Stable Invariant Models via Koopman Spectra

Takuya Konishi^{a,b,*}, Yoshinobu Kawahara^{a,b}

 ^aGraduate School of Information Science and Technology, Osaka University, 1-5 Yamadaoka, Suita, Osaka, Japan
 ^bCenter for Advanced Intelligence Project, RIKEN, 1-4-1 Nihonbashi, Chuo-ku, Tokyo, Japan

Abstract

Weight-tied models have attracted attention in the modern development of neural networks. The deep equilibrium model (DEQ) represents infinitely deep neural networks with weight-tying, and recent studies have shown the potential of this type of approach. DEQs are needed to iteratively solve root-finding problems in training and are built on the assumption that the underlying dynamics determined by the models converge to a fixed *point*. In this paper, we present the stable invariant model (SIM), a new class of deep models that in principle approximates DEQs under stability and extends the dynamics to more general ones converging to an invariant set (not restricted in a fixed point). The key ingredient in deriving SIMs is a representation of the dynamics with the spectra of the Koopman and Perron–Frobenius operators. This perspective approximately reveals stable dynamics with DEQs and then derives two variants of SIMs. We also propose an implementation of SIMs that can be learned in the same way as feedforward models. We illustrate the empirical performance of SIMs with experiments and demonstrate that SIMs achieve comparative or superior performance against DEQs in several learning tasks.

Keywords: Neural networks, Deep learning, Dynamical systems, Spectral analysis

© 2023. This manuscript version is made available under the CC-BY-NC-ND 4.0 license http://creativecommons.org/licenses/by-nc-nd/4.0/

^{*}Corresponding author

Email address: konishi@ist.osaka-u.ac.jp (Takuya Konishi)

1. Introduction

A feedforward neural network learns a representation by explicitly iterating a number of layer-by-layer computations. Each layer performs a transformation of outputs from the previous layer, which is typically characterized by different sets of parameters among the layers. However, several recent studies have shown that models with weight-tying, i.e. the ones employing the same transformation in each layer, achieve results competitive with stateof-the-art performances (Bai et al., 2019b; Dabre & Fujita, 2019; Dehghani et al., 2019). Motivated from this fact, Bai et al. (2019a) recently proposed the deep equilibrium model (DEQ), which is equal to running an *infinitely* deep feedforward model with weight-tying instead of using a finite number of layers. The models compute a representation by finding a fixed point (or an equilibrium point) with root-finding in practice, and are thus regarded as an instance of the so-called implicit-depth models such as neural ordinary differential equations (Chen et al., 2018). The following studies on DEQs have shown capability in this type of approach in several learning tasks (Bai et al., 2020; Winston & Kolter, 2020).

The forward pass of DEQs and their variants involves solving root-finding problems, which can lead to a high computational time and tends to be unstable (regarding, for example, the sensitivity in hyper-parameter tuning and initialization). Hence, DEQs sometimes require extensive and timeconsuming tuning to achieve strong performance and convergence to solutions (Linsley et al., 2020). Additionally, DEQs assume that the underlying dynamics determined by the models converge to a fixed *point*. However, as is often reported in papers on sequential neural models such as recurrent neural networks and also known in the scientific studies of brain activity, a broader class of convergence such as nonlinear oscillations could convey preferred capabilities in learning (Selverston & Moulins, 1985; Townley et al., 2000; Chang et al., 2019; Kag et al., 2020).

In this study, we propose a novel class of deep models, referred to as the *stable invariant model (SIM)*. The key insight behind SIMs is to interpret the underlying dynamics regarding DEQs through the Koopman operator. The Koopman operator is a linear operator over functions defined on latent states of dynamics (Koopman, 1931; Mezić, 2005). Because of its linearity, we can capture inherent temporal and spatial patterns of dynamics by the representation of the spectra, i.e. eigenvalues and eigenvectors, of the corresponding Koopman operator. We first show that the spectra of the Koopman



Figure 1: Comparison of DEQs and single-tier SIMs. The notation is introduced in Sections 2 and 3.

operator clarify stable dynamics with DEQs and then develop two variants of SIMs. The first models, single-tier SIMs, principally approximate DEQs under the stability. The resulting models, somewhat surprisingly, consist of only three-step transformations, although they approximate DEQs that are infinitely deep (see Figure 1). Moreover, the second models, two-tier SIMs, extend the dynamics to broader ones converging to an invariant *set* (e.g. a set of points, curve, and more general manifold) by employing the connection between the Koopman and Perron–Frobenius operators. We further provide a practical scheme to implement SIMs so that they can be learned in the same manner as feedforward models. Finally, we illustrate the behaviors of SIMs with numerical experiments in supervised learning tasks. We demonstrate that our models achieve competitive or superior performances compared to DEQs with less computational time.

The remainder of this paper is organized as follows. First, in Section 2, we briefly review DEQs, and the Koopman and Perron-Frobenius operators. In Section 3, we propose SIMs along with the description of their resulting architectures and characteristics. We describe the related works in Section 4 and investigate the empirical performance of our models in three learning tasks in Section 5. We conclude this paper in Section 6. The details of some equation derivations and experiments are presented in Appendix A and Appendix B, respectively.

2. Background

2.1. Deep Equilibrium Models

One of the core ideas in DEQs is weight-tying, i.e. the same set of parameters is shared across the layers of a deep network. Formally, DEQs consider an *L*-layer weight-tied transformation with shared parameters θ :

$$\boldsymbol{z}_{l+1} = \boldsymbol{f}_{\theta}(\boldsymbol{z}_l, \boldsymbol{x}), \quad l = 0, 1, \dots, L-1,$$
(1)

where $\boldsymbol{x} \in \mathbb{R}^{D}$ is the input to the model, $\boldsymbol{z}_{l} \in \mathbb{R}^{d}$ is the hidden state of the *l*-th layer, and $\boldsymbol{f}_{\theta} \colon \mathbb{R}^{d+D} \to \mathbb{R}^{d}$ is a continuous function. DEQs suppose that stacking such layers infinitely tends to a fixed point:

$$\lim_{l\to\infty} \boldsymbol{z}_l = \lim_{l\to\infty} \boldsymbol{f}_{\theta}^l(\boldsymbol{z}_0, \boldsymbol{x}) \coloneqq \boldsymbol{f}_{\theta}(\boldsymbol{z}^*, \boldsymbol{x}) = \boldsymbol{z}^*.$$
(2)

The forward pass of DEQs uses root-finding algorithms to directly compute the fixed point z^* by solving the equation $z^* = f_{\theta}(z^*, x)$. Then, z^* is transformed to an output y by a function h as $y = h(z^*)$. We can train DEQs with backpropagation by computing the gradient of the fixed point through implicit differentiation (Krantz & Parks, 2013).

A series of transformations in Eq. (1) can be viewed as a discrete-time nonlinear dynamical system where the hidden state z_l is a state vector at step l. The underlying dynamics are determined by the transformation f_{θ} , which is affected by the input x at every step.

2.2. Koopman and Perron–Frobenius Operators

We briefly overview the Koopman and Perron–Frobenius operators. Please see other references, e.g. (Mauroy et al., 2020), for more information.

Definitions. Consider a discrete-time nonlinear dynamical system: $\mathbf{z}_{t+1} = \mathbf{f}(\mathbf{z}_t)$, defined on a state space $\mathbb{S} \subset \mathbb{R}^d$, where $\mathbf{z}_t \in \mathbb{S}$ is the state vector at time t, and $\mathbf{f} \colon \mathbb{S} \to \mathbb{S}$ is a (possibly, nonlinear) state-transition function. Let $g \in \mathcal{G}$ be an observable, which is a scalar complex-valued function on \mathbb{S} in some (Banach) space \mathcal{G} . The Koopman operator $\mathcal{K} \colon \mathcal{G} \to \mathcal{G}$ is defined through the following composition:

$$(\mathcal{K}g)(\boldsymbol{z}) = g(\boldsymbol{f}(\boldsymbol{z})),$$

where $z \in S$ is a state vector. \mathcal{K} acts on observables and maps g to a new function $\mathcal{K}g$. Although the dynamics described by f may be nonlinear, \mathcal{K} is linear and infinite-dimensional.

