overall community growth and starch degradation between the single-agent and multiagent cases. Starch utilization significantly decreases when the number of agents increases (Fig. 2G–I). This trend suggests that when more agents are present in the environment, they become more conservative in terms of secreting amylase, granted spatial homogeneity.

To look deeper into this observation, we examined the policy of the agents trained in different cases. The policy of the agent in Toy-Exoenzyme-Single-Agent is significantly different from the agents in Toy-Exoenzyme-Two-Agents (Fig. 3A–C). When the agent is trained in Exoenzyme-Single-Agent, the glucose resulted from breaking down starch can be utilized only by the agent itself. However, when other agents exist in the environment this is not the case. Other agents can learn to cheat because of the randomness in their behavior and take up the available glucose without paying the cost for building the amylase molecules. As a result, in the single-agent environment amylase secretion is negatively correlated with the glucose level, the optimal policy instructs the agent to secrete amylase when the glucose level is low, so that the low glucose concentration is compensated by breaking down the available starch. On the other hand, secreting amylase when the glucose level is low in environments with more than one agent is risky. In this case, at low glucose level cheating can have a more deteriorative effect on the amylase producer organism. For this reason, in environments with more than one agent amylase production increases with glucose level.

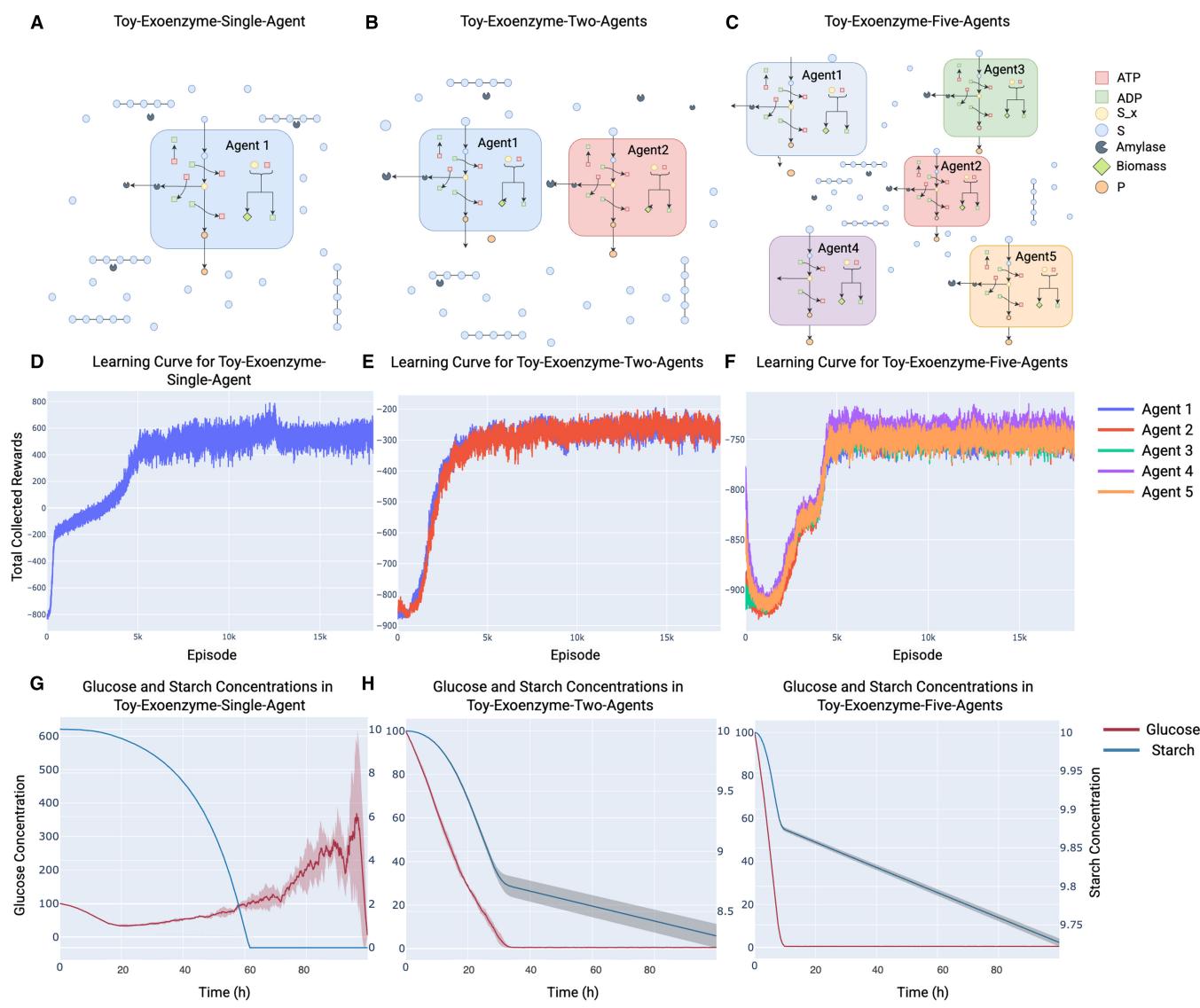


Figure 2. Training agents in a well-mixed chemostat system with starch and glucose initially present in the system. (A–C) Schemas for the toy communities: Toy-Exoenzyme-Single-Agent, Toy-Exoenzyme-Two-Agents, and Toy-Exoenzyme-Five-Agents, respectively. (D–F) Learning curve of the agents in Toy-Exoenzyme-Single-Agent, Toy-Exoenzyme-Two-Agents, and Toy-Exoenzyme-Five-Agents after training on 5000 batches of four episodes, respectively. Learning curves show the total collected rewards during an episode change over the course of the training process. Due to the energy required for maintenance, in multiagent environments most agents receive negative rewards with the absence of sugar. (G–I) Starch and Glucose concentration over time in Toy-Exoenzyme-Single-Agent, Toy-Exoenzyme-Two-Agents, and Toy-Exoenzyme-Five-Agents, respectively. The solid lines represent the mean value across all episodes in a batch and the shades represent 1 SD across all episodes in a batch. Note that each actor acts randomly around the mean of the actor network output with standard deviation of 0.1. The left axis in each plot shows glucose concentrations and the right axes show starch concentrations.

