

Prediction of fitness in bacteria with causal jump dynamic mode decomposition

Shara Balakrishnan, Aqib Hasnain, Nibodh Boddupalli, Dennis M. Joshy, Robert G. Egbert, and Enoch Yeung

Abstract—In this paper, we consider the problem of learning a predictive model for population cell growth dynamics as a function of the media conditions. We first introduce a generic data-driven framework for training operator-theoretic models to predict cell growth rate. We then introduce the experimental design and data generated in this study, namely growth curves of *Pseudomonas putida* as a function of casein and glucose concentrations. We use a data driven approach for model identification, specifically the nonlinear autoregressive (NAR) model to represent the dynamics. We show theoretically that Hankel DMD can be used to obtain a solution of the NAR model. We show that it identifies a constrained NAR model and to obtain a more general solution, we define a causal state space system using 1-step, 2-step, ..., τ -step predictors of the NAR model and identify a Koopman operator for this model using extended dynamic mode decomposition. The hybrid scheme we call causal-jump dynamic mode decomposition, which we illustrate on a growth profile or fitness prediction challenge as a function of different input growth conditions. We show that our model is able to recapitulate training growth curve data with 96.6% accuracy and predict test growth curve data with 91% accuracy.

I. INTRODUCTION

One of the most fundamental processes in life is the ability to replicate and pass on hereditary material [1]. From viral particles to bacteria to mammalian cells, cell division is fundamental to growth, maintenance of physiological health, and intrinsically tied to the notion of senescence [2].

The mechanisms for controlling growth in organisms are determined by metabolic networks [3], [4], namely their topological structure and parametric realization. Known metabolic networks in well studied model organisms such as *E. coli* [5] and *S. cerevisiae* [6], [7] have given rise to predictive models that translate environmental activity to metabolic network state, and ultimately to predictions of growth rate. For canonical biological model systems, these models have been highly accurate in predicting growth rate and found utility in industrial microbiology applications, e.g. in the design of bioreactors or informing best practices in food safety.

Shara Balakrishnan (sbalakrishnan@ucsb.edu) is with the Department of Electrical and Computer Engineering, University of California, Santa Barbara

Aqib Hasnain, Nibodh Boddupalli and Dennis Joshy are with the Department of Mechanical Engineering, University of California, Santa Barbara, Robert G. Egbert is with the Environmental and Biological Sciences Directorate at the Pacific Northwest National Laboratory, Richland, WA.

Enoch Yeung (eyeung@ucsb.edu) is with the Department of Mechanical Engineering, Center for Control, Dynamical Systems, and Computation, and Biomolecular Science and Engineering, University of California, Santa Barbara

For many biological life forms, relatively little is known about their metabolic network or structure. This is especially the case when developing bioengineering tools in novel host microbes [8], [9]. For new organisms, canonical metabolic networks are lacking and often obtained through a process of sequence alignment and comparative analysis with existing metabolic network models in relative species. However, many novel strains do not exhibit significant similarity, and even in the case of sequence similarity, small mutations can lead to dramatically different growth phenotypes, e.g. growth of non-pathogenic soil strains [10], [11] versus pathogenic counterparts [12]. The absence of predictive cross-species models, as well as the inability to predict growth phenotype wholly from sequence data, motivates the need for data-driven methods to accelerate the discovery of metabolic models and growth rate prediction models.

Due to advances in high-throughput experimental techniques, it is relatively easy to characterize growth rates as a function of exposure to environment. Liquid and acoustic-liquid handling robotics enables interrogation of thousands of growth conditions in a single microtiter plate, which in turn opens the door for using data-driven approaches [13] to predict growth rate as a function of environmental state. Is it possible to accurately predict the growth rate of a microbe, entirely from the chemical composition and environmental parameters of its growth condition? In this paper we explore a data-driven operator theoretic approach that utilizes microtiter plate reader data, and more generally multi-variate time-series data, to develop predictive models of growth rate in *Pseudomonas putida*, a broadly used strain for commercial bioreactors and a target workhorse for tractable genetic engineering [14].

A broadly successful class of data-driven modeling approaches stem from the study of Koopman operators, a mathematical construct for representing the time-evolution of nonlinear dynamical systems. In Koopman operator theory, the time-evolution of a nonlinear system is defined on a function space, acting on the original state of the system. In this function space, known the observables space, the Koopman operator is a linear operator, enabling spectral analysis, the decomposition of eigenspaces, and study of nonlinear structure [15]. The Koopman operator framework has been developed for continuous [16] and discrete time systems [17], [18], for open-loop [17] and input-controlled [19], [20] dynamic systems. Thus, Koopman operators present a powerful framework for analyzing the behavior of nonlinear

systems, including predicting how experimental conditions regulate growth dynamics.

Many numerical methods for identifying Koopman operators from data have been developed in the last two decades [21]–[28]. The most common approach is to use dynamic mode decomposition (DMD), which models nonlinear dynamics via an approximate local linear expansion [25]. In [17] an extended dictionary of basis functions with universal function approximation properties is used to discover an approximation of the lifting map or observables. These techniques suffer from combinatorial explosion, which generally has prohibited analysis of high-dimensional nonlinear systems [29]. The most recent developments in the field of DMD integrate established advances in deep learning with DMD [27], [30]–[32] where the deep neural networks have the capacity to approximate exponentially many distinct observable functions. Recent work has shown deep Koopman learning algorithms can be extended to synthesize controllers for systems subject to uncertainty [33], suggesting that deep Koopman learning can be used broadly for robust controller synthesis.

The existing algorithms assume a full state measurement and construct observables from that. With partial state observables like in biological systems, we construct a state space model for an output nonlinear difference equation model and identify a Koopman operator for that model. In Section II we briefly introduce the Koopman operator and DMD and the existing literature. In Section III we describe the experimental setup for obtaining the growth curve data of *Pseudomonas putida* by adding different concentrations of casein and glucose substrates to the media. In Section IV, we justify nonlinear autoregressive difference equation model is an appropriate choice for this system. In Section V, we formulate the Hankel DMD as a solution of the NAR model and bring out its issues and in Section VI, we formulate a state space model for the NAR model and use extended dynamic mode decomposition to identify a Koopman operator for the NAR model. In Section VII, we show that the algorithm is able to train a predictive Koopman operator, that predicts with 3.4% on the training data and 9% on the test data on extended forecasting tasks approximately 500 time steps ahead.