The Perron–Frobenius operator is often used to describe the transition of the density over the state of dynamical systems (Lasota & Mackey, 1994; Gaspard, 1998; Cvitanović et al., 2020). Given a measure space $(\mathbb{S}, \mathbb{A}, \mu)$ that consists of a state space \mathbb{S} , σ -algebra \mathbb{A} , and measure μ , we suppose that \boldsymbol{f} is nonsingular if $\mu(\boldsymbol{f}^{-1}(A)) = 0$ for a Borel set $A \in \mathbb{A}$ such that $\mu(A) = 0$, where $\boldsymbol{f}^{-1}(A)$ is the preimage of \boldsymbol{f} given A. Let $p \in \mathcal{L}^1$ denote a density function on \mathbb{S} in the space of absolutely integrable functions \mathcal{L}^1 . The *Perron–Frobenius operator* $\mathcal{P} \colon \mathcal{L}^1 \to \mathcal{L}^1$ acts on densities and is defined as

$$\int_{A} (\mathcal{P}p)(\boldsymbol{z}) d\mu = \int_{\boldsymbol{f}^{-1}(A)} p(\boldsymbol{z}) d\mu, \quad A \in \mathbb{A}.$$

It should be noted that the Koopman (or Perron–Frobenius) operator is the adjoint of the Perron–Frobenius (or Koopman) operator for appropriately defined spaces.

Koopman Spectrum. Because the Koopman operator \mathcal{K} (and also Perron– Frobenius operator) is linear, it can be characterized by spectral properties. We assume \mathcal{K} has only point spectra and also has non-trivial eigenfunctions. Let $\lambda_j \in \mathbb{C}, j = 1, 2, ...,$ be the eigenvalue of \mathcal{K} . The eigenfunction $\phi_j \colon \mathbb{S} \to \mathbb{C}$ for λ_j satisfies the relation

$$(\mathcal{K}\phi_j)(\boldsymbol{z}) = \phi_j(\boldsymbol{f}(\boldsymbol{z})) = \lambda_j \phi_j(\boldsymbol{z}).$$

 λ_j and ϕ_j are called the Koopman eigenvalue and Koopman eigenfunction, respectively. If an observable g is in the subspace of \mathcal{G} spanned by the Koopman eigenfunctions $\{\phi_j\} := \{\phi_j \mid j = 1, 2, \ldots\}$, the observable can be represented as

$$(\mathcal{K}g)(\boldsymbol{z}) = g(\boldsymbol{f}(\boldsymbol{z})) = \sum_{j=1}^{\infty} \lambda_j v_j \phi_j(\boldsymbol{z}),$$

where the coefficient $v_j \in \mathbb{C}$, j = 1, 2, ..., is referred to as the Koopman mode associated with g. The subspace spanned by $\{\phi_j\}$ is invariant under the Koopman operator, i.e., the observables in the subspace remain in the subspace after being acted by \mathcal{K} .

Finite-Dimensional Approximation. Consider a subspace of observables which is spanned by N basis functions $\{\varphi_i\} \coloneqq \{\varphi_i \colon \mathbb{S} \to \mathbb{C} \mid j = 1, 2, ..., N\}$. If g exists on the subspace, then g is represented as a linear combination of the basis functions, i.e. $g(\boldsymbol{z}) = \boldsymbol{w}^{\top} \boldsymbol{\varphi}(\boldsymbol{z})$, where we denote the concatenation of the basis functions as a vector-valued one $\boldsymbol{\varphi} = (\varphi_1, \ldots, \varphi_N)^{\top} \colon \mathbb{S} \to \mathbb{C}^N$, and $\boldsymbol{w} \in \mathbb{C}^N$ is the coordinate of g on the subspace. By projecting the action of \mathcal{K} onto the span of $\{\varphi_j\}$, we approximate \mathcal{K} with another linear operator \mathcal{K}_N . This approximates the Koopman operator \mathcal{K} and satisfies

$$(\mathcal{K}_N g)(\boldsymbol{z}) = (\boldsymbol{K} \boldsymbol{w})^\top \boldsymbol{\varphi}(\boldsymbol{z}),$$

where $\mathbf{K} \in \mathbb{C}^{N \times N}$ is referred to as the Koopman matrix. The above shows that a temporal evolution of g with \mathcal{K}_N is represented by applying \mathbf{K} to the coordinate \mathbf{w} . Therefore, the Koopman matrix owns the one-to-one correspondence to \mathcal{K}_N . Moreover, the Koopman eigenvalues, eigenfunctions, and modes of \mathcal{K}_N can be obtained from the eigenvalues, right-eigenvectors, and left-eigenvectors of \mathbf{K} , respectively.

Additionally, the Koopman matrix provides an approximation of the dynamics through the basis functions:

$$\varphi(\boldsymbol{f}(\boldsymbol{z})) \approx (\boldsymbol{K}^{\top} \varphi)(\boldsymbol{z}) \coloneqq (\boldsymbol{A} \varphi)(\boldsymbol{z}),$$
 (3)

where \boldsymbol{A} is the transpose of the Koopman matrix, i.e. $\boldsymbol{A} = \boldsymbol{K}^{\top}$. Eq. (3) indicates that a temporal evolution with \boldsymbol{f} can be approximated by a temporal evolution with the finite and linear dynamics described by \boldsymbol{A} over a *lifted* space with $\boldsymbol{\varphi}$. The equation holds if $\mathcal{K} = \mathcal{K}_N$.

It should be noted that, if g is a real-valued function, then it is sufficient to consider real-valued ones for the corresponding quantities (such as φ_j , \boldsymbol{w} and \boldsymbol{K}).

3. Stable Invariant Models

In this section, we introduce our proposed SIMs. The models are motivated by two problems concerning DEQs. First, DEQs work under the assumption that the underlying dynamics are stable, i.e. Eq. (2) holds for any input and initial state. However, it is in general hard to assess whether nonlinear dynamics satisfy the assumption. We address this problem by approximately specifying the stable dynamics using the Koopman spectrum, which leads our first single-tier SIMs by identifying the convergent behavior for a fixed point. The second problem is that DEQs only consider convergence to a fixed point. In the literature on dynamical systems, one often considers more general convergence to an invariant set: for a dynamical system on a state space S, a set $S \subset S$ is said to be an invariant set if any trajectory starting in S remain in S. The convergence to an invariant set allows for the dynamics that oscillate on the set. Notable examples include limit cycles, tori, and other nonlinear oscillations (Strogatz, 2015). We also refer to dynamics as stable if the dynamics converge to a bounded invariant set for any input and initial state and construct our second two-tier SIMs by incorporating only the dynamics converging to an invariant set via the Koopman and Perron– Frobenius operators.

3.1. Approximating DEQs

We begin by considering the representations of DEQs with the Koopman operator. However, it should be noted that f_{θ} in Eq. (1) has another vector \boldsymbol{x} as an input differently from \boldsymbol{f} in Section 2.2. Although there are several ways to define the Koopman operator for such a case, we consider the following Koopman operator \mathcal{K} acting on the observable g of both the hidden state \boldsymbol{z} and input \boldsymbol{x} (Proctor et al., 2018):

$$(\mathcal{K}g)(\boldsymbol{z},\boldsymbol{x}) = g(\boldsymbol{f}_{\theta}(\boldsymbol{z},\boldsymbol{x}),\boldsymbol{x}).$$
(4)

The action of this operator is restricted so that \boldsymbol{x} is maintained at the same point. Even if \boldsymbol{x} is injected, the properties of the Koopman operator discussed in Section 2.2 still hold.

Hereafter, we consider a real-valued observable g. First, let $\varphi_j(\boldsymbol{z}, \boldsymbol{x}) \colon \mathbb{R}^{d+D} \to \mathbb{R}$ be N real-valued basis functions $(j = 1, \ldots, N)$. The concatenation is given by a vector-valued basis function $\boldsymbol{\varphi}$. If Eq. (1) converges to a fixed point \boldsymbol{z}^* as in Eq. (2), we can approximate \boldsymbol{z}^* with \boldsymbol{x} by the Koopman operator of Eq. (4) as

$$(\boldsymbol{z}^{*}, \boldsymbol{x}) = \left(\lim_{l \to \infty} \boldsymbol{f}_{\theta}(\boldsymbol{z}_{l}, \boldsymbol{x}), \boldsymbol{x}\right)$$
$$= \left(\lim_{l \to \infty} \boldsymbol{f}_{\theta}^{l}(\boldsymbol{z}_{0}, \boldsymbol{x}), \boldsymbol{x}\right)$$
$$\approx \widehat{\boldsymbol{\varphi}^{-1}} \left(\boldsymbol{\varphi} \left(\lim_{l \to \infty} \boldsymbol{f}_{\theta}^{l}(\boldsymbol{z}_{0}, \boldsymbol{x}), \boldsymbol{x}\right)\right)$$
$$\approx \widehat{\boldsymbol{\varphi}^{-1}} \left(\lim_{l \to \infty} \boldsymbol{A}^{l} \boldsymbol{\varphi}(\boldsymbol{z}_{0}, \boldsymbol{x})\right).$$
(5)