What follows from this observation is that in a multiagent environment, the agents that are trained in multiagent environments should perform better than when trained in a single-agent environment as they have experienced different aspects of coexisting with other agents, such as cheating by other agents. To test this statement, we created a new two-agent environment, Toy-Exoenzyme-Single-Two-Comb. One agent is selected from Toy-Exoenzyme-Single-Agent and the other from Toy-Exoenzyme-Two-Agents. *Supplementary Fig. S2* shows that Agent 1 that is trained in Toy-Exoenzyme-Single-Agent achieves significantly lower return compared to Agent 2 that is trained in Toy-Exoenzyme-Two-Agents. Furthermore, the level of starch utilization is higher than Toy-Exoenzyme-Two-Agents but lower than Toy-Exoenzyme-Single-Agent which points to the fact that high amylase

production by Agent 1 is exploited by Agent 2. This is an example of an emergent property of microbiomes which is made possible by considering the long-term effect of actions and explains the observed deterioration of performance in systems where large molecules are broken down by microbial cells.

3.2 Amylase secretion with mass transfer considerations

An important question is that how mass transfer rate in an environment can change this cheating behavior. Higher mass transfer increases the possibility of components moving away from the producers and being used by the cheaters. To test if our algorithm predicts lower mass transfer rate favors amylase producers, we simulated two five-agent environments: a

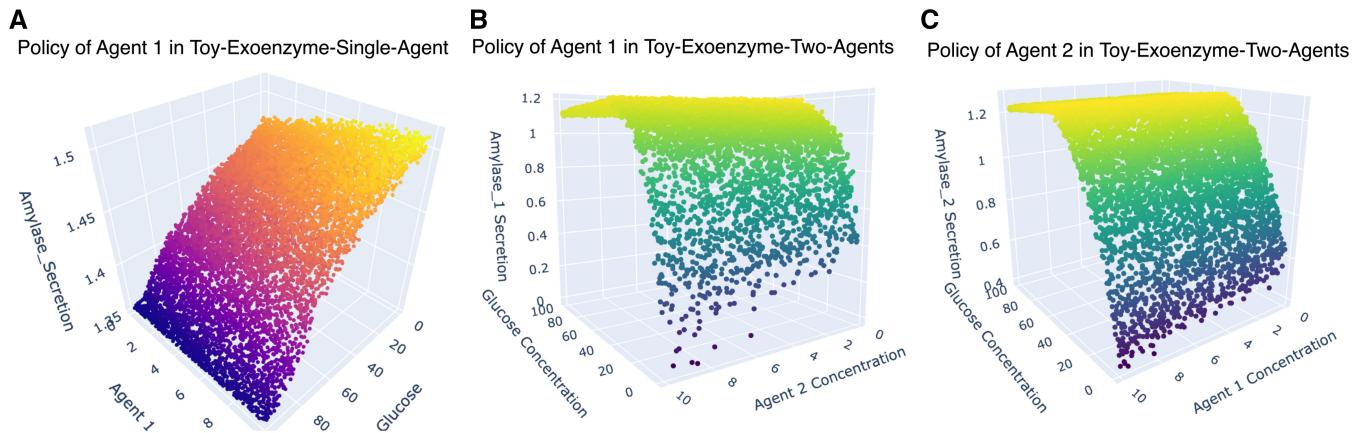


Figure 3. The policy profiles of agents in Toy-Exoenzyme-Single-Agent and Toy-Exoenzyme-Two-Agents environments after training on 5000 batches of four episodes. These plots are created by randomly generating 10 000 points in the policy space of the trained agents. (A) Policy of Agent 1 in Toy-Exoenzyme-Single-Agent with respect to glucose and Agent 1's concentration. (B) Policy of Agent 1 in Toy-Exoenzyme-Two-Agents with respect to glucose and Agent 2's concentration. (C) Policy of Agent 2 in Toy-Exoenzyme-Two-Agents with respect to glucose and Agent 1's concentration.

high versus low mass transfer rate. Note that SPAM-DFBA does not consider spatial variations. We simply added pseudo-reactions to distinguish between private and public goods and limit the generation rate of the public goods. More detailed information about this implementation can be found in the Jupyter Notebook titled “Case_Study_1_Starch_Amylase” in the documentation website and the GitHub repository for this project. The result of this experiment agreed with our hypothesis and the agents in the environment with lower mass transfer rate achieved higher return and higher starch utilization, i.e. lower final starch concentration (Supplementary Fig. S3).

3.3 Toy-NECOM-Auxotroph

Metabolite exchange between two auxotroph strains is another type of interaction that has been observed frequently in nature (Mee *et al.* 2014). Modeling a community containing such strains with DFBA is problematic. In many cases, such as amino acid exchange, biomass maximization assumption in DFBA does not allow secretion of an amino acid. However, secretion of such compounds can be beneficial in the long run as the auxotrophs rely on the metabolic product of the other strains to survive. This problem becomes even more interesting since usually different strains of same species could compete for the same resources. To see if SPAM-DFBA can predict metabolite exchange between auxotrophs, Toy-NECOM-Auxotrophs environment was created. The schematic description of this environment is provided in Fig. 4A. Although they cannot survive on their own, the agents can grow by exchanging A and B. Each agent can sense its own concentration, the concentration of the other agent, concentration of S, concentration of extracellular A, and the concentration of extracellular B. Figure 4 shows the result of training the agents in this environment.