II. MATHEMATICAL PRELIMINARIES

Consider a discrete-time autonomous nonlinear dynamical system

$$x_{k+1} = f(x_k) \quad (1)$$

with $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is analytic. There exists a Koopman operator [34] of (1), which acts on a function space \mathcal{F} as $\mathcal{K} : \mathcal{F} \rightarrow \mathcal{F}$. This action can be given by

$$\mathcal{K}\psi(x_k) = \psi \circ f(x_k). \quad (2)$$

where the function $\psi : \mathbb{R}^n \rightarrow \mathbb{R}$ is called an *observable* of the system and the set of all observables $\psi \triangleq \{\psi_i\}_{i=1}^p, p \leq \infty$ on the system. Here \mathcal{F} is invariant under the action of \mathcal{K} .

The most important property of the Koopman operator that we utilize is the linearity of the operator, in other words,

$$\mathcal{K}(\alpha\psi_1 + \beta\psi_2) = \alpha\psi_1 \circ f + \beta\psi_2 \circ f = \alpha\mathcal{K}\psi_1 + \beta\mathcal{K}\psi_2$$

which follows from (2) since the composition operator is linear. Thus, we have that the Koopman operator of (1) is a linear operator that acts on observable functions $\psi(x_k)$ and propagates them forward in time.

A. DMD and relevant variants

The practical identification of Koopman operator for a nonlinear system from input-output data is commonly done using DMD [25] or extended DMD [17] which constructs an approximate Koopman operator K . Rowley et. al showed that the finite-dimensional approximation to the Koopman operator obtained from DMD is closely related to a spectral analysis of the linear but infinite-dimensional Koopman operator [18]. The approach taken to compute an approximation to the Koopman operator in both DMD and extended DMD is as follows

$$K = \min_K \|\Psi(X_f) - K\Psi(X_p)\| = \Psi(X_f)\Psi(X_p)^\dagger \quad (3)$$

where $X_f \equiv [x_1 \ \dots \ x_{N-1}]$, $X_p \equiv [x_2 \ \dots \ x_N]$ are snapshot matrices formed from the discrete-time dynamical system (1), $\Psi(X) \equiv [\psi_1(x) \ \dots \ \psi_R(x)]$ is the mapping from physical space into the space of observables and † denotes the Moore-Penrose pseudoinverse. Here N is the number of snapshots i.e. timepoints. We note that DMD is a special case of extended DMD where $\psi(x) = x$. Throughout the rest of the paper, when we refer to the Koopman operator we mean the finite dimensional approximation to the infinite-dimensional Koopman operator.

III. EXPERIMENTAL SETUP

We describe the procedure adopted to obtain *P. putida*'s growth curve for varying concentrations of glucose and casein substrates in the media.

Incubation: We revived *P. putida* cryopreserved at -80°C in 30% (vol/vol) glycerol stock by suspending a small portion into a polypropylene test tube containing 4 mL of Lysogeny Broth (LB). We cultured it at 30°C spinning with a speed of 200 revolutions per minute (rpm) for 12 hours overnight. A visual inspection of the culture tube resulted in a cloudy culture medium, suggesting subsequent growth with the seed *P. putida* culture in a plate reader was feasible.

Solution Preparation: We prepared 300 g/L solution of glucose solution and 225 g/L casein acid hydrolysate solution. Once the bacteria reached a certain optical density(OD), we shifted the culture from LB media to R2A 2x media obtained from Teknova Inc to 2x the required initial OD.

Serial dilution setup for *P. putida* culture: We use a 630 μL 96 well plate to create media with different substrate concentrations. Each well of this plate contained 500 μL of modified media - 250 μL of culture in 2x R2A at 0.4 OD and 120 μL containing a mixture of casein and glucose solutions. To vary casein and glucose across the 96 well plate, we perform 2D serial dilution such that the concentration of

glucose was halved across columns and concentration of casein is halved across rows as shown in Figure 1. Then, the culture was mixed into each well to get a starting OD of 0.2 in 1x R2A media since equal volumes of culture media and substrate solutions were added.

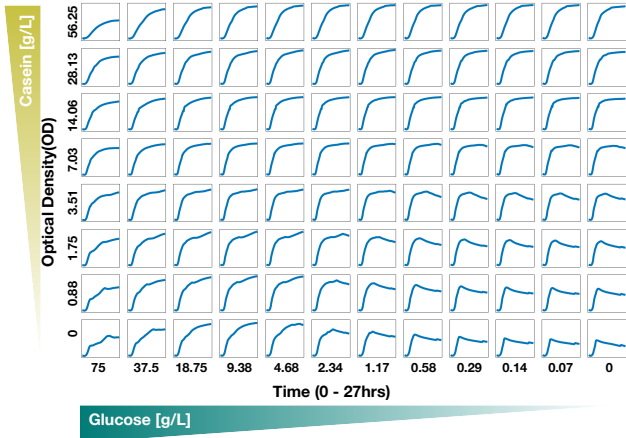


Fig. 1: Different initial conditions of substrates obtained by two dimensional serial dilution of casein and glucose and the corresponding growth curves are obtained for a period of 27 hours.

Data Collection: The microplate reader was set to 30°C and the shaker to 807 cycles per minute, with continuous double orbital mixing. The absorbance at 600 nanometers (nm), which is termed as the Optical Density at 600 nm (OD_{600}), was measured as a function of time for 27 hours. We assume in this work, as is widely accepted, that OD_{600} measurements were collected in a linear regime, where cell population is proportional OD_{600} measurements. The obtained data along with the varying substrate concentrations are shown in Figure. 1.

IV. GROWTH CURVE DYNAMICS MODEL

The dynamics of the bacterial cell growth can be represented by

$$\begin{bmatrix} N_{k+1}^{(b)} \\ C_{k+1} \\ G_{k+1} \end{bmatrix} = f(N_k^{(b)}, C_k, G_k) \quad (4)$$

where the bacterial cell count ($N^{(b)}$), casein substrate concentration (C) and glucose substrate concentration (G) are the states of the system, f is the nonlinear dynamics and k is the discrete time index. We measure the OD_{600} data as mentioned in section III and the output equation is given by

$$y_k = h(N_k^{(b)}) \quad (5)$$

as OD_{600} is a function of only the number of cells. Some of the existing empirical nonlinear models for growth curve dynamics include the Monod's model [35] which uses a single substrate to form the foundation of the growth curve dynamics and in [36] and [37] multiple substrates are incorporated. Monod's model is a two-state nonlinear dynamical

system comprising of the substrate (S) and the number of bacteria ($N^{(b)}$):

$$\begin{aligned} \dot{N}^{(b)}(t) &= r_{max} \frac{S(t)N^{(b)}(t)}{K_s + S(t)} \\ \dot{S} &= -\gamma \dot{N}^{(b)} \end{aligned} \quad (6)$$

where r_{max} is the maximum growth rate and K_s is the half velocity constant. As $N^{(b)}$ is the only variable of measurement in (5), we convert the model to a single differential equation containing only $N^{(b)}$

$$\begin{aligned} \ddot{N}^{(b)}(t) &= \frac{1}{r_{max}K_sN^{(b)}}(K_sr_{max}\dot{N}^{(b)2} - \gamma\dot{N}^{(b)3}) \\ &+ 2\gamma r_{max}N^{(b)}\dot{N}^{(b)2} - \gamma r_{max}^2N^{(b)2}\dot{N}^{(b)} \end{aligned} \quad (7)$$