Here, we define a function $\widehat{\varphi^{-1}}$ in the first approximation. The equation holds if there exists the inverse function φ^{-1} and $\widehat{\varphi^{-1}} = \varphi^{-1}$. Strictly invertible φ with respect to all possible inputs and outputs is rather restrictive and expensive in practice. We relax the invertibility by supposing $\widehat{\varphi^{-1}}$ as a surrogate function that approximately models the output-input relations in the subspace where most of the inputs and outputs are distributed. The second approximation of Eq. (5) follows the finite-dimensional approximation of the Koopman operator through the lifted dynamics as described in Eq. (3). The basis functions need higher expressiveness to better approximate the subspace of the observables that the Koopman operator acts. By dividing the surrogate function $\widehat{\varphi^{-1}}$ into two parts corresponding to z and x, i.e. $\widehat{\varphi^{-1}} = (\widehat{\varphi_z^{-1}}, \widehat{\varphi_x^{-1}})$, the fixed point z^* can be approximated more directly as

$$\boldsymbol{z}^* \approx \widehat{\boldsymbol{\varphi}_{\boldsymbol{z}}^{-1}} \left(\lim_{l \to \infty} \boldsymbol{A}^l \boldsymbol{\varphi}(\boldsymbol{z}_0, \boldsymbol{x}) \right).$$
 (6)

3.2. Convergent Behavior via Koopman Spectra

The approximation (6) assumes that DEQs converge to a fixed point. However, whether a DEQ converges to a fixed point depends on the behavior of the underlying dynamics with the DEQ. The approximation (6) allows us to characterize the convergent behavior of a DEQ via eigenvalues of the corresponding A.

First, we denote by $\lambda_j \in \mathbb{C}$ for j = 1, 2, ..., N the eigenvalues of \boldsymbol{A} . It should be noted that the eigenvalues can be complex values even though \boldsymbol{A} is a real matrix (because it is not necessarily symmetric). For any real matrix \boldsymbol{A} , there exists a nonsingular matrix $\boldsymbol{U} \in \mathbb{R}^{N \times N}$ that consists of the generalized-eigenvectors including the ordinal eigenvectors of \boldsymbol{A} . We represent \boldsymbol{U} and \boldsymbol{U}^{-1} with N vectors, respectively, as $\boldsymbol{U} = (\boldsymbol{u}_1, \ldots, \boldsymbol{u}_N)$ and $\boldsymbol{U}^{-1} = (\boldsymbol{v}_1, \ldots, \boldsymbol{v}_N)^{\top}$, where \boldsymbol{u}_j is associated with a generalized-eigenvector of λ_j if λ_j is real, and the real or imaginary part of a generalized-eigenvector of λ_j if λ_j is non-real. Additionally, let $\rho(\boldsymbol{A})$ be the spectral radius of \boldsymbol{A} , i.e. $\rho(\boldsymbol{A}) \coloneqq \max\{|\lambda_1|, \ldots, |\lambda_N|\}$. We can then classify the convergent behavior of the lifted dynamics into the following four cases:

(i). If $\rho(\mathbf{A}) < 1$, then the dynamics converge to the origin:

$$\lim_{l o\infty} oldsymbol{A}^l oldsymbol{arphi}(oldsymbol{z}_0,oldsymbol{x}) = oldsymbol{0}$$

(ii). If $\rho(\mathbf{A}) = 1$, all the eigenvalues with $|\lambda_j| = 1$ take the values 1, and their corresponding eigenvectors are linearly independent, then the lifted dynamics converge to a fixed point. That is, if we denote by $J_1 = \{j \mid \lambda_j = 1\}$ the index set of such eigenvalues, then we have

$$\lim_{l\to\infty} \boldsymbol{A}^{l} \boldsymbol{\varphi}(\boldsymbol{z}_{0}, \boldsymbol{x}) = \sum_{j\in J_{1}} \boldsymbol{u}_{j} \boldsymbol{v}_{j}^{\top} \boldsymbol{\varphi}(\boldsymbol{z}_{0}, \boldsymbol{x}).$$
(7)

(iii). If $\rho(\mathbf{A}) = 1$ and the eigenvectors with eigenvalues of $|\lambda_j| = 1$ are linearly independent, then the lifted dynamics do not converge to a point but oscillates in the state space. More concretely, if we denote $\lambda_j = \alpha_j + i\beta_j, J_2 = \{j \mid \lambda_j = -1\}, \text{ and } J_3 = \{(j,k) \mid |\lambda_j| = |\lambda_k| = 1, \beta_j = \beta_k \neq 0, \lambda_k = \overline{\lambda_j}\},$ then we have

$$\lim_{l \to \infty} \boldsymbol{A}^{l} \boldsymbol{\varphi}(\boldsymbol{z}_{0}, \boldsymbol{x}) = \lim_{l \to \infty} \left(\sum_{j \in J_{1}} \boldsymbol{u}_{j} \boldsymbol{v}_{j}^{\top} + \sum_{j \in J_{2}} (-1)^{l} \boldsymbol{u}_{j} \boldsymbol{v}_{j}^{\top} + \sum_{(j,k) \in J_{3}} \left(\cos(l\Delta_{j}) \boldsymbol{u}_{j} - \sin(l\Delta_{k}) \boldsymbol{u}_{k} \right) \boldsymbol{v}_{j}^{\top} + \left(\sin(l\Delta_{j}) \boldsymbol{u}_{j} + \cos(l\Delta_{k}) \boldsymbol{u}_{k} \right) \boldsymbol{v}_{k}^{\top} \right) \boldsymbol{\varphi}(\boldsymbol{z}_{0}, \boldsymbol{x}),$$

$$(8)$$

where *i* is the imaginary unit, and $\Delta_j = \arctan(\beta_j/\alpha_j)$. J_3 is the set of index pairs whose eigenvalues are conjugated. It should be noted that there always exists a conjugate eigenvalue for every non-real eigenvalue when \boldsymbol{A} is real.

(iv). Otherwise, at least one element among the states of the lifted dynamics diverges.

We describe the derivations of the four cases in Appendix A. Figure 2 shows the complex plane and coordinates of the eigenvalues described in the above cases. The red and green dots, and the blue line denote the locations of the eigenvalues that correspond to J_1 , J_2 , and J_3 , respectively. The gray area shows the area where the absolute values of the eigenvalues are less than 1. The first three cases approximately correspond to the stable dynamics with the DEQ. Case (i) converges to the origin regardless of \boldsymbol{x} , which is useless for any learning problems. In contrast, case (ii) converges to a fixed point that reflects \boldsymbol{x} . Case (iii) does not converge to any fixed point but oscillates on a manifold in the state space, where the terms corresponding to



Figure 2: Areas where eigenvalues corresponding to J_1 , J_2 , and J_3 are located in the complex plane.

 J_2 and J_3 respectively include coefficients such as $(-1)^l$ and $\cos(l\Delta_j)$ that neither diverge nor converge. Lastly, case (iv) diverges and the corresponding dynamics are not stable.

3.3. Model Description

We derive SIMs based on the above analysis of DEQs from the perspective of the Koopman spectrum. We first describe the variant that approximates DEQs by leveraging case (ii). Moreover, we present another variant of SIMs that incorporates a broader class of dynamics by utilizing case (iii).

Fixed Point. Case (ii) represents the dynamics that converge to a fixed point as in Eq. (7). Because DEQs assume that the underlying dynamics converge towards a fixed point, which are covered by case (ii), we can approximate the fixed point of a DEQ by plugging Eq. (7) into Eq. (6):

$$\boldsymbol{z}^* \approx \widehat{\boldsymbol{\varphi}_{\boldsymbol{z}}^{-1}} \left(\sum_{j \in J_1} \boldsymbol{u}_j \boldsymbol{v}_j^\top \boldsymbol{\varphi}(\boldsymbol{z}_0, \boldsymbol{x}) \right) = \widehat{\boldsymbol{\varphi}_{\boldsymbol{z}}^{-1}} \left(\hat{\boldsymbol{U}} \hat{\boldsymbol{V}} \boldsymbol{\varphi}(\boldsymbol{z}_0, \boldsymbol{x}) \right), \quad (9)$$

where $\hat{U} \in \mathbb{R}^{N \times K}$ and $\hat{V} \in \mathbb{R}^{K \times N}$ consist of the vectors that are the rows of U and V whose indices belong to J_1 and K is the size of J_1 . Interestingly, this approximation implies that a fixed point of DEQs can be represented as a finite-depth model that consists of three-step transformations if the basis functions can approximate the subspace of the observables that the corresponding Koopman operator acts. We refer to the right-hand side of Eq. (9) single-tier SIMs. This comes from the fact that a state of the dynamics is lifted to another tier with φ .



Figure 3: Overview of SIMs.