After training in this environment, the agents learn exchanging precursors A and B with each other. This prediction is qualitatively similar to the prediction by NECom (Cai *et al.* 2020). To contrast with this result, we also simulated a community where two agents can synthesize both A and B but with different efficiencies (Supplementary Fig. S4). Mutual crossfeeding is beneficial but not obligatory for survival. Both agents learn to selfishly take up any A and B and not to cooperate. This is consistent with the previous results and was

demonstrated previously to be not captured by maximizing community biomass (Khandelwal *et al.* 2013).

3.3.1 Simulating adaptation in new environments

If metabolite exchange in this environment is dictated merely by the fact that cells rely on each other for survival, what happens if we supplement the Toy-NECOM-auxotroph environment agents with A and B externally after training? The underlying hypothesis is that this richer environment should discourage the cooperation between the two complementary auxotrophs. To answer this question, we created a new environment, Toy-NECOM-Auxotrophs-Shift. We used the two auxotroph agents from Toy-NECOM-Auxotrophs and simulated a scenario where A and B initially is supplemented in the environment. This case is designed to predict how changes in environment can shape the trophic behavior of microbial communities.

Figure 5 shows that the auxotrophic agents shift their policy from metabolic exchange to selfishly taking up A and B (negative median flux) when A and B is supplemented externally. This prediction is consistent with previous experimental results that supplementation of the needed metabolites discourages the cross-feeding interactions between auxotrophic mutants (Hoek *et al.* 2016). This also shows an intriguing capability of SPAM-DFBA to predict how a certain microbial population adapted to an environment might evolve in a new environment.

3.4 *Escherichia coli* auxotrophs

So far, we only discussed small toy examples. IJO1366-Tyr-Phe-Auxotrophs environment was created to prove the scalability of our approach to a community of genome-scale models, IJO1366 for *Escherichia coli* K-12 MG1655 (Orth *et al.* 2011). In this environment, two *E.coli* auxotrophs were made. One mutant could not synthesize tyrosine and the other could not synthesize phenylalanine. Neither of the mutants could grow alone on M9 minimal medium. However, after training the agents converged to a policy that they would exchange the amino acid that they can produce and grow by using the amino acid secreted by the other agent. In this environment, phenylalanine mutant grew significantly more than the tyrosine mutant (Fig. 6). This observation is experimentally validated for the same system by Mee *et al.* (2014) where the exchange of amino acids between the strains and