The existing models though heuristic, suggest that $N^{(b)}$ at any point in time is a function of the past

$$N_{k+1}^{(b)} = f(N_k^{(b)}, N_{k-1}^{(b)}, \dots)$$

$N_k^{(b)}$ to be a function of its finite past. This is the general structure of the discrete nonlinear autoregressive (NAR) model given by

$$\begin{aligned} y_k &= f(y_{k-1}, y_{k-2}, \dots, y_{k-\tau}) \\ y_i &\in \mathbb{R}^p \quad \forall i \in \mathbb{Z}_{>0} \\ f &: \underbrace{\mathbb{R}^p \times \mathbb{R}^p \times \dots \times \mathbb{R}^p}_{\tau \text{ times}} \rightarrow \mathbb{R}^p \end{aligned} \quad (8)$$

where the current output is a function of the past τ outputs.

V. HANKEL DYNAMIC MODE DECOMPOSITION

Given the nonlinear system (4) with the state measurement given by (5) and modeled by the discrete time difference equation (8), Hankel DMD [38] is a suitable algorithm to solve the model identification problem with the NAR structure. The promising feature of using a DMD algorithm is that it identifies a linear state space representation which has a theoretical foundation in Koopman operator theory.

Given the autonomous state space system

$$\begin{aligned} \tilde{x}_{k+1} &= \tilde{f}(\tilde{x}_k) \\ y_k &= h(\tilde{x}_k) \end{aligned} \quad (9)$$

where $x_k \in \mathbb{R}^n$ is the state, $\tilde{f}: \mathbb{R}^n \rightarrow \mathbb{R}^n$ is the dynamics, $y_k \in \mathbb{R}^p$ is the output and $h: \mathbb{R}^n \rightarrow \mathbb{R}^p$ is a nonlinear function that maps the state directly to itself, i.e. x is identical to the output y , Hankel DMD constructs a Koopman model of the form

$$\begin{bmatrix} \psi(y_{k+1}) \\ \psi(y_{k+2}) \\ \vdots \\ \psi(y_{k+\tau}) \end{bmatrix} = K \begin{bmatrix} \psi(y_k) \\ \psi(y_{k+1}) \\ \vdots \\ \psi(y_{k+\tau-1}) \end{bmatrix} \quad (10)$$

such that $\psi: \mathbb{R}^p \rightarrow \mathbb{R}^{N_p}$ is the dictionary of state inclusive observables of the state \tilde{x}_k constructed by a nonlinear transformation of the corresponding output y_k and K is the Koopman operator. Regardless of full-state measurements,

we nonetheless cast Hankel DMD in this form to compare it with our subsequent causal jump DMD algorithm.

Given the output measurements $\{y_1, y_2, \dots, y_N\}$, to identify an approximate Koopman operator K using Hankel DMD, the time shifted Hankel matrices are constructed as

$$\Psi(Y_p) = \begin{bmatrix} \psi(y_1) & \psi(y_2) & \dots & \psi(y_{N-\tau}) \\ \psi(y_2) & \psi(y_3) & \dots & \psi(y_{N-\tau+1}) \\ \vdots & \vdots & \ddots & \vdots \\ \psi(y_\tau) & \psi(y_{\tau+1}) & \dots & \psi(y_{N-1}) \end{bmatrix} \quad (11)$$

$$\Psi(Y_f) = \begin{bmatrix} \psi(y_2) & \psi(y_3) & \dots & \psi(y_{N-\tau+1}) \\ \psi(y_3) & \psi(y_4) & \dots & \psi(y_{N-\tau+2}) \\ \vdots & \vdots & \ddots & \vdots \\ \psi(y_{\tau+1}) & \psi(y_{\tau+2}) & \dots & \psi(y_N) \end{bmatrix}$$

and the optimization problem

$$\min_K \|\Psi(Z_f) - K\Psi(Z_p)\| \quad (12)$$

is solved using the Moore-Penrose pseudoinverse method mentioned in section II. This yields a solution of the form

$$\begin{bmatrix} \psi(y_{k+1}) \\ \psi(y_{k+2}) \\ \vdots \\ \psi(y_{k+\tau}) \end{bmatrix} = \begin{bmatrix} 0 & \mathbb{I}_{N_p} & \dots & 0 & 0 \\ 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & \mathbb{I}_{N_p} \\ k_1 & k_2 & \dots & k_{\tau-1} & k_\tau \end{bmatrix} \begin{bmatrix} \psi(y_k) \\ \psi(y_{k+1}) \\ \vdots \\ \psi(y_{k+\tau-1}) \end{bmatrix}$$

Other than the last N_p equations, the others are trivial. To construct an output predictor, we take the component $y_{k+\tau}$ of $\psi(y_{k+\tau})$ to get

$$y_{k+\tau} = \tilde{k}_1\psi(y_k) + \tilde{k}_2\psi(y_{k+1}) + \dots + \tilde{k}_\tau\psi(y_{k+\tau-1}) \quad (13)$$

where \tilde{k}_i are the components of k_i that map $\psi(y_{k+\tau-1})$ to $y_{k+\tau}$. More generally, this yields a nonlinear equation of the form

$$y_k = \tilde{f}_1(y_{k-1}) + \tilde{f}_2(y_{k-2}) + \dots + \tilde{f}_\tau(y_{k-\tau}) \quad (14)$$

where the functions $\tilde{f}_1, \tilde{f}_2, \dots, \tilde{f}_\tau$ have the same basis functions with different coefficients. This identifies a constrained NAR model as it imposes an additive structure on the basis of nonlinear models across time.

VI. DYNAMIC MODE DECOMPOSITION OF NONLINEAR AUTOREGRESSIVE MODELS

To identify a Koopman operator for the unconstrained NAR model (8), we formulate a state space representation for the NAR model with full state observation and identify an approximate Koopman operator for that model using the general class of DMD algorithms like extended DMD and deep DMD.