Invariant Set. Case (iii) represents the dynamics that converge to an invariant set given by Eq. (8). While such dynamics do not converge to any single fixed point, it is necessary to characterize case (iii) with some representative point to obtain a trainable model. We here focus on invariant sets being characterized by spectra with eigenvalue 1 of the corresponding Perron–Frobenius operator (Billings & Schwartz, 2008; Froyland & Padberg, 2009). This viewpoint suggests that the spectra can encode the information of convergent trajectories in case (iii).

Now, let $\boldsymbol{\psi}$ be the (finite-dimensional) basis functions of an embedding of the state into some Hilbert space that encodes the density defining the corresponding Perron–Frobenius operator. Building upon the above insight, we propose to represent case (iii) as the following extended equilibrium over the states on one more lifted space similar to case (ii):

$$\boldsymbol{z}^{\star} \coloneqq \widehat{\boldsymbol{\varphi}_{\boldsymbol{z}}^{-1}} \left(\widehat{\boldsymbol{\psi}^{-1}} \left(\sum_{j \in J_{1}'} \boldsymbol{u}_{j}' \boldsymbol{v}_{j}'^{\top} \boldsymbol{\psi} \left(\boldsymbol{\varphi}(\boldsymbol{z}_{0}, \boldsymbol{x}) \right) \right) \right)$$

$$= \widehat{\boldsymbol{\varphi}_{\boldsymbol{z}}^{-1}} \left(\widehat{\boldsymbol{\psi}^{-1}} \left(\widehat{\boldsymbol{U}}' \widehat{\boldsymbol{V}}' \boldsymbol{\psi} \left(\boldsymbol{\varphi}(\boldsymbol{z}_{0}, \boldsymbol{x}) \right) \right) \right),$$
(10)

where J'_1 is the index set of the spectra of the Perron–Frobenius operator analogous to J_1 , and \hat{U}' and \hat{V}' are the corresponding matrices. $\widehat{\psi}^{-1}$ is another surrogate function for approximating the inverse function of ψ . The new point z^* reflects the convergent behavior of the dynamics over an invariant set. We refer to this type of models as *two-tier* SIMs, which come from the structure involving two-tier maps by φ and ψ . Two-tier SIMs are also represented as finite-depth models.

Figure 3 shows the overview of SIMs. The black arrows denote the first transformation with the basis functions φ . The blue and red arrows correspond to the transformations of the single-tier and two-tier SIMs, respectively. Each model considers the convergent behavior of the dynamics in the first-tier or second-tier space and then returns to the original space with the surrogate functions.

Although DEQs are defined by the transformation f_{θ} , SIMs require the specification of the basis functions φ and ψ . One natural question may be how the two approaches and convergence properties relate to each other. Generally, if two dynamical systems converge to similar fixed points or invariant sets, the corresponding dynamics also have similar topological properties. However, the actual transformations are determined by many factors (e.g. the training algorithms and choice of basis functions), and the convergence regions of the two approaches will be different. Investigating how f_{θ} can be associated with φ and ψ is an important future work.

3.4. Practical Implementation

The key to encode high expressiveness in the dynamics with SIMs is to utilize rich basis functions. Therefore, we leverage neural networks with structures to construct the components in our models (9) and (10) as follows.

First, we utilize a neural network whose inputs are only \boldsymbol{x} , which we denote by $\boldsymbol{\mu}_{\text{NN}}(\boldsymbol{x})$, to approximate $\boldsymbol{\varphi}(\boldsymbol{z}_0, \boldsymbol{x})$. This is because the initial state \boldsymbol{z}_0 can be basically assumed to be zero. Although we can prepare a non-zero \boldsymbol{z}_0 (e.g. the original implementation of DEQs allows to use the final state of the previous step in the training loop), the effectiveness of such a biased initial state is unclear. Therefore, we assume that \boldsymbol{z}_0 is always zero and encoded in the basis functions, and thus drop \boldsymbol{z}_0 from the input of $\boldsymbol{\varphi}$.

Next, because \hat{U} and \hat{V} in Eq. (9) are simply real matrices, we can deal with them as two consecutive linear layers although \hat{U} and \hat{V} respectively consist of linearly independent vectors. Although we could consider the constraint, we ignore it because the vectors in a learned matrix will rarely be linearly dependent. \hat{V} has another implicit constraint: \hat{V} originates from U^{-1} and thus is affected by \hat{U} . However, \hat{V} is also influenced by the vectors of U which correspond to eigenvalues with an absolute value less than 1. Although these vectors disappear by taking the limit and do not appear in the model, \hat{V} will have degrees of freedom thanks to the vectors. From this observation, we model \hat{V} independently from \hat{U} .

Further, we prepare another neural network for the surrogate functions, $\widehat{\varphi_z^{-1}}$ and $\widehat{\psi^{-1}}$. This network is, principally, required to approximate the inverse function corresponding to the respective part of φ and ψ . However, even approximately modeling a neural network to be invertible is costly and difficult unless the dimension of the lifted space is identical to that of the original space (Papamakarios et al., 2021). The invertibility will also require SIMs to consider $\widehat{\varphi_x^{-1}}$ that is omitted from Eq. (6). This is because it is necessary to construct $\widehat{\varphi_{-1}^{-1}}$ including $\widehat{\varphi_x^{-1}}$ to approximate the invertibility even though the output of $\widehat{\varphi_x^{-1}}$ is never used in subsequent transformations. In this paper, we focus on the practical aspect of our implementation and approximate $\widehat{\varphi_z^{-1}}$ in Eq. (9) or the composition of $\widehat{\varphi_z^{-1}}$ and $\widehat{\psi^{-1}}$ in Eq. (10) by a neural network $\nu_{\rm NN}$ without the restriction of invertibility. This approach follows the same manner as common encoder-decoder models; an encoder model is usually designed independently of the corresponding decoder model despite their close connection of going back and forth between original and latent spaces.

For two-tier SIMs, we employ random Fourier features (RFFs) (Rahimi & Recht, 2007) to approximate the embedding ψ , i.e.

$$\boldsymbol{\psi}(\boldsymbol{x}) = \frac{1}{\sqrt{M/2}} \left(\sin(\boldsymbol{\omega}_1^{\top} \boldsymbol{x}), \ \cos(\boldsymbol{\omega}_1^{\top} \boldsymbol{x}), \ \cdots, \ \sin(\boldsymbol{\omega}_{M/2}^{\top} \boldsymbol{x}), \ \cos(\boldsymbol{\omega}_{M/2}^{\top} \boldsymbol{x}) \right)^{\top},$$
$$\boldsymbol{\omega}_j \sim P(\boldsymbol{\omega}),$$

where $\boldsymbol{\omega}_j \in \mathbb{R}^N$ for $j = 1, \ldots, M/2$ are random vectors drawn from a probability distribution $P(\boldsymbol{\omega})$, to avoid the increase of the computational cost along the sample size when using reproducing kernels. Because $\boldsymbol{\psi}$ is the M-dimensional real-valued function, $\hat{\boldsymbol{U}}'$ and $\hat{\boldsymbol{V}}'$ are defined as $M \times K'$ and $K' \times M$ real matrices, respectively, where K' is the size of J'_1 . We model those matrices in the same way as $\hat{\boldsymbol{U}}$ and $\hat{\boldsymbol{V}}$.

Putting the above pieces together, we model Eqs. (9) and (10) by $\boldsymbol{z}_{\text{single-tier}}$ and $\boldsymbol{z}_{\text{two-tier}}$, respectively, as follows:

$$\begin{aligned} \boldsymbol{z}_{\text{single-tier}} &= \boldsymbol{\nu}_{\text{NN}} \left(\hat{\boldsymbol{U}} \hat{\boldsymbol{V}} \boldsymbol{\mu}_{\text{NN}}(\boldsymbol{x}) \right) \quad \text{and} \\ \boldsymbol{z}_{\text{two-tier}} &= \boldsymbol{\nu}_{\text{NN}} \left(\hat{\boldsymbol{U}}' \hat{\boldsymbol{V}}' \boldsymbol{\psi} \left(\boldsymbol{\mu}_{\text{NN}}(\boldsymbol{x}) \right) \right). \end{aligned}$$
(11)

In the end, those implementations of SIMs are realized as feedforward models. The implementations no longer require any root-finding algorithms for the forward pass and implicit differentiation for the backward pass. In the following experiments, we will instantiate more specific examples of the implementations for each task. However, the implementations do not limit themselves to such particular forms and have the flexibility to take various configurations by changing $\mu_{\rm NN}$ and $\nu_{\rm NN}$.