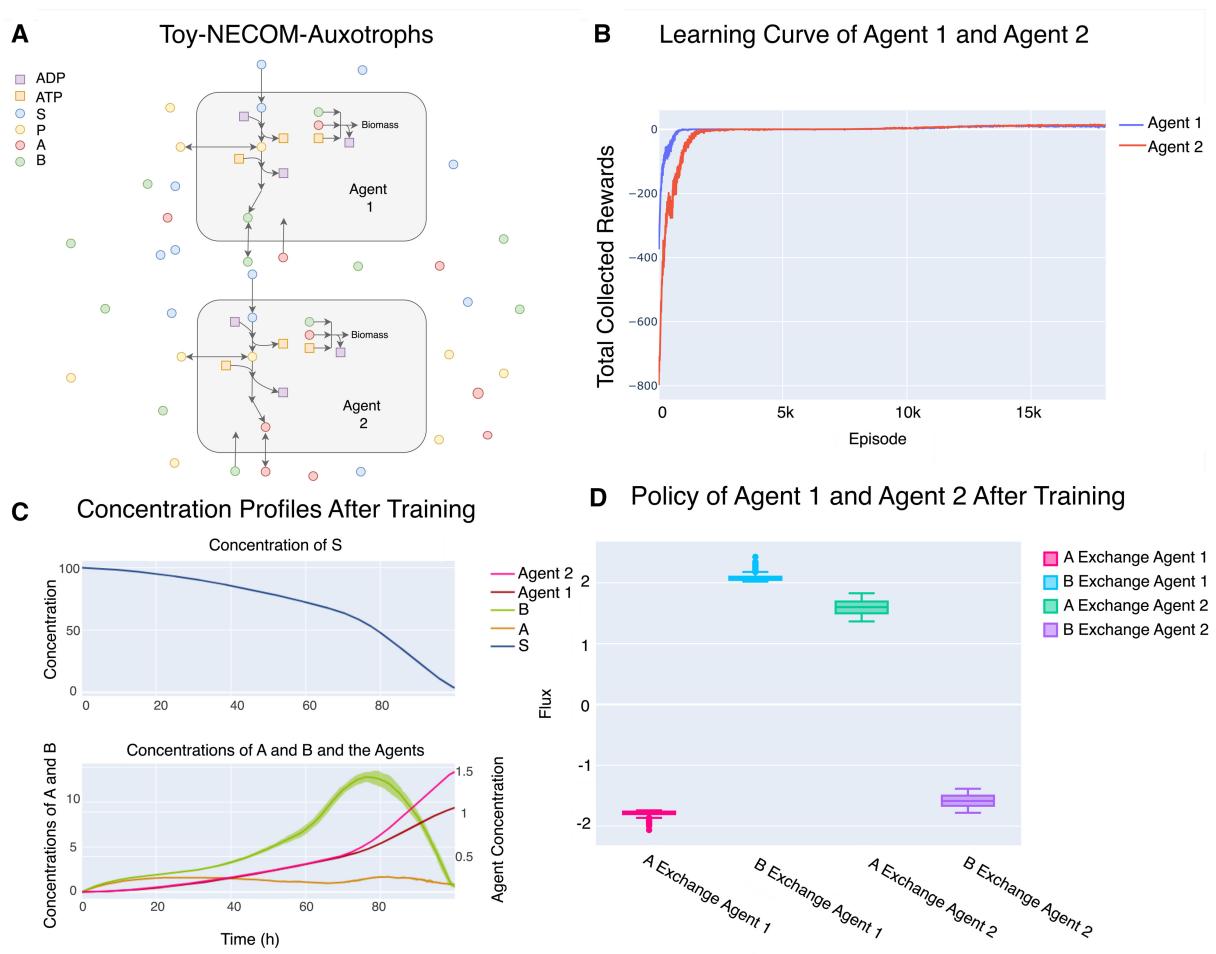


Figure 4. Training two agents in Toy-NECOM-Auxotrophs environment. (A) A schematic view of the environment. In this case Agent 1 and Agent 2 have similar metabolic network with only one difference. Agent 1 cannot produce the biomass precursor A and Agent 2 cannot produce the precursor B. (B) Learning curve of the two agents during 5000 batches of training. (C) Concentration profile of the species in this environment over time in a batch of four episodes. (D) Policies learned by the agents. Negative sign for fluxes means uptake and positive sign means secretion. Agent 1 learns to uptake whatever (A) that exists in the environment, while secreting (B) for agent 2. Agent 2 has learned the opposite strategy which agrees with their mutation.

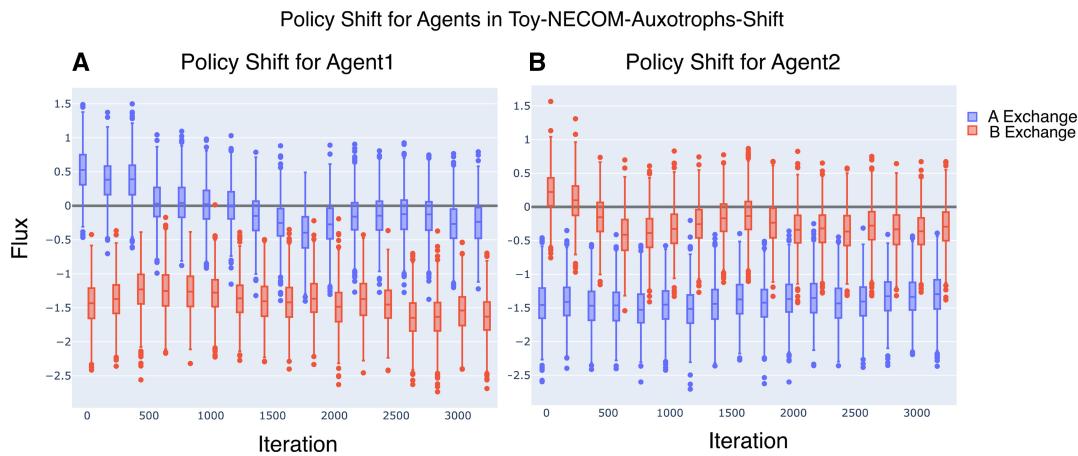


Figure 5. Policy shift after 3000 iteration of training agents of Toy-NECOM-Auxotrophs environment in Toy-NECOM-Auxotrophs-Shift. Both agents shift their policy from cross-feeding to taking up both (A) and (B) as much as possible. Here, the box plots show the range of actions across randomly generated states to represent the policy function.

the prevalence of the phenylalanine mutant has been verified (Agent 2 in our simulation). No phenotypic data from this or other experiments was used during the training process,

which shows the promise of the assumption that evolution favors traits at individual level that leads to high long-term fitness of the cells.

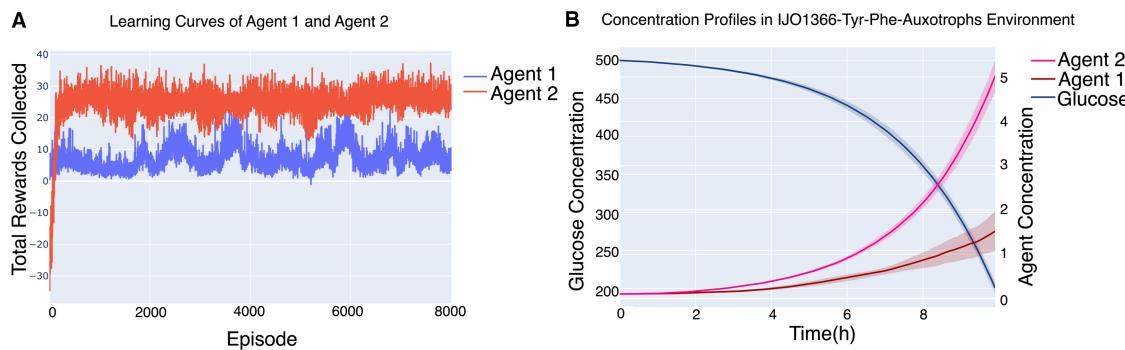


Figure 6. Training *E.coli* mutants in IJO1366-Tyr-Phe-Auxotrophs environment. Agent 1 is Tyrosine mutant and Agent 2 is Phenylalanine mutant. (A) Learning curve of Agent 1 and Agent 2. Agent 2 receives higher return than Agent 1 although it starts from worse policy. (B) Concentration profiles for glucose and the agents. It is also obvious from B that Agent 2, phenylalanine, achieves superior growth compared to Agent 1, tyrosine agent.

4 Discussion and conclusion

In this article, we presented a novel algorithm that provides insights into microbial interactions by allowing the microbial agents to freely explore flux regulation strategies and select metabolism regulation strategies that lead to their higher long-term fitness of the agents. This way we can explain the observed phenotypes in multiple communities that current algorithms fail to explain. Defining this problem in a FBA framework forces the strategies to be inside a space where mass balance and flux constraints are still satisfied. Another advantage of using FBA is that the underlying LP problems can be solved efficiently. We examined this algorithm on multiple test scenarios that emulate biologically relevant scenarios.