In this methodology, the problem is broken into two pieces the system identification aspect where we select the model structure and the dynamic mode decomposition aspect where we have to construct the dictionary of observables. We define a window parameter $\tau \in \mathbb{Z}_{>0}$ indicating how many past output snapshots are used to define a new extended dictionary of monomial observable functions, up to order $n_o \in \mathbb{Z}_{>0}$.

The new τ -dictionary defines a general extended dynamic mode decomposition problem, which we then solve using classical methods.

We proceed as follows: given the NAR model (8) with the system identification parameter τ , we construct a state defined by

$$z_k := [y_{k+1} \quad y_{k+2} \quad \dots \quad y_{k+\tau}]^T \quad (15)$$

with $z_k \in \mathbb{R}^{p\tau}$. This yields the state space representation

$$z_{k+1} = \begin{bmatrix} y_{k+2} \\ y_{k+3} \\ \vdots \\ y_{k+\tau} \\ y_{k+\tau+1} \end{bmatrix} := \begin{bmatrix} f_1(y_{k+1}, y_{k+2}, \dots, y_{k+\tau}) \\ f_2(y_{k+1}, y_{k+2}, \dots, y_{k+\tau}) \\ \vdots \\ f_{\tau-1}(y_{k+1}, y_{k+2}, \dots, y_{k+\tau}) \\ f_\tau(y_{k+1}, y_{k+2}, \dots, y_{k+\tau}) \end{bmatrix}$$

$$:= \begin{bmatrix} y_{k+2} \\ y_{k+3} \\ \vdots \\ y_{k+\tau} \\ f(y_{k+1}, y_{k+2}, \dots, y_{k+\tau}) \end{bmatrix} = F(z_k)$$

$$\Rightarrow z_{k+1} = \tilde{F}(z_k) \quad (16)$$

where $\tilde{F} : \mathbb{R}^{p\tau} \rightarrow \mathbb{R}^{p\tau}$ represents the dynamics of the lifted "state" model. The approximate EDMD model for the full output observable model is given by

$$\psi(z_{k+1}) = K\psi(z_k) \quad (17)$$

where $\psi(z_k)$ is the state inclusive dictionary of observables defined as

$$\psi(z_k) = \begin{bmatrix} z_k \\ \varphi(z_k) \end{bmatrix} \quad (18)$$

with $\varphi : \mathbb{R}^{p\tau} \rightarrow \mathbb{R}^{N_p}$ being a nonlinear transformation that constructs the nonlinear observables. Since the only additional information in the state z_{k+1} when compared to the state z_k is $y_{k+\tau+1}$, the output predictor form for the Koopman model can be identified considering the complete Koopman model and extracting the equation that corresponds to $y_{k+\tau+1}$ given by

$$\psi(z_{k+1}) = K\psi(z_k)$$

$$\Rightarrow \begin{bmatrix} y_{k+2} \\ \vdots \\ y_{k+\tau} \\ y_{k+\tau+1} \\ \varphi(z_{k+1}) \end{bmatrix} = \begin{bmatrix} \bullet & \dots & \bullet & \bullet & \bullet \\ \vdots & \ddots & \vdots & \vdots & \vdots \\ \bullet & \dots & \bullet & \bullet & \bullet \\ k_1 & \dots & k_{\tau-1} & k_\tau & k_{11} \\ \bullet & \dots & \bullet & \bullet & \bullet \end{bmatrix} \begin{bmatrix} y_{k+1} \\ \vdots \\ y_{k+\tau-1} \\ y_{k+\tau} \\ \varphi(z_k) \end{bmatrix}$$

$$\Rightarrow y_{k+\tau+1} = k_1 y_{k+1} + \dots + k_\tau y_{k+\tau} + k_{11}^T \varphi(y_{k+1}, \dots, y_{k+\tau}) \quad (19)$$

The output predictor form keeps the general structure of the NAR model intact as opposed to the predictor identified by Hankel DMD which identified a constrained model. But, the issue with this model is the causality. It can be seen from (19) that the Koopman model is non causal due to the overlap of outputs y_k between the states z_{k+1} and z_k . This identifies models that use future outputs to predict past outputs which

are inadmissible as our system is causal. To identify a causal model, the property of (16) proved in proposition 1 is very important.

Proposition 1: Given the state space model (16) for the nonlinear autoregressive (NAR) model (8), if the state is propagated i time steps where $i \in \{1, 2, \dots, \tau\}$

$$z_{k+i} = \tilde{F}^i(z_k) = \underbrace{\tilde{F} \circ \tilde{F} \circ \dots \circ \tilde{F}}_{i \text{ times}}(z_k),$$

then the last i functions of $\tilde{F}_H^i(z_k)$ are such that

$$\begin{aligned} (\tilde{F}^i)^{(\tau-i+j)}(z_k) &= f^{(j)}(z_k) \quad j \in \{1, 2, \dots, i\} \\ y_{y+\tau+j} &= f^{(j)}(z_k) = f^{(j)}(y_{k+1}, y_{k+2}, \dots, y_{k+\tau}) \end{aligned} \quad (20)$$

where $(\tilde{F}^i)^{(b)}(z_k)$ corresponds to the b^{th} function of $\tilde{F}^i(z_k)$ and $f^{(j)}(z_k)$ is the j -step predictor of the NAR model (8).

Proof: Given the state z_k defined in (15), the state propagated i time steps $\forall i \in \mathbb{Z}_{\geq 0}$ is given by

$$z_{k+i} = [y_{k+i+1} \quad y_{k+i+2} \quad \dots \quad y_{k+i+\tau}]^T$$

and the m^{th} component of z_{k+i} is given by $z_{k+i}^{(m)} = y_{k+i+m}$ where $m \in \{1, 2, \dots, \tau\}$.