A practical merit of DEQs compared to feedforward models is the memory efficiency: the required memory does not depend on the number of transformations by f_{θ} owing to root-finding and implicit differentiation. Because we realize SIMs as feedforward models, they do not fully inherit this efficiency. However, DEQs also model f_{θ} with feedforward architectures such as Transformers (Vaswani et al., 2017; Bai et al., 2019a) and thus need the memory for f_{θ} , which could be potentially large. Additionally, we notice that DEQs are infinitely deep. The naive approximation with a feedforward model may need to increase the depth of the feedforward model. Our derivation showed that it is critical to approximate the DEQs according to the form of Eq. (11). If the basis functions of SIMs could be represented by shallower models, it will save a lot of memory compared to naive approximations. We can regard SIMs as an intermediate approach that can approximate infinitely deep models while avoiding the simple dependence on depth.

3.5. Relation to Implicit Neural Representation

Implicit neural representation has recently gained interest in the machine learning community. The aim is to represent complex natural signals (e.g. images, 3D objects, and audio waves) with a function modeled by a neural network (Sitzmann et al., 2020; Tancik et al., 2020).

SIMs have an interesting connection to implicit neural representation by (Tancik et al., 2020). Tancik et al. (2020) proposed to use a Fourier feature mapping as a pre-processing. If the basis functions modeled by $\mu_{\rm NN}$ of single-tier SIMs are modeled by an RFF ψ , then the implementation of the SIMs is the same as the Fourier feature mapping by (Tancik et al., 2020) up to the term $\hat{U}\hat{V}$:

$$\boldsymbol{z}_{\text{RFF}} = \boldsymbol{\nu}_{\text{NN}} \left(\hat{\boldsymbol{U}} \hat{\boldsymbol{V}} \boldsymbol{\psi}(\boldsymbol{x}) \right).$$
 (12)

As in Eq. (11), z_0 is omitted from $\psi(x)$ because it is supposed to be zero and does not affect the output of ψ . Alternatively, this RFF only model can be interpreted as an instance of two-tier SIMs when $\mu_{\rm NN}$ is an identity mapping. RFFs have been initially proposed as an approximation for kernel methods, which are also applied for basis functions in Koopman operator analysis (Kawahara, 2016). Hence, RFFs will be a natural choice as a class of basis functions for SIMs. Eq. (12) indicates that such a reasonable choice leads to a similar method to (Tancik et al., 2020).

4. Related Works

The origin of DEQs dates back to the work on recurrent backpropagation (RBP) (Almeida, 1987; Pineda, 1987). Those studies proposed to utilize early implicit-depth models and have been applied to other studies, e.g. graph neural networks (Gori et al., 2005; Scarselli et al., 2008). Recently, Liao et al. (2018) revisited the RBP algorithm and improved it with the conjugate gradient method and Neumann series. Bai et al. (2019a) introduced a perspective of the use of a fixed point as a replacement for depth and called the approaches the DEQ. This study also proved the universality of a single DEQ layer and developed a practical quasi-Newton method that works for large-scale sequential tasks. Subsequent studies have reported the theoretical analysis (Kawaguchi, 2021; Pabbaraju et al., 2021) and proposed improved architectures (Bai et al., 2020; Xie et al., 2021). Several relevant studies focused on the stability issue of DEQs. Winston & Kolter (2020) proposed monotone DEQs (monDEQs) with guaranteed convergence to a fixed point based on monotone operators. Bai et al. (2021) also proposed to stabilize the training using Jacobian regularization.

The Koopman operator has been known for a long time as a tool for analyzing dynamics in physics (Koopman, 1931). It has recently received attention that the spectral properties of the Koopman operator play an important role in revealing global characteristics of the underlying dynamics (Mezić, 2005). Because the Koopman operator is linear but infinite-dimensional, a major challenge is how to compute the spectra of the Koopman operator in practice. Dynamic mode decomposition (DMD) (Schmid, 2010) has gained popularity as a data-driven approach to computing a reasonable finitedimensional approximation of the Koopman spectra (Rowley et al., 2009). While the original DMD algorithm supposes that the basis functions are linear, several subsequent studies have been proposed to utilize nonlinear basis functions (Williams et al., 2015), a kernel-based approach (Kawahara, 2016), a Bayesian formulation (Takeishi et al., 2017a), and learning from data (Takeishi et al., 2017b). The Koopman operator is also applied to machine learning and the related fields. Dogra & Redman (2020) leveraged the representation with the Koopman operator to accelerate the training of neural networks. Manojlović et al. (2020) analyzed the dynamics in the training of neural networks with the Koopman operator and proposed a method to characterize the architectures of neural networks with the Koopman spectrum. Takeishi & Kawahara (2021) proposed a neural network to learn stable dynamical systems by utilizing a map obtained by, for example, an eigenfunction of the Koopman operator (although it is not described explicitly in the paper) by extending stable deep dynamics models by Manek & Kolter (2019).

5. Experiments

We evaluated SIMs through experiments on three supervised learning tasks. For all tasks, we first trained models with candidates of hyperparameters on training data and evaluated them on validation data. We then selected the best hyper-parameter, re-trained the model with the best one, and evaluated the trained model on test data. We implemented our proposed models and training algorithms using Pytorch (Paszke et al., 2019). For the basic building blocks and optimizers, we used the existing implementation in Pytorch. We tuned the hyper-parameters using Tune (Liaw et al., 2018). In all experiments, we set the search algorithm to the random search and trial scheduler to the asynchronous successive halving algorithm (ASHA) (Li et al., 2020). We conducted the experiments on an internal computing server with Intel Xeon Bronze 3104 CPUs and NVIDIA V100 GPUs.

We set $\mu_{\rm NN}$ as a task-specific neural network for each task. For ψ , we used a normal distribution to sample random vectors. For $\nu_{\rm NN}$, we used a three-layer fully connected network (FCN) with ReLU activation (Nair & Hinton, 2010) for all tasks in common. We used a linear function as h to fit the dimension of the hidden state to the output size and then the softmax function if the addressed task is a classification problem. For SIMs, we set K = N/2 and K' = M/2, respectively.

The details of the datasets, model architectures, and training algorithms are described in Appendix B.

5.1. Copy Memory Task

We first report the results of the copy memory task (Hochreiter & Schmidhuber, 1997) to evaluate the effectiveness of SIMs against DEQs. The goal

of the copy memory task is to predict a sequence of digits of length T + 20from another input sequence of the same length. The first ten elements of the input sequence consist of digits randomly drawn from 1 to 8, the subsequent T-1 elements are filled with 0, and the last eleven elements are all 9. Given this input, the first T+10 elements of the output sequence are 0 and the last ten elements are the same as the first ten ones of the input sequence. Hence, this task evaluates how well a model can remember the first elements of an input. In the experiment, we set T to 500 and followed the experimental procedure of (Bai et al., 2018).

For SIMs, we applied the temporal convolutional network (TCN) architecture to model $\mu_{\rm NN}$ by following the implementation of (Bai et al., 2018): the architecture includes 1D dilated causal convolution, ReLU activation, and residual connection (He et al., 2016). Because the method of applying the TCN is rather complicated compared to other tasks, we explain the case of the single-tier SIM as an example. Formally, if we denote $\boldsymbol{x} \in \{0, 1, \ldots, 9\}^{520}$ as an input sequence, the TCN outputs the sequence of states:

$$\boldsymbol{W} = \mathrm{TCN}(\boldsymbol{x}),$$

where $\boldsymbol{W} = (\boldsymbol{w}_1, \dots, \boldsymbol{w}_{520}) \in \mathbb{R}^{N \times 520}$ and $\boldsymbol{w}_j \in \mathbb{R}^N$ denotes a state on the lifted space at the *j*-th position in the sequence. Each state feeds the common architecture:

$$\boldsymbol{z}_{\text{single},j} = \boldsymbol{\nu}_{\text{NN}} \left(\hat{\boldsymbol{U}} \hat{\boldsymbol{V}} \boldsymbol{w}_j \right) \quad (j = 1, ..., 520),$$

where $\boldsymbol{z}_{\text{single},j} \in \mathbb{R}^d$ is the hidden state of the *j*-th position. The prediction at the *j*-th position of the output sequence is obtained from $\boldsymbol{z}_{\text{single},j}$. We can interpret this architecture as follows; for this task, the model has different $\boldsymbol{\mu}_{\text{NN}}$ for each position, but it transforms all the positions to the hidden states at once by the TCN, and then the same subsequent transformation is applied to all the positions. Although we could use different $\boldsymbol{\nu}_{\text{NN}}$ and $\hat{\boldsymbol{U}}\hat{\boldsymbol{V}}$ for each position, the above architecture can retain memory in practice while keeping the model size small.

For the DEQ, we applied the Universal Transformer (Dehghani et al., 2019) as the base function f_{θ} and mostly adopted the original implementation of (Bai et al., 2019a). However, the public source code was optimized for the tasks of language modeling; hence, we mainly modified the following two parts of the implementation. First, the implementation uses the adaptive softmax function (Baevski & Auli, 2019; Grave et al., 2017) to address the



Figure 4: Training losses along (a) run-time, and (b) epochs for three compared methods in the copy memory task.