In scenarios where agents should coexist with other agents, the agents learned aspects of interacting with others while still tried to maximize their own return. The outcome of starch-amylase system has an interesting interpretation. Cheating in microbial communities can significantly affect the amount which large molecules such as starch are degraded in hydrolysis (Velicer *et al.* 2000, Rainey and Rainey 2003, Greig and Travisano 2004, Yurtsev *et al.* 2013, Harrington and Sanchez 2014, Popat *et al.* 2015, Szilágyi *et al.* 2017, Abisado *et al.* 2018, Heyer *et al.* 2019, Morales *et al.* 2021, Han and Liang 2022). Taking spatial heterogeneity into consideration revealed that in communities with higher mass transfer limits, the agents secrete more amylase and starch utilization becomes higher (Supplementary Fig. S3). The reason behind this observation is that low mass transfer implies that the glucose that is produced by an agent will stay away from the other agent that could possibly cheat, and in turn, the agents will see more positive signal by secreting amylase.

With this algorithm, we were able to explore other types of microbial interactions in a dynamic context. One problem that we were interested in was that whether we can explain metabolite exchange between auxotrophic strains through this framework (Mee *et al.* 2014, Zengler and Zaramela 2018). We hypothesized that without any predetermined exchange strategies or community level objectives, the optimal agents can find metabolite exchange with other agents strategy to maximize their own long-term fitness. Optimal agents in Toy-NECOM-Auxotrophs learned that exchanging A and B will increase their long-term fitness.

To see if this algorithm can be used for genome-scale model in real environments, we created an environment of two *E.coli* auxotrophs, tyrosine and phenylalanine, inspired by

the experiments in Mee *et al.* (2014). Although we did not use any sort of experimentally observed phenotypic data, the agents learned to exchange the amino acid that they can produce, and the other agent cannot. Interestingly, our simulations indicate that the phenylalanine mutant achieves superior growth compared to the tyrosine mutant, Fig. 6, which follows the same trend as is experimentally observed and reported in Mee *et al.* (2014). Being able to predict such emergent behaviors of microbiomes by purely relying on metabolic capability of the cells and ecological first principles is what distinguishes SPAM-DFBA from the other existing algorithms.

An interesting study (Hoek *et al.* 2016) reported the behavior change of auxotrophs when inserted in an environment that supplies all the components that they need for growth. In this scenario, they shift their exchange strategy to uptake all the compounds from the environment and stop secreting the metabolites further. Our simulations showed similar shift for auxotrophic agents which reflects that the agents adapt their strategies according to the changes in the environment, Fig. 5, and shows assuming that cells are *maximizing their own long-term fitness* can reproduce several real scenarios is missed by simple DFBA.

Previous cases revealed that agents that depend on each other for survival will evolve to exchange metabolites with each other and when this strict dependence does not exist anymore selfish behaviors emerge. Toy-NECOM-Facultative-Exchange provides more evidence for this trend. In this case, if a community level objective such as, total community biomass maximization, is defined then there will be A and B exchange between Agent 1 and Agent 2 (Khandelwal *et al.* 2013, Cai *et al.* 2020). This is the result predicted by the direct extension of FBA where a microbial community is optimized as one compartmentalized model. However, this is not what SPAM-DFBA predicts. In this case, A and B exchange strategy is exploitable by the agents. Since the agents do not rely on each other for survival, any exchange of A and B is exploitable by either of the agents in the case of resource limitation. Consequently, the agents finally adhere to taking up any A and B, limited by the kinetic rules provided for the model, which exist in the environment which is shown in Supplementary Fig. S4. This is consistent with the previous NECom prediction and game-theoretical analysis (Cai *et al.* 2020), and exactly matches simple DFBA prediction.

SPAM-DFBA is well suited for answering important questions in the field of microbiology by predicting the emergent behavior of microbiomes using metabolic capability of the

cells in contrast with commonly used ecological models such as Generalized Lotka-Volterra (Bomze 1983, Hofbauer and Sigmund 1998, Venturelli *et al.* 2018), which do not base their predictions on the metabolic network of the microbes. SPAM-DFBA is a dynamic framework, and the environment changes such as resources limitations can be simulated, while methods based on FBA, and not DFBA, cannot make such considerations which is critical and can significantly shape microbial interactions (Hoek *et al.* 2016).

Another advantage of this approach is that optimization is done at individual model level instead of community level objectives. This means that unrealistic interactions discussed in detail in (Cai *et al.* 2020) are avoided. If a particular random action is advantageous to the long-term fitness of an agent, this behavior gets reinforced in the policy of the agent using the PPO algorithm. We believe that this has a lot of similarities to the process of natural selection. We would like to emphasize that our method does not imply that the microbial cells are intelligently seeking the optimal behavior in their environment. However, it is the resemblance of this algorithm to the process of natural selection that results in more realistic predictions for a given environment.

There are multiple future directions to further improve SPAM-DFBA. While the current implementation can scale to multiple GEMs in a manageable amount of time, more efficient implementations will help simulate more complex microbiomes with hundreds of taxa. Another particularly interesting venue is to completely remove the need for LP solvers by letting the agents sample the feasible action space. This not only could improve the speed dramatically, but it also relaxes any assumption imposed on the metabolism of the agents under optimization formulation. Improvements in sample efficiency of RL algorithms can also improve the efficiency of this algorithm in future.

Hyperparameters such as clipping threshold or learning rates also affect the efficiency and stability of the learning process for the agents. Although we used same hyperparameters for all the case studies, optimal combination of the hyperparameters can be explored either by exhaustive search or using appropriate optimization techniques (Boroujeni and Pashaei 2021, Kiran and Ozyildirim 2022). As an example, we have examined the effect of important hyperparameters on the learning process for “Toy-Exoenzyme-Single-Agent” (Supplementary Fig. S5).