A function $f^{(j)} : \underbrace{\mathbb{R}^p \times \mathbb{R}^p \times \dots \times \mathbb{R}^p}_{\tau \text{ times}} \rightarrow \mathbb{R}^p$ is a j -step predictor of the NAR model (8) if it has the following form

$$y_{k+\tau+j} = f^{(j)}(x_k) = f^{(j)}(y_{k+1}, y_{k+2}, \dots, y_{k+\tau}).$$

Now that we have the state definitions and the predictor function definitions in place, we prove (20) by induction. For $i = 1$,

$$\begin{aligned} z_{k+1} &= \tilde{F}^{(1)}(z_k) \\ (\tilde{F}^1)^{(\tau-i+j)}(z_k) &= (\tilde{F}^1)^{(\tau)}(z_k) = f^{(1)}(x_k) \quad j \in \{1\} \\ \Rightarrow z_{k+1}^{(\tau)} &= y_{k+\tau+1} = f^{(1)}(z_k) \end{aligned}$$

Hence (20) is satisfied for $i = 1$. We assume the result is true for $i = p$. This yields

$$\begin{aligned} (\tilde{F}^p(z_k))^{(\tau-p+j)}(z_k) &= f^{(j)}(z_k) \quad j \in \{1, 2, \dots, p\} \\ \Rightarrow z_{k+p} &= \begin{bmatrix} y_{k+p+1} \\ \vdots \\ y_{k+\tau} \\ y_{k+\tau+1} \\ \vdots \\ y_{k+\tau+p} \end{bmatrix} = \tilde{F}^p(z_k) = \begin{bmatrix} y_{k+p+1} \\ \vdots \\ y_{k+\tau} \\ f^{(1)}(z_k) \\ \vdots \\ f^{(p)}(z_k) \end{bmatrix}. \end{aligned}$$

For $i = p + 1$, the state z_{k+p+1} becomes

$$\begin{aligned} z_{k+p+1} &= \tilde{F}^{p+1}(z_k) = \tilde{F} \circ \tilde{F}^p(z_k) \\ \Rightarrow \begin{bmatrix} y_{k+p+2} \\ \vdots \\ y_{k+\tau} \\ y_{k+\tau+1} \\ \vdots \\ y_{k+\tau+p} \\ y_{k+\tau+p+1} \end{bmatrix} &= \tilde{F} \left(\begin{bmatrix} y_{k+p+1} \\ \vdots \\ y_{k+\tau} \\ f^{(1)}(z_k) \\ \vdots \\ f^{(p-1)}(z_k) \\ f^{(p)}(z_k) \end{bmatrix} \right) = \begin{bmatrix} y_{k+p+2} \\ \vdots \\ y_{k+\tau} \\ f^{(1)}(z_k) \\ \vdots \\ f^{(p)}(z_k) \\ g \end{bmatrix} \end{aligned}$$

where

$$\begin{aligned} g &= f(y_{k+p+1}, \dots, y_{k+\tau}, f^{(1)}(z_k), \dots, f^{(p)}(z_k)) \\ &= f(z_k^{(p+1)}, \dots, z_k^{(\tau)}, f^{(1)}(z_k), \dots, f^{(p)}(z_k)) \\ &:= g(z_k). \end{aligned}$$

Since g is a function of only z_k and since $y_{k+\tau+p+1} = g, g(z_k)$ satisfies the definition of a predictor function and hence is a $(p + 1)$ -step predictor of (8)

$$y_{k+\tau+p+1} = z_{k+p+1}^{(\tau)} = (\tilde{F}^{p+1}(z_k))^{(\tau)} = f^{(p+1)}(z_k).$$

Therefore, for $i = p + 1$,

$$(\tilde{F}^i(z_k))^{(\tau-p-1+j)} = f^{(j)}(z_k) \quad j \in \{1, 2, \dots, (p + 1)\}$$

stating that the last $(p + 1)$ entries of z_{k+p+1} are $f^{(1)}(x), f^{(2)}(x), \dots, f^{(p+1)}(x)$. Hence the proof. \blacksquare

To identify a causal Koopman model for the NAR system (8), we propagate the model (16) by τ time steps to ensure no intersection of outputs between the states z_{k+1} and z_k . We define a new state $x_k = z_{k\tau}$ which yields

$$\begin{aligned} x_k &= z_{k\tau} \\ \Rightarrow x_{k+1} &= z_{k\tau+\tau} = \tilde{F}^{\tau}(z_{k\tau}) = F(x_k) \quad (21) \\ \text{where } F &= \tilde{F}^{\tau} = \underbrace{\tilde{F} \circ \tilde{F} \circ \dots \circ \tilde{F}}_{\tau \text{ times}} \end{aligned}$$

Using proposition 1, we can say that the nonlinear state space model contains functions that are 1-step, 2-step, ..., τ -step predictors in the following form

$$\begin{aligned} x_{k+1} &= \begin{bmatrix} y_{k\tau+\tau+1} \\ y_{k\tau+\tau+2} \\ \vdots \\ y_{k\tau+2\tau-1} \\ y_{k\tau+2\tau} \end{bmatrix} := \begin{bmatrix} f_1(y_{k\tau+1}, y_{k\tau+2}, \dots, y_{k\tau+\tau}) \\ f_2(y_{k\tau+1}, y_{k\tau+2}, \dots, y_{k\tau+\tau}) \\ \vdots \\ f_{\tau-1}(y_{k\tau+1}, y_{k\tau+2}, \dots, y_{k\tau+\tau}) \\ f_{\tau}(y_{k\tau+1}, y_{k\tau+2}, \dots, y_{k\tau+\tau}) \end{bmatrix} \\ &:= \begin{bmatrix} f^{(1)}(y_{k\tau+1}, y_{k\tau+2}, \dots, y_{k\tau+\tau}) \\ f^{(2)}(y_{k\tau+1}, y_{k\tau+2}, \dots, y_{k\tau+\tau}) \\ \vdots \\ f^{(\tau-1)}(y_{k\tau+1}, y_{k\tau+2}, \dots, y_{k\tau+\tau}) \\ f^{(\tau)}(y_{k\tau+1}, y_{k\tau+2}, \dots, y_{k\tau+\tau}) \end{bmatrix} = F(x_k) \quad (22) \end{aligned}$$

where $f^{(i)}$ is the i -step predictor of the NAR model. We prove the existence of a Koopman operator for this model in Proposition 2.

Proposition 2: If the function $f(x)$ in the NAR model (8) is analytic, then a Koopman operator exists for (16) and (22).

Proof: Since f in (8) is analytic, \tilde{F} in (16) is analytic since all the entries of \tilde{F} are either linear functions or are equal to f . Since F is obtained by the composition of \tilde{F} τ times, F is also analytic.

$F(x)$ admits a countable-dimension Koopman operator K_x , with an invariant subspace isomorphic to either a finite or an infinite Taylor polynomial basis [34]. Moreover, isomorphism with a Taylor polynomial basis ensures that the

Koopman observable space contains the full state observable, i.e. it is state inclusive.