Table 1: Test cross-entropy loss (test loss) and the number of learnable parameters (#params) in the copy memory task.

	Test loss	#params
DEQ	2.24e-09	17,010
SIM (single)	7.03e-09	17,294
SIM (two)	5.06e-08	17,294

large word vocabulary. Because the copy memory task considers only ten digits, we did not use this architecture in the experiment. Second, the implementation uses the memory padding and nonzero initial hidden states that employ the final states of the previous step in a training loop. We tested the two techniques but observed that the test loss was considerably worse even though the training loss was fine. Hence, we omit the techniques and instead used the empty memory and zero initial hidden states in the experiment.

Table 1 shows the test cross-entropy loss per sequence and the number of learnable parameters for each model. Although the DEQ achieved the lowest loss, the single-tier SIM obtained was comparable with almost the same number of learnable parameters. Figure 4 (a) shows the progress of the training losses along the run-time for 20 epochs when we re-trained the models. The training speed depends on the batch size. During the hyperparameter search, 1, 1, and 10 were selected as the batch size for the DEQ, single-tier SIM, and two-tier SIM, respectively. Hence, the two-tier SIM was the fastest to complete the training. For the DEQ and single-tier SIM, the single-tier model was ten times faster than the DEQ under the same batch size. Figure 4 (b) shows the progress of the training losses along the epochs. We can observe that the single-tier SIM converged faster than the DEQ.

parameters (#params) in the image classification task. We showed the results of monDEQs reported in the paper (Winston & Kolter, 2020).

Table 2: Mean of the test accuracy over three runs (test acc.) and the number of learnable

	CIFA	CIFAR-10		SVHN		MNIST	
	Test acc.	#params	Test acc.	#params	Test acc.	#params	
monDEQ (single) monDEQ (multi) SIM (single) SIM (two)	74.0 ± 0.1 72.0 ± 0.3 79.4 ± 0.2 78.2 ± 0.2	$172,218 \\ 170,194 \\ 168,694 \\ 168,264$	88.7 ± 1.1 92.4 \pm 0.1 91.8 \pm 0.2 92.4 \pm 0.1	172,218 170,194 168,694 168,264	99.1 ± 0.1 99.0 ± 0.1 99.4 ± 0.0 99.4 ± 0.0	84,460 81,394 81,480 80,466	

5.2. Image Classification

We next report the results of the image classification to compare SIMs to monDEQs. We followed the experiment of image classification conducted in (Winston & Kolter, 2020). We prepared the CIFAR-10 (Krizhevsky, 2009), SVHN (Netzer et al., 2011), and MNIST (LeCun et al., 1998) datasets, which contain images in 10 different classes and evaluated the classification performance in the standard setting. Following (Winston & Kolter, 2020), we evaluated the models on test data three times with different initialization and reported the averaged performance.

For SIMs, we employed convolutional neural networks for $\mu_{\rm NN}$. Following the VGG models (Simonyan & Zisserman, 2015), $\mu_{\rm NN}$ consists of two convolutional layers each of which has two convolution filters with ReLU activation and batch normalization and one max pooling. The output of the convolution layers is additionally transformed by a linear layer to fit the dimension of the lifted space. We constrained the number of learnable parameters of SIMs to be comparable to the one of monDEQs in (Winston & Kolter, 2020).

Table 2 lists the means of the test accuracy over three runs of monDEQs and SIMs with different initialization. SIMs showed comparable or better performances against monDEQs for all datasets. Particularly, the performance was improved for the CIFAR-10 dataset even though the SIMs have a similar number of learnable parameters to monDEQ ones.

5.3. Image Regression

Finally, we report the results of the image regression task, which is an example of implicit neural representation tasks. In Section 3.5, we found that SIMs have a close connection to the work of (Tancik et al., 2020). The purpose of the experiment is to verify 1) that the RFF only model actually works well for this task, and 2) how well the other types of SIMs perform. The

	Test PSNR				
	Natural	Text			
SIM (single)	19.67 ± 2.87	17.54 ± 2.07			
SIM (two) SIM (RFF only)	22.24 ± 3.03 25.19 ± 3.92	25.50 ± 2.72 27.69 ± 1.63			

Table 3: Mean of the test PSNR over 16 images of two types of datasets in the image regression task.

goal of the task is to obtain a neural network where the input is a 2-D pixel coordinate and the output is its 3-D RGB value. Following (Tancik et al., 2020), we evaluate SIMs with 32 datasets, where 16 are natural images ¹ and the rest are text images ². This task considers each of the images to be one dataset. For an image, we picked 1/4 pixels as training data and other 1/4 as test data. Additionally, we prepared another 1/4 pixels as validation data to select the hyper-parameters.

We compare three instances of SIMs. The first two models employ an FCN for $\mu_{\rm NN}$: it consists of two linear layers with ReLU activation and one linear layer to fit the dimension of the lifted space. The last one is the RFF only model (12).

Table 3 shows the mean of the test PSNR over 16 images of two types of datasets. The results indicate that the RFF only model was better than the single-tier SIM. This result is consistent with the original work on (Tancik et al., 2020). The two-tier SIM was also effective although the test PSNRs were slightly worse than the RFF only model. Figure 5 illustrates examples of prediction for two images. Although the single-tier SIM produced blurred images, the RFF only model can generate sharper ones. Although the images of the two-tier SIM are a little blurry, the detail can be recognized compared to the single-tier SIM.

5.4. Discussion

The first two results indicate that SIMs can provide time-effective alternatives to DEQs. This is because SIMs are implemented as feedforward models, which do not require implicit differentiation even though they can

¹https://drive.google.com/uc?id=1TtwlEDArhOMoH18aUyjIMSZ3WODFmUab

²https://drive.google.com/uc?id=1V-RQJcMuk9GD4JCUn70o7nwQE0hEzHoT



Figure 5: Examples of prediction in the image regression task.

approximate DEQs. While we compared SIMs to DEQs as fairly as possible by keeping the number of parameters nearly the same, SIMs performed slightly better than DEQs in many experiments. It may be because the chosen basis functions were suitable for the tasks, and finite-depth models were easier to evaluate and optimize. The results also showed that two-tier SIMs did not necessarily outperform single-tier SIMs for the first two tasks. One of the reasons would be that even single-tier SIMs suffice for the tasks where the models can capture the essential feature of sequences and images.

The third result demonstrated that SIMs are also effective in the task of implicit neural representation. Particularly, we observed that the two-tier SIM also obtained moderate performance. The work of (Tancik et al., 2020) does not cover the form of the two-tier models, and we believe that this result is an interesting finding.

6. Conclusions

In this study, we considered DEQs from the viewpoint of the Koopman operator, which enables us to identify stable dynamics described by DEQs via the representation of the spectra. This perspective yielded our proposed SIMs that approximated DEQs and exploited more general dynamics converging to an invariant set. Despite having such noteworthy properties, the resulting models can be represented as simple feedforward models, which will provide new insights into the studies on DEQs.

A promising future work will be to consider more theoretical analysis. Investigating expressive powers, sample efficiency, and approximation error bounds would help to further understand SIMs. Another direction will be to explore ways to exploit dynamics converging to an invariant set. While we proposed the second-tier SIMs in this paper, other approaches will be possible and could improve the performance.

Acknowledgements

This work was supported by JST CREST Grant Number JPMJCR1913 and JSPS KAKENHI Grant Numbers 18H03287, 22H00516, and 22K17950.

Appendix A. Convergent Behavior on Lifted Space

We clarify the convergent behavior of the lifted dynamics in Section 3. It should be noted that the following result is not novel; it employs a known consequence of linear algebra.

First, we define the l-th step of the lifted dynamics as

$$oldsymbol{arphi}^l\coloneqqoldsymbol{A}^loldsymbol{arphi}(oldsymbol{z}_0,oldsymbol{x}),$$

where $\varphi^0 \coloneqq \varphi(z_0, x)$. To reveal the property of φ^l , we decompose A^l with the eigenvalues and eigenvectors.