In this article, we just showed the potentials of formulating DFBA as a RL problem and discussed how this approach can predict microbial interactions in simple communities. Applying this approach to more complex ecosystems and validation with experimentally observed phenotypes is a natural next step for future studies.

Supplementary data

Supplementary data are available at *Bioinformatics* online.

Conflict of interest

None declared.

Funding

This work was supported by the U.S. Army Research Office and U.S. Army Research Laboratory and was accomplished

under Cooperative Agreement Number W911NF-07-2-0055. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Office, Army Research Laboratory, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation hereon.

Data availability

No new data were generated or analysed in support of this research.

References

- Abisado Rhea G, Benomar S, Klaus JR *et al.* Bacterial quorum sensing and microbial community interactions. *mBio* 2018;9:e02331–17. <https://doi.org/10.1128/mBio.02331-17>.
- Amundson KK, Borton MA, Daly RA *et al.* Microbial colonization and persistence in deep fractured shales is guided by metabolic exchanges and viral predation. *Microbiome* 2022;10:5. <https://doi.org/10.1186/s40168-021-01194-8>.
- Arulkumaran K, Deisenroth MP, Brundage M *et al.* Deep reinforcement learning: a brief survey. *IEEE Signal Process Mag* 2017;34:26–38. <https://doi.org/10.1109/MSP.2017.2743240>.
- Bauer E, Zimmermann J, Baldini F *et al.* BacArena: individual-based metabolic modeling of heterogeneous microbes in complex communities. *PLoS Comput Biol* 2017;13:e1005544.
- Bomze IM. Lotka-Volterra equation and replicator dynamics: A two-dimensional classification. *Biological Cybernetics* 1983;48:201–11.
- Boroujeni SPH, Pashaei E. A novel hybrid gene selection based on random forest approach and binary dragonfly algorithm'. In: 2021 18th International Conference on Electrical Engineering, Computing Science and Automatic Control (CCE), Mexico City, Mexico, 2021, 1–8.
- Brown N, Bakhtin A, Lerer A. *et al.* Combining Deep Reinforcement Learning and Search for Imperfect-Information Games. In: *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS'20*. Red Hook, NY, USA: Curran Associates Inc. 2020.
- Cai J, Tan T, Chan SHJ. Predicting nash equilibria for microbial metabolic interactions. *Bioinformatics* 2020;36:5649–55. <https://doi.org/10.1093/bioinformatics/btaa1014>.
- Chan SHJ, Simons MN, Maranas CD. SteadyCom: predicting microbial abundances while ensuring community stability. *PLoS Comput Biol* 2017;13:e1005539. <https://doi.org/10.1371/journal.pcbi.1005539>.
- de Oliveira RD, Le Roux GAC, Mahadevan R. Nonlinear programming reformulation of dynamic flux balance analysis models. *Comput Chem Eng* 2023;170:108101. <https://doi.org/10.1016/j.compchemeng.2022.108101>.
- Dukovski I, Bajić D, Chacón JM *et al.* A metabolic modeling platform for the computation of microbial ecosystems in time and space (COMETS). *Nat Protoc* 2021;16:5030–82. <https://doi.org/10.1038/s41596-021-00593-3>.
- Ebrahim A, Lerman JA, Palsson BO *et al.* COBRApY: COnstraints-based reconstruction and analysis for python. *BMC Syst Biol* 2013; 7:74. <https://doi.org/10.1186/1752-0509-7-74>.
- Far BE, Ahmadi Y, Yari Khosrourshahi A *et al.* Microbial alpha-amylase production: progress, challenges and perspectives. *Adv Pharm Bull* 2020;10:350–8.
- Gomez JA, Höffner K, Barton PI. DFBAlab: a fast and reliable MATLAB code for dynamic flux balance analysis. *BMC Bioinformatics* 2014;15:409. <https://doi.org/10.1186/s12859-014-0409-8>.
- Gore J, Youk H, Van Oudenaarden A. Snowdrift game dynamics and facultative cheating in yeast. *Nature* 2009;459:253–6. <https://doi.org/10.1038/nature07921>.