There are two easy arguments to conclude the proof. First, note that since f is analytic, f^τ is analytic and thus by the same reasoning as in [34], f^τ thus must admit a Koopman operator. The second argument is a constructive one, noting that equation

$$\psi(x[(k)\tau]) = K^\tau \psi(x[(k-1)\tau]) \quad (23)$$

must hold due to τ applications of the 1-step Koopman equation. This means therefore that the following *matrix* equation must hold

$$\psi \left(\begin{bmatrix} x[(k)\tau] \\ x[k\tau+1] \\ \vdots \\ x[(k+1)\tau-1] \end{bmatrix} \right) = \mathbf{K}_J \psi \left(\begin{bmatrix} x[(k-1)\tau] \\ x[(k-1)\tau+1] \\ \vdots \\ x[(k)\tau-1] \end{bmatrix} \right) \quad (24)$$

where $\mathbf{K}_J = \text{diag}(K^\tau, K^\tau, \dots, K^\tau)$. This concludes the proof. ■

Since the existence of a Koopman operator has been proved for the model (22) in Proposition 2, we construct a state inclusive dictionary of observables

$$\psi(x_k) = \begin{bmatrix} x_k \\ \varphi(x_k) \end{bmatrix} \quad (25)$$

with $\varphi: \mathbb{R}^{p\tau} \rightarrow \mathbb{R}^{N_p}$ to define a Koopman model

$$\psi(x_{k+1}) = K\psi(x_k) \quad (26)$$

This Koopman model is causal since there is no intersection of outputs between x_{k+1} and x_k . The added feature of this model is that the DMD algorithm while identifying a Koopman operator, also simultaneously minimizes the 1-step, 2-step, ..., τ -step prediction error of the NAR model.

Now that we have a theoretical state space representation of a NAR model and established the conditions under which a Koopman operator exists, we turn our attention to the algorithm for identification of the Koopman operator. Given the data with M data sets and N data points in each data set $\{y_1^{(i)}, y_2^{(i)}, \dots, y_N^{(i)}\}$ where $i \in \{1, 2, \dots, M\}$ is the index of the data set, we construct the Hankel states z_k and the dictionary of observables allowing the intermixing of states. We compile the observables into snapshot matrices $\tilde{\Psi}_f(z)$ and $\tilde{\Psi}_p(z)$ with a τ time step jump and solve the Koopman learning problem

$$\|\tilde{\Psi}_f(z) - K\tilde{\Psi}_p(z)\|_F$$

using the methodology in Algorithm 1.

VII. RESULTS

From the data-sets obtained in the plate reader experiments shown in Fig. 1, we used Algorithm 1 to implement extended DMD using monomials as the dictionary of observables

$$\psi(z_k) = [y_{k+1}, \dots, y_{k+\tau}, y_{k+1}^2, y_{k+1}y_{k+2}, \dots, y_{k+\tau}^2, y_{k+1}^3, y_{k+1}^2y_{k+2}, \dots]^T.$$

Algorithm 1 Extended DMD for NAR models

- 1: Get NAR model parameter τ
- 2: Get extended DMD parameter n_o for monomial observables
- 3: **for** dataset $i = 1, 2, \dots, M$ **do**
- 4: **for** time index $j = 1, 2, \dots, N - \tau$ **do**
- 5: Construct the Hankel state

$$z_j^{(i)} = [y_{j+1}^{(i)} \quad y_{j+2}^{(i)} \quad \dots \quad y_{j+\tau}^{(i)}]$$

- 6: Construct the dictionary of observables $\psi(z_j^{(i)})$
- 7: **end for**
- 8: Construct the snapshot matrices for each data set with the τ -jump

$$\Psi_p^{(i)}(x) = [\psi(z_1^{(i)}) \quad \psi(z_2^{(i)}) \quad \dots \quad \psi(z_{N-2\tau}^{(i)})]$$

$$\Psi_f^{(i)}(x) = [\psi(z_{1+\tau}^{(i)}) \quad \psi(z_{2+\tau}^{(i)}) \quad \dots \quad \psi(z_{N-\tau}^{(i)})]$$

- 9: **end for**
- 10: Compile the snapshot matrices across data sets

$$\tilde{\Psi}_p(x) = [\Psi_p^{(1)}(x) \quad \Psi_p^{(2)}(x) \quad \dots \quad \Psi_p^{(M)}(x)]$$

$$\tilde{\Psi}_f(x) = [\Psi_f^{(1)}(x) \quad \Psi_f^{(2)}(x) \quad \dots \quad \Psi_f^{(M)}(x)]$$

- 11: Compute the SVD of $\tilde{\Psi}_p(x) = USV^*$
- 12: Truncate to required number of singular values and identify the Koopman operator

$$\hat{K} = \tilde{\Psi}_f(x)\tilde{V}\tilde{S}^{-1}\tilde{U}^*$$

to identify an approximate Koopman operator for the state space model (22) as a solution to the identification of the NAR model (8).

We use all the datasets in Fig. 1 to find a Koopman operator invariant to the substrate concentrations. They are broken equally into training, validation and test set. Given the two parameters τ (NAR model parameter) and n_o (extended DMD parameter), we can find the optimal approximate Koopman operator by cumulatively iterating through the principal components and evaluating the summation of the mean squared error(MSE) of training and validation data. The number of principal components corresponding to the minimum MSE yields the optimal Koopman operator for a given τ and n_o . We then iterate through the two parameters to find the optimal model that minimizes the

By choosing $\tau = 9$ and keeping the maximum order of monomials to 3, the Koopman operator has been identified and the prediction on the training data is shown in Figure 2 and it has an MSE of 3.4%. The identified Koopman operator has an MSE of 9% and the fit is shown in Figure 3.

The results on the experimental data suggest that Causal Jump DMD is a suitable candidate algorithm for identifying the Koopman operator of the population growth dynamics of bacteria and can also be extended in general to identify Koopman operators for NAR models.