We begin with the decomposition of A. If A is diagonalizable, it can be represented by the set of the eigenvalues and corresponding eigenvectors. However, if A has repeated eigenvalues, A is not necessarily diagonalizable. Although the Jordan canonical form can be used in such a case, this form is constructed by complex matrices if A contains complex eigenvalues. To represent A by real matrices for convenience, we consider the real Jordan canonical form: any real square matrix can be written as

$$\boldsymbol{A} = \boldsymbol{U}\boldsymbol{J}\boldsymbol{U}^{-1},$$

where \boldsymbol{U} is defined in Section 3.2, and $\boldsymbol{J} \in \mathbb{R}^{N \times N}$ is the following block diagonal matrix:

$$\boldsymbol{J} = \operatorname{diag}(\boldsymbol{J}_1, \boldsymbol{J}_2, \ldots, \boldsymbol{J}_R),$$

where $R \leq N$ corresponds to the number of linearly independent eigenvectors. Each block $J_r \in \mathbb{R}^{s_r \times s_r}$ $(r = 1, \ldots, R)$ is associated with s_r repeated eigenvalues with a value $\lambda_{(r)}$ if $\lambda_{(r)}$ is real or $s_r/2$ repeated eigenvalues with a value $\lambda_{(r)}$ and their $s_r/2$ complex conjugates with a value $\overline{\lambda_{(r)}}$ if $\lambda_{(r)}$ is nonreal. If $\lambda_{(r)}$ is real, J_r is also associated with one ordinal eigenvector of $\lambda_{(r)}$ and $s_r - 1$ non-ordinal generalized-eigenvectors. If $\lambda_{(r)}$ is nonreal, J_r is associated with a complex conjugate pair of ordinal eigenvectors of $\lambda_{(r)}$ and $\overline{\lambda_{(r)}}$ and $(s_r - 1)/2$ complex conjugate pairs of non-ordinal generalized-eigenvectors. It should also be noted that there always exists a conjugate eigenvalue for every nonreal eigenvalue when A is real. J_r takes one of the following two forms:

$$oldsymbol{J}_r = egin{cases} oldsymbol{J}_{r,\mathrm{R}} & ext{if } \lambda_{(r)} ext{ is real} \ oldsymbol{J}_{r,\mathrm{C}} & ext{if } \lambda_{(r)} ext{ is nonreal}, \end{cases}$$

where

$$m{J}_{r,\mathrm{R}} = egin{pmatrix} \lambda_{(r)} & 1 & & \ & \lambda_{(r)} & 1 & & \ & & \ddots & \ddots & \ & & & \lambda_{(r)} & 1 \ & & & & \lambda_{(r)} \end{pmatrix},$$

$$m{J}_{r, ext{C}} = egin{pmatrix} m{C}_r & m{I} & & \ & m{C}_r & m{I} & & \ & & \ddots & \ddots & \ & & & m{C}_r & m{I} \ & & & m{C}_r & m{I} \ & & & m{C}_r \end{pmatrix}.$$

Here, C_r and I are defined by

$$\boldsymbol{C}_r = \begin{pmatrix} \alpha_{(r)} & \beta_{(r)} \\ -\beta_{(r)} & \alpha_{(r)} \end{pmatrix}, \quad \boldsymbol{I} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix},$$

where

$$\frac{\lambda_{(r)}}{\overline{\lambda_{(r)}}} = \alpha_{(r)} + i\beta_{(r)}$$
$$\overline{\lambda_{(r)}} = \alpha_{(r)} - i\beta_{(r)}.$$

As a special case, $J_{r,R}$ can be a scalar $\lambda_{(r)}$ when $s_r = 1$, and $J_{r,C}$ can be C_r when $s_r = 2$. In this case, $J_{r,R}$ and $J_{r,C}$ are only associated with one ordinal eigenvector and a complex conjugate pair of ordinal eigenvectors, respectively. Moreover, the matrix C_r can be rewritten as a polar form:

$$\boldsymbol{C}_{r} = r_{(r)} \begin{pmatrix} \cos \Delta_{(r)} & \sin \Delta_{(r)} \\ -\sin \Delta_{(r)} & \cos \Delta_{(r)} \end{pmatrix},$$

where

$$r_{(r)} = |\lambda_{(r)}| = \sqrt{\alpha_{(r)}^2 + \beta_{(r)}^2},$$

$$\Delta_{(r)} = \arctan(\beta_{(r)}/\alpha_{(r)}).$$

The real Jordan canonical form can represent $\boldsymbol{\varphi}^l$ as

$$oldsymbol{arphi}^l = oldsymbol{U}oldsymbol{J}^loldsymbol{U}^{-1}oldsymbol{arphi}(oldsymbol{z}_0,oldsymbol{x}),$$

where

$$\boldsymbol{J}^l = \operatorname{diag}(\boldsymbol{J}_1^l, \boldsymbol{J}_2^l, \dots, \boldsymbol{J}_R^l)$$

Each J_r^l can be written as

$$\boldsymbol{J}_{r,\mathrm{R}}^{l} = \begin{pmatrix} \lambda_{(r)}^{l} & \binom{l}{1} \lambda_{(r)}^{l-1} & \cdots & \binom{l}{s_{r}-1} \lambda_{(r)}^{l-s_{r}+1} \\ \lambda_{(r)}^{l} & \cdots & \binom{l}{s_{r}-2} \lambda_{(r)}^{l-s_{r}+2} \\ & \ddots & \ddots & \vdots \\ & & \lambda_{(r)}^{l} & \binom{l}{1} \lambda_{(r)}^{l-1} \\ & & & \lambda_{(r)}^{l} \end{pmatrix}, \quad (A.1)$$

or

$$\boldsymbol{J}_{r,\mathrm{C}}^{l} = \begin{pmatrix} \boldsymbol{C}_{r}^{l} & {\binom{l}{1}} \boldsymbol{C}_{r}^{l-1} & \cdots & {\binom{l}{s_{r}/2-1}} \boldsymbol{C}_{r}^{l-s_{r}/2+1} \\ \boldsymbol{C}_{r}^{l} & {\binom{l}{1}} \boldsymbol{C}_{r}^{l-1} & \cdots & {\binom{l}{s_{r}/2-2}} \boldsymbol{C}_{r}^{l-s_{r}/2+2} \\ & \ddots & \ddots & \ddots & \vdots \\ & & \boldsymbol{C}_{r}^{l} & {\binom{l}{1}} \boldsymbol{C}_{r}^{l-1} \\ & & & \boldsymbol{C}_{r}^{l} & {\binom{l}{1}} \boldsymbol{C}_{r}^{l-1} \\ & & & \boldsymbol{C}_{r}^{l} \end{pmatrix}, \quad (A.2)$$

where

$$C_r^l = \left(r_{(r)} \begin{pmatrix} \cos \Delta_{(r)} & \sin \Delta_{(r)} \\ -\sin \Delta_{(r)} & \cos \Delta_{(r)} \end{pmatrix} \right)^l$$
$$= r_{(r)}^l \begin{pmatrix} \cos l \Delta_{(r)} & \sin l \Delta_{(r)} \\ -\sin l \Delta_{(r)} & \cos l \Delta_{(r)} \end{pmatrix}.$$

This form makes it easy to evaluate the convergent behavior of the lifted dynamics: because the number of steps l only depends on the block diagonal matrix J^l , it suffices to focus on each block of J^l that takes the form of Eq. (A.1) or (A.2).

If the spectral radius $\rho(\mathbf{A}) < 1$, all the elements in every block of \mathbf{J}^l converge to 0 by taking the limit $l \to \infty$. Hence, the lifted dynamics in case (i) converges to the origin:

$$egin{aligned} &\lim_{l o\infty}oldsymbol{arphi}^l = oldsymbol{U}oldsymbol{O}_Noldsymbol{U}^{-1}oldsymbol{arphi}(oldsymbol{z}_0,oldsymbol{x}) \ &= oldsymbol{0}, \end{aligned}$$

where $\boldsymbol{O}_N \in \mathbb{R}^{N \times N}$ is a zero square matrix of order N.

If $\rho(\mathbf{A}) = 1$, we can classify the limit of the *r*-th block \mathbf{J}_r^l into the following six cases:

 $\langle 1 \rangle$. If $|\lambda_{(r)}| < 1$, then J_r^l converges to a zero matrix:

$$\lim_{l\to\infty}\boldsymbol{J}_r^l=\boldsymbol{O}_{s_r}.$$

 $\langle 2 \rangle$. If $\lambda_{(r)} = 1$ and $s_r = 1$, then \boldsymbol{J}_r^l converges to 1:

$$\lim_{l \to \infty} \boldsymbol{J}_r^l = \lim_{l \to \infty} \lambda_{(r)}^l = 1.$$

(3). If $\lambda_{(r)} = -1$ and $s_r = 1$, then J_r^l neither converges nor diverges and takes the form of

$$\boldsymbol{J}_r^l = \lambda_{(r)}^l = (-1)^l.$$

(4). If $\lambda_{(r)} \in \{1, -1\}$, and $s_r > 1$, then the non-diagonal elements of $\boldsymbol{J}_r^l (= \boldsymbol{J}_{r,R}^l)$, e.g. $\binom{l}{1} \lambda_{(r)}^{l-1}$, diverge.