- Greig D, Travisano M. The Prisoner's Dilemma and polymorphism in yeast SUC genes. *Proceedings of the Royal Society of London. Series B: Biological Sciences* 2004;271:S25–S26.
- Gurobi Optimization, LLC. Beaverton, Oregon. *Gurobi Optimizer Reference Manual*. 2023.
- Han SY, Liang T. Reinforcement-learning-based vibration control for a vehicle semi-active suspension system via the PPO approach. *Appl Sci (Switzerland)* 2022;12:3078. <https://doi.org/10.3390/app12063078>.
- Harrington KI, Sanchez A. Eco-evolutionary dynamics of complex social strategies in microbial communities. *Commun Integr Biol* 2014;7:e28230. <https://doi.org/10.4161/cib.28230>.
- Haruta S, Yamamoto K. Model microbial consortia as tools for understanding complex microbial communities. *Curr Genomics* 2018;19:723–33. <https://doi.org/10.2174/138920291966180911131206>.
- Henson MA, Hanly TJ. Dynamic flux balance analysis for synthetic microbial communities. *IET Syst Biol* 2014;8:214–29. <https://doi.org/10.1049/iet-syb.2013.0021>.
- Heyer R, Schallert K, Siewert C et al. Metaproteome analysis reveals that syntrophy, competition, and Phage-Host interaction shape microbial communities in biogas plants. *Microbiome* 2019;7:69. <https://doi.org/10.1186/s40168-019-0673-y>.
- Hock TA, Axelrod K, Biancalani T et al. Resource availability modulates the cooperative and competitive nature of a microbial cross-feeding mutualism. *PLoS Biol* 2016;14:e1002540. <https://doi.org/10.1371/journal.pbio.1002540>.
- Hofbauer J, Sigmund K. *Evolutionary Games and Population Dynamics*. Cambridge: Cambridge University Press, 1998.
- Höffner K, Harwood SM, Barton PI. A reliable simulator for dynamic flux balance analysis. *Biotechnol Bioeng* 2013;110:792–802. <https://doi.org/10.1002/bit.24748>.
- Holubar MS, Wiering MA. Continuous-action reinforcement learning for playing racing games: Comparing SPG to PPO. *ArXiv*, abs/2001.05270 2020.
- Jebellat I, Pishkenari HN, Jebellat E. Training microrobots via reinforcement learning and a novel coding method. In: 2021 9th RSI International Conference on Robotics and Mechatronics (ICRoM), Tehran, Islamic Republic of Iran, 2021, 105–111.
- Kargar M, Sardarmehni T, Song X. Optimal powertrain energy management for autonomous hybrid electric vehicles with flexible driveline power demand using approximate dynamic programming. *IEEE Trans Veh Technol* 2022;71:12564–75. <https://doi.org/10.1109/TVT.2022.3199681>.
- Khandelwal RA, Olivier BG, Röling WFM et al. Community flux balance analysis for microbial consortia at balanced growth. *PLoS One* 2013;8:e64567. <https://doi.org/10.1371/journal.pone.0064567>.
- Khodayari A, Maranas CD. A genome-scale *Escherichia coli* kinetic metabolic model k-Ecoli457 satisfying flux data for multiple mutant strains. *Nat Commun* 2016;7:13806. <https://doi.org/10.1038/ncomms13806>.
- Kiran BR, Sobh I, Talpaert V et al. Deep reinforcement learning for autonomous driving: a survey. *IEEE Trans Intell Transport Syst* 2022;23:4909–26. <https://doi.org/10.1109/TITS.2021.3054625>.
- Kiran M, Ozyildirim M. Hyperparameter tuning for deep reinforcement learning applications. *ArXiv*, abs/2201.11182. 2022.
- Kumar M, Ji B, Zengler K et al. Modelling approaches for studying the microbiome. *Nat Microbiol* 2019;4:1253–67.
- Laterre A, Fu Y, Jabri MK et al. Ranked reward: enabling self-play reinforcement learning for combinatorial optimization. *ArXiv*, abs/1807.01672. 2018.
- Lotfi F, Semiaci O, Afghah F et al. Evolutionary deep reinforcement learning for dynamic slice management in O-RAN. In: 2022 IEEE Globecom Workshops (GC Wkshps), Rio de Janeiro, Brazil, 2022, 227–232.
- Mahadevan R, Edwards JS, Doyle FJ. Dynamic flux balance analysis of diauxic growth in *Escherichia coli*. *Biophys J* 2002;83:1331–40. [https://doi.org/10.1016/S0006-3495\(02\)73903-9](https://doi.org/10.1016/S0006-3495(02)73903-9).
- Mee MT, Collins JJ, Church GM et al. Syntrophic exchange in synthetic microbial communities. *Proc Natl Acad Sci USA* 2014;111:E2149–E2156. <https://doi.org/10.1073/pnas.1405641111>.
- Mnih V, Kavukcuoglu K, Silver D et al. Playing atari with deep reinforcement learning. *ArXiv*, abs/1312.5602. 2013.
- Morales Neydis M, Patel M, Stewart CJ et al. Optogenetic tools for control of public goods in *saccharomyces cerevisiae*. *mSphere* 2021;6. <https://doi.org/10.1128/msphere.00581-21>.
- Moritz P, Nishihara R, Wang S et al. Ray: a distributed framework for emerging ai applications. *ArXiv*, abs/1712.05889 2017.
- Mousavi SS, Schukat M, Howley E. Traffic light control using deep policy-gradient and value-function based reinforcement learning. *ArXiv*, abs/1704.08883 2017.
- Oriano M, Zorzetto L, Guagliano G et al. The open challenge of *in vitro* modeling complex and multi-microbial communities in three-dimensional niches. *Front Bioeng Biotechnol* 2020;8:539319.
- Orth JD, Conrad TM, Na J et al. A comprehensive genome-scale reconstruction of *Escherichia coli* metabolism. *Mol Syst Biol* 2011;7:535. <https://doi.org/10.1038/msb.2011.65>.
- Orth JD, Thiele I, Palsson BO. What is flux balance analysis? *Nat Biotechnol* 2010;28:245–8. <https://doi.org/10.1038/nbt.1614>.
- Paszke A, Gross S, Massa F et al. (2019) PyTorch: An imperative style, high-performance deep learning library. In: Wallach H., et al. (eds), *Advances in Neural Information Processing Systems*. Red Hook, NY, USA: Curran Associates, Inc. 2019.
- Phelan VV, Liu W-T, Pogliano K et al. Microbial metabolic exchange—the chemotype-to-Phenotype link. *Nat Chem Biol* 2012;8:26–35. <https://doi.org/10.1038/nchembio.739>.
- Plotly Technologies Inc. (2015). Collaborative data science. Montreal, QC: Plotly Technologies Inc. Retrieved from <https://plot.ly>.
- Popat R, Pollitt EJG, Harrison F et al. Conflict of interest and signal interference lead to the breakdown of honest signaling. *Evolution* 2015;69:2371–83. <https://doi.org/10.1111/evo.12751>.
- Rainey PB, Rainey K. Evolution of cooperation and conflict in experimental bacterial populations. *Nature* 2003;425:72–4. <https://doi.org/10.1038/nature01906>.
- Ramachandran P, Zoph B, Le QV. Swish: a self-gated activation function. *arXiv: Neural and Evolutionary Computing* 2017.
- Rousk J, Bengtson P. Microbial regulation of global biogeochemical cycles. *Front Microbiol* 2014;5:103.
- Schmidt Caleb M, Ghadermazi P, Chan SHJ. Predicting microbiome metabolism and interactions through integrating multidisciplinary principles. *mSystems* 2021;6. <https://doi.org/10.1128/msystems.00768-21>.
- Schroeder WL, Saha R. Introducing an optimization- and explicit Runge-Kutta-based approach to perform dynamic flux balance analysis. *Sci Rep* 2020;10:9241. <https://doi.org/10.1038/s41598-020-65457-4>.
- Schulman J, Wolski F, Dhariwal P et al. Proximal policy optimization algorithms. *ArXiv*, abs/1707.06347 2017.
- Scott F, Wilson P, Conejeros R et al. Simulation and optimization of dynamic flux balance analysis models using an interior point method reformulation. *Comput Chem Eng* 2018;119:152–70. <https://doi.org/10.1016/j.compchemeng.2018.08.041>.
- Segata N, Waldron L, Ballarini A et al. Metagenomic microbial community profiling using unique clade-specific marker genes. *Nat Methods* 2012;9:811–4. <https://doi.org/10.1038/nmeth.2066>.
- Silver D, Hubert T, Schrittwieser J et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science* 2018;362:1140–4. <https://doi.org/10.1126/science.aar6404>.
- Song HS, Cannon WR, Beliaev AS et al. Mathematical modeling of microbial community dynamics: a methodological review. *Processes* 2014;2:711–52.
- Song J, Tang H, Liang H et al. Effect of bioaugmentation on biochemical characterisation and microbial communities in Daqu using *Bacillus*, *Saccharomyces* and *Absidia*. *Int J Food Sci Tech* 2019;54:2639–51. <https://doi.org/10.1111/ijfs.14176>.
- Sutton RS, Barto AG. *Reinforcement Learning: An Introduction* Second. The MIT Press 2018.
- Szilágyi A, Boza G, Scheuring I. Analysis of stability to cheaters in models of antibiotic degrading microbial communities. *J Theor Biol* 2017;423:53–62. <https://doi.org/10.1016/j.jtbi.2017.04.025>.
- Treloar NJ, Fedorec AJH, Ingalls B et al. Deep reinforcement learning for the control of microbial co-cultures in bioreactors. *PLoS Comput Biol* 2020;16:e1007783.