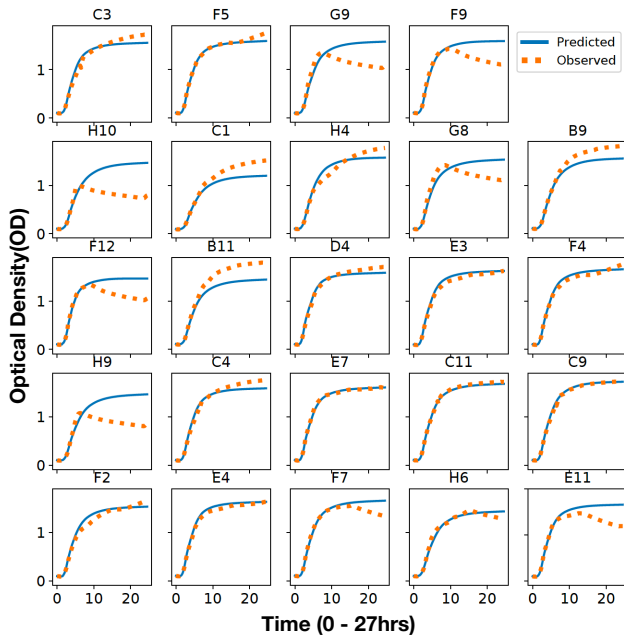


Fig. 2: The identified Koopman operator is tested on the training sets with 9 point initial condition and up to 3^{rd} order monomials to get a MSE of 3.4%

VIII. CONCLUSION

In this paper, we introduced the microbial growth curve dynamics to motivate the usage of DMD algorithms to identify Koopman operators for NAR models. We formulated Hankel DMD as a state space representations of the NAR model and showed that it is restrictive in its structure. We construct a causal state space model for the NAR model and identify a Koopman operator for it using extended dynamic mode decomposition with a monomial dictionary of observables. We showed that it does a good job in predicting the population growth dynamics of *Pseudomonas putida* invariant to substrate concentrations. The future goals of this work is to use this model to identify the optimal media conditions for maximal and minimal growth of the microbe thereby enabling us to develop a general methodology to develop an external growth harness for microbes for dynamic growth control. To achieve this, we need to extend the mathematical models to allow for inputs and extend the identification to NARX and NARMAX models. Further, if we integrate this framework with deepDMD which aids in finding the observable functions in a parsimonious fashion, it renders a useful tool for identifying high dimensional linear models for nonlinear systems.

IX. ACKNOWLEDGMENTS

The authors gratefully acknowledge the funding of DARPA grants FA8750-17-C-0229, HR001117C0092, HR001117C0094, DEAC0576RL01830. The authors would also like to thank Professors Arun K. Tangirala, Igor Mezic, Milan Korda, Alexandre Mauroy, Nathan Kutz, Steve Haase, John Harer, Devin Strickland, and Eric Klavins for insightful discussions. Any opinions, findings, conclusions,

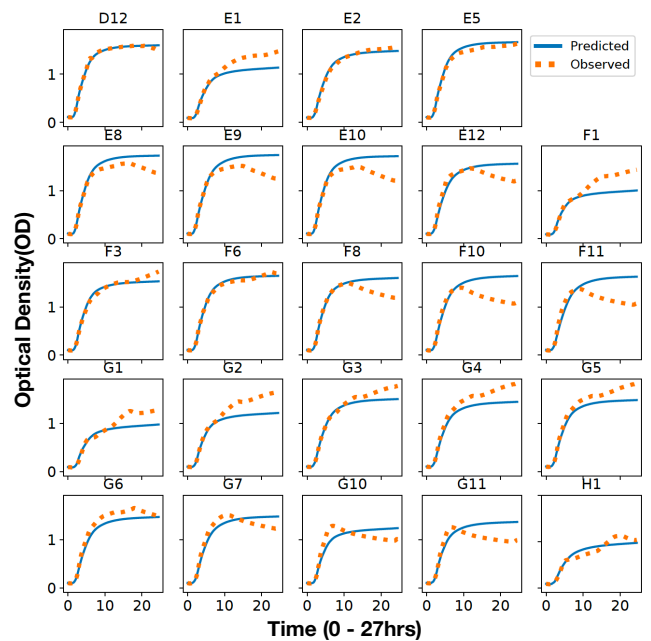


Fig. 3: The identified Koopman operator is tested on the test sets by using the initial observables $\psi(x_0)$ and the mean squared error remains the same as that of the training set.

or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the Defense Advanced Research Project Agency, the Department of Defense, or the United States government. This material is based on work supported by DARPA and AFRL under contract numbers FA8750-17-C-0229, HR001117C0092, HR001117C0094, DEAC0576RL01830.

REFERENCES

- [1] D. Kornberg and D. TA, "Replication," *San Francisco: W H. Freeman*, 1980.
- [2] N. F. Mathon and A. C. Lloyd, "Cell senescence and cancer," *Nature Reviews Cancer*, vol. 1, no. 3, p. 203, 2001.
- [3] G. Wu, Q. Yan, J. A. Jones, Y. J. Tang, S. S. Fong, and M. A. Koffas, "Metabolic burden: cornerstones in synthetic biology and metabolic engineering applications," *Trends in biotechnology*, vol. 34, no. 8, pp. 652–664, 2016.
- [4] D. S. Glazier, "Is metabolic rate a universal pacemaker for biological processes?" *Biological Reviews*, vol. 90, no. 2, pp. 377–407, 2015.
- [5] D. De Martino, F. Capuani, and A. De Martino, "Growth against entropy in bacterial metabolism: the phenotypic trade-off behind empirical growth rate distributions in e. coli," *Physical biology*, vol. 13, no. 3, p. 036005, 2016.
- [6] B. J. Sanchez, C. Zhang, A. Nilsson, P.-J. Lahtvee, E. J. Kerkhoven, and J. Nielsen, "Improving the phenotype predictions of a yeast genome-scale metabolic model by incorporating enzymatic constraints," *Molecular systems biology*, vol. 13, no. 8, 2017.
- [7] M. Zwietering, I. Jongenburger, F. Rombouts, and K. Van't Riet, "Modeling of the bacterial growth curve," *Appl. Environ. Microbiol.*, vol. 56, no. 6, pp. 1875–1881, 1990.
- [8] T. Tschirhart, V. Shukla, E. E. Kelly, Z. Schultzhaus, E. NewRingeisen, J. S. Erickson, Z. Wang, W. Garcia, E. Curl, R. G. Egbert *et al.*, "Synthetic biology tools for the fast-growing marine bacterium *Vibrio natriegens*," *ACS synthetic biology*, 2019.
- [9] N. Khan, E. Yeung, Y. Farris, S. J. Fansler, and H. C. Bernstein, "A broad-host-range event detector: expanding and quantifying performance across bacterial species," *bioRxiv*, p. 369967, 2018.
- [10] C. Gill and K. Tan, "Effect of carbon dioxide on growth of *Pseudomonas fluorescens*," *Appl. Environ. Microbiol.*, vol. 38, no. 2, pp. 237–240, 1979.