- Uygun K, Matthew HWT, Huang Y. DFBA-LQR: an optimal control approach to flux balance analysis. *Ind Eng Chem Res* 2006;45: 8554–64. <https://doi.org/10.1021/ie060218f>.
- Velicer GJ, Kroos L, Lenski RE. Developmental cheating in the social bacterium *Myxococcus Xanthus*. *Nature* 2000;404:598–601. <https://doi.org/10.1038/35007066>.
- Venturelli OS, Carr AV, Fisher G et al. Deciphering microbial interactions in synthetic human gut microbiome communities. *Mol Syst Biol* 2018;14:e8157. <https://doi.org/10.15252/msb.20178157>.
- Willemse AM, Hendrickx DM, Hoefsloot HCJ et al. MetDFBA: incorporating time-resolved metabolomics measurements into dynamic flux balance analysis. *Mol Biosyst* 2015;11:137–45. <https://doi.org/10.1039/c4mb00510d>.
- Wu G, Fang W, Wang J et al. Dyna-PPO reinforcement learning with gaussian process for the continuous action decision-making in autonomous driving. *Appl Intell* 2022;53:16893–907. <https://doi.org/10.1007/s10489-022-04354-x>.
- Yu C, Velu A, Vinitsky E. et al. The surprising effectiveness of PPO in cooperative multi-agent games. In: *Neural Information Processing Systems*. Red Hook, NY, USA: Curran Associates, Inc. 2021.
- Yurtsev EA, Chao HX, Datta MS et al. Bacterial cheating drives the population dynamics of cooperative antibiotic resistance plasmids. *Mol Syst Biol* 2013;9:683. <https://doi.org/10.1038/msb.2013.39>.
- Zengler K, Zaramela LS. The social network of microorganisms—how auxotrophies shape complex communities. *Nat Rev Microbiol* 2018; 16:383–90. <https://doi.org/10.1038/s41579-018-0004-5>.
- Zhang D, Wang Y, Zheng D et al. New combination of xylanolytic bacteria isolated from the lignocellulose degradation microbial consortium XDC-2 with enhanced xylanase activity. *Bioresour Technol* 2016;221:686–90. <https://doi.org/10.1016/j.biotech.2016.09.087>.
- Zhao X, Noack S, Wiechert W et al. Dynamic flux balance analysis with nonlinear objective function. *J Math Biol* 2017;75:1487–515. <https://doi.org/10.1007/s00285-017-1127-4>.
- Zomorrodi AR, Islam MM, Maranas CD. D-OptCom: dynamic multi-level and multi-objective metabolic modeling of microbial communities. *ACS Synth Biol* 2014;3:247–57. <https://doi.org/10.1021/sb4001307>.
- Zomorrodi AR, Maranas CD. OptCom: a multi-level optimization framework for the metabolic modeling and analysis of microbial communities. *PLoS Comput Biol* 2012;8:e1002363. <https://doi.org/10.1371/journal.pcbi.1002363>.
- Zomorrodi AR, Segrè D. Genome-driven evolutionary game theory helps understand the rise of metabolic interdependencies in microbial communities. *Nat Commun* 2017;8:1563. <https://doi.org/10.1038/s41467-017-01407-5>.