- [11] D. M. Gulliver, G. V. Lowry, and K. B. Gregory, "Comparative study of effects of co2 concentration and ph on microbial communities from a saline aquifer, a depleted oil reservoir, and a freshwater aquifer," *Environmental Engineering Science*, vol. 33, no. 10, pp. 806–816, 2016.
- [12] A. E. LaBauve and M. J. Wargo, "Growth and laboratory maintenance of *Pseudomonas aeruginosa*," *Current protocols in microbiology*, vol. 25, no. 1, pp. 6E–1, 2012.
- [13] A. P. Palacios, J. M. Marín, E. J. Quinto, M. P. Wiper *et al.*, "Bayesian modeling of bacterial growth for multiple populations," *The Annals of Applied Statistics*, vol. 8, no. 3, pp. 1516–1537, 2014.
- [14] H. H. Lee, N. Ostrov, B. G. Wong, M. A. Gold, A. Khalil, and G. M. Church, "Vibrio natriegens, a new genomic powerhouse," *bioRxiv*, p. 058487, 2016.
- [15] I. Mezic, "Spectral properties of dynamical systems, model reduction and decompositions," *Nonlinear Dynamics*, vol. 41, no. 1-3, pp. 309–325, 2005.
- [16] M. Budišić, R. Mohr, and I. Mezić, "Applied koopmanism," *Chaos: An Interdisciplinary Journal of Nonlinear Science*, vol. 22, no. 4, p. 047510, 2012.
- [17] M. O. Williams, I. G. Kevrekidis, and C. W. Rowley, "A data-driven approximation of the koopman operator: Extending dynamic mode decomposition," *Journal of Nonlinear Science*, vol. 25, no. 6, pp. 1307–1346, 2015.
- [18] C. W. Rowley, I. Mezić, S. Bagheri, P. Schlatter, and D. S. Henningson, "Spectral analysis of nonlinear flows," *Journal of fluid mechanics*, vol. 641, pp. 115–127, 2009.
- [19] J. L. Proctor, S. L. Brunton, and J. N. Kutz, "Dynamic mode decomposition with control," *SIAM Journal on Applied Dynamical Systems*, vol. 15, no. 1, pp. 142–161, 2016.
- [20] M. O. Williams, M. S. Hemati, S. T. Dawson, I. G. Kevrekidis, and C. W. Rowley, "Extending data-driven koopman analysis to actuated systems," *IFAC-PapersOnLine*, vol. 49, no. 18, pp. 704–709, 2016.
- [21] T. Askham and J. N. Kutz, "Variable projection methods for an optimized dynamic mode decomposition," *SIAM Journal on Applied Dynamical Systems*, vol. 17, no. 1, pp. 380–416, 2018.
- [22] Y. Kaneko, S. Muramatsu, H. Yasuda, K. Hayasaka, Y. Otake, S. Ono, and M. Yukawa, "Convolutional-sparse-coded dynamic mode decomposition and its application to river state estimation," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 1872–1876.
- [23] O. Azencot, W. Yin, and A. Bertozzi, "Consistent dynamic mode decomposition," *arXiv preprint arXiv:1905.09736*, 2019.
- [24] K. Manohar, E. Kaiser, S. L. Brunton, and J. N. Kutz, "Optimized sampling for multiscale dynamics," *Multiscale Modeling & Simulation*, vol. 17, no. 1, pp. 117–136, 2019.
- [25] P. J. Schmid, "Dynamic mode decomposition of numerical and experimental data," *Journal of fluid mechanics*, vol. 656, pp. 5–28, 2010.
- [26] S. Sinha and E. Yeung, "On computation of koopman operator from sparse data," *arXiv:1901.03024*, 2019.
- [27] E. Yeung, S. Kundu, and N. Hodas, "Learning deep neural network representations for koopman operators of nonlinear dynamical systems," in *2019 American Control Conference (ACC)*. IEEE, 2019, pp. 4832–4839.
- [28] A. Hasnain, S. Sinha, Y. Dorfan, A. E. Borujeni, Y. Park, P. Maschhoff, U. Saxena, J. Urrutia, N. Gaffney, D. Becker *et al.*, "A data-driven method for quantifying the impact of a genetic circuit on its host," *arXiv preprint arXiv:1909.06455*, 2019.
- [29] C. A. Johnson and E. Yeung, "A class of logistic functions for approximating state-inclusive koopman operators," in *2018 Annual American Control Conference (ACC)*. IEEE, 2018, pp. 4803–4810.
- [30] S. E. Otto and C. W. Rowley, "Linearly recurrent autoencoder networks for learning dynamics," *SIAM Journal on Applied Dynamical Systems*, vol. 18, no. 1, pp. 558–593, 2019.
- [31] N. Takeishi, Y. Kawahara, and T. Yairi, "Learning koopman invariant subspaces for dynamic mode decomposition," in *Advances in Neural Information Processing Systems*, 2017, pp. 1130–1140.
- [32] Q. Li, F. Dietrich, E. M. Bollt, and I. G. Kevrekidis, "Extended dynamic mode decomposition with dictionary learning: A data-driven adaptive spectral decomposition of the koopman operator," *Chaos: An Interdisciplinary Journal of Nonlinear Science*, vol. 27, no. 10, p. 103111, 2017.
- [33] P. You, J. Pang, and E. Yeung, "Deep koopman controller synthesis for cyber-resilient market-based frequency regulation," *IFAC-PapersOnLine*, vol. 51, no. 28, pp. 720–725, 2018.
- [34] E. Yeung, Z. Liu, and N. O. Hodas, "A koopman operator approach for computing and balancing gramians for discrete time nonlinear systems," in *2018 Annual American Control Conference (ACC)*. IEEE, 2018, pp. 337–344.
- [35] J. Monod, "The growth of bacterial cultures," *Annual review of microbiology*, vol. 3, no. 1, pp. 371–394, 1949.
- [36] B. W. Brandt, I. M. van Leeuwen, and S. A. Kooijman, "A general model for multiple substrate biodegradation. application to co-metabolism of structurally non-analogous compounds," *Water research*, vol. 37, no. 20, pp. 4843–4854, 2003.
- [37] D. S. Kompala, D. Ramkrishna, N. B. Jansen, and G. T. Tsao, "Investigation of bacterial growth on mixed substrates: experimental evaluation of cybernetic models," *Biotechnology and Bioengineering*, vol. 28, no. 7, pp. 1044–1055, 1986.
- [38] H. Arbabi and I. Mezić, "Ergodic theory, dynamic mode decomposition, and computation of spectral properties of the koopman operator," *SIAM Journal on Applied Dynamical Systems*, vol. 16, no. 4, pp. 2096–2126, 2017.