 $\langle 5 \rangle$. If $|\lambda_{(r)}| = 1$, $\lambda_{(r)}$ is nonreal, and $s_r = 2$, then J_r^l neither converges nor diverges and takes the form of

$$\boldsymbol{J}_{r}^{l} = \boldsymbol{C}_{r}^{l} = \begin{pmatrix} \cos l\Delta_{(r)} & \sin l\Delta_{(r)} \\ -\sin l\Delta_{(r)} & \cos l\Delta_{(r)} \end{pmatrix}.$$
 (A.3)

(6). If $|\lambda_{(r)}| = 1$, $\lambda_{(r)}$ is nonreal, and $s_r > 2$, then the non-block-diagonal parts of $J_r^l (= J_{r,C}^l)$, e.g. $\binom{l}{l} C_r^{l-1}$, diverge.

Hence, if at least one block falls into cases $\langle 4 \rangle$ or $\langle 6 \rangle$, the lifted dynamics diverges: the complexity of the non-diagonal elements or non-block-diagonal parts is at least $\mathcal{O}(l)$. This case falls into case (iv).

Case (ii) means all the blocks are either case $\langle 1 \rangle$ or $\langle 2 \rangle$, and each block in case $\langle 2 \rangle$ corresponds to one of the eigenvalues in J_1 . Hence, the eigenvectors with eigenvalues of $\lambda_j = 1$ are linearly independent. In the end, J^l converges to a diagonal matrix $P_{J_1} \in \mathbb{R}^{N \times N}$ where the diagonal elements corresponding to J_1 are 1 and the rest are 0, and we obtain Eq. (6):

$$egin{aligned} &\lim_{l o\infty}oldsymbol{arphi}^l = oldsymbol{U}oldsymbol{P}_{J_1}oldsymbol{U}^{-1}oldsymbol{arphi}(oldsymbol{z}_0,oldsymbol{x}) \ &= \sum_{j\in J_1}oldsymbol{u}_joldsymbol{v}_j^ opoldsymbol{arphi}(oldsymbol{z}_0,oldsymbol{x}). \end{aligned}$$

Case (iii) means that all the blocks are either case $\langle 1 \rangle$, $\langle 2 \rangle$, $\langle 3 \rangle$, or $\langle 5 \rangle$. Each block in cases $\langle 3 \rangle$ and $\langle 5 \rangle$ corresponds to one of the eigenvalues of J_2 and the complex conjugate pairs of the eigenvalues of J_3 , respectively. Hence, the eigenvectors with eigenvalues of $|\lambda_j| = 1$ are linearly independent. In addition to \mathbf{P}_{J_1} , if we denote $\mathbf{P}_{J_2} \in \mathbb{R}^{N \times N}$ as a diagonal matrix where the diagonal elements corresponding to J_2 are -1 and the rest are 0 and $\mathbf{P}_{J_3} \in \mathbb{R}^{N \times N}$ as a block diagonal matrix where the diagonal blocks corresponding to the pairs of J_3 take the form of Eq. (A.3) and the rest are zero, we obtain Eq. (7):

$$egin{aligned} &\lim_{l o\infty}oldsymbol{arphi}^l = \lim_{l o\infty}oldsymbol{U}(oldsymbol{P}_{J_1}+oldsymbol{P}_{J_2}+oldsymbol{P}_{J_3})oldsymbol{U}^{-1}oldsymbol{arphi}(oldsymbol{z}_0,oldsymbol{x}) \ &= \lim_{l o\infty}igg(\sum_{j\in J_1}oldsymbol{u}_joldsymbol{v}_j^\top+\sum_{j\in J_2}(-1)^loldsymbol{u}_joldsymbol{v}_j^\top \ &+ \sum_{(j,k)\in J_3}igl(\cos(l\Delta_j)oldsymbol{u}_j-\sin(l\Delta_k)oldsymbol{u}_kigr)oldsymbol{v}_j^\top \ &+ igl(\sin(l\Delta_j)oldsymbol{u}_j+\cos(l\Delta_k)oldsymbol{u}_kigr)oldsymbol{arphi}(oldsymbol{z}_0,oldsymbol{x}). \end{aligned}$$

Table B.4: Settings of Tune and ASHA for each task.

	Copy memory	Image classification	Image regression
#trials	100	200	200
Scope	last	last	last
$\#$ max_t (ASHA)	6	80	100
Metric (ASHA)	loss	accuracy	PSNR
Grace period (ASHA)	3	5	5
Reduction factor (ASHA)	2	2	2

Lastly, if $\rho(\mathbf{A}) > 1$, at least one block of \mathbf{J}^l diverges. This case falls into case (iv).

Appendix B. Experimental Details

We show the detailed settings and configurations of the experiments in Section 5. It should be noted that we used the default values of Pytorch for the arguments not mentioned below.

Table B.4 shows that the settings of Tune and ASHA for each task. #trials denotes the number of trials in running a search algorithm. Scope denotes the method of selecting the best hyper-parameters. We used **last** that selects the best one by comparing the last performance at the end of training. #max_t denotes the number of maximum epochs or iterations in a trial. ASHA periodically stops trials when the specified metric is poor and reduces them by a factor of the reduction factor. However, any trial is run until the grace period. Those settings are partially different for each task but the same as compared methods in a task.

Table B.5 shows the detailed architecture of the reverse model $\nu_{\rm NN}$ in SIMs. The first column denotes the layers of $\nu_{\rm NN}$, which transform the input in order from the top to the bottom layers. "× 2" means that the same transformation has been repeated twice. The second column denotes detailed information about the layers. The variable #hidden units is set to a different value for each task.

Appendix B.1. Copy Memory Task

Table B.6 lists the statistics of the dataset in the copy memory task. #train, #valid, and #test denote the number of training, validation, and test data, respectively. We first trained the models with the candidates of hyper-parameters on the training data and selected the best hyper-parameter

Table B.5: Architecture of the reverse model $\nu_{\rm NN}$ in SIMs for all tasks.

	Detail
Fully connected layer $\times 2$ Linear	output features: #hidden units
ReLU Linear layer Linear	output features: #hidden units

-

on the validation data. We then used both training and validation data to re-train the model with the best hyper-parameter and evaluated it on the test data. A data point of this dataset consists of the pair of an input sequence of length 520 and an output sequence of length 520. We generated all the data points by following the procedure of (Bai et al., 2018).

Table B.7 shows the model architecture of the TCN. The first building block is the temporal convolution layer where we adopted the architecture of (Bai et al., 2018): it mainly consists of the 1D dilated causal convolution, ReLU activation, and residual connection, and down-sampling with 1D convolution is applied to the residual connection (He et al., 2016) except for the first layer. We stacked this layer eight times and then connected a linear layer to fit the dimension of the lifted space.

Table B.8 lists the configurations of the compared models. The upper and lower rows denote the configuration of the model architectures and training algorithms, respectively. The configuration with the set notation (e.g. [1e-3, 1e+3 and $\{1, 2, 5, 10, 20\}$ means that the corresponding variable was tuned as a hyper-parameter within the range of the set. For each trial, ASHA sampled a candidate from the set uniformly at random. For the configuration with a real interval, it was sampled in the logarithmic space with base 10. In this task, we tuned the bandwidth of the RFF, batch size, and learning rate. The configuration of the DEQ indicates that the variables of the implementation of (Bai et al., 2019a): we set it to ensure almost the same number of learnable parameters as SIMs. For the training algorithm of the DEQ, we adopted Adam (Kingma & Ba, 2015) with the step-wise cosine decaying schedule for the learning rate as in (Bai et al., 2019a). For SIMs, we used Adam with a constant learning rate schedule and slightly larger ϵ to improve the stability of the training algorithm. We also used gradient clipping with the default values of the DEQ and TCN. For all compared models, we minimized the cross-entropy loss, i.e. the negative log-likelihood with the softmax function. Although every model contains the Dropout

Table B.6: Statistics of the dataset in the copy memory task.

#train	#valid	#test	Input dim.	Output dim.
4,500	500	500	520×1	520×1

Table B.7: Architecture of the TCN in SIMs for the copy memory task.

	Details
Temporal convolution layer $\times 8$	
1D dilated causal convolution	kernel size: 8, output channels: 10
Weight normalization	
Chomp	
ReLU	
Dropout	
1D dilated causal convolution	kernel size: 8, output channels: 10
Weight normalization	
Chomp	
ReLU	
Dropout	
Residual connection	
ReLU	
Linear layer	
Linear	out features: N

layer (Srivastava et al., 2014), we did not apply it for all models.

Appendix B.2. Image Classification

Table B.9 shows the statistics of three datasets in the image classification task. We pre-processed these datasets in almost the same way as (Winston & Kolter, 2020). The difference from (Winston & Kolter, 2020) is to additionally prepare the validation data. We initially used 90% of the prepared training data as training data and the remaining 10% as validation data. We then used both of them to train models with the best hyper-parameters for evaluating the performance on the test data.

Table B.10 shows the architecture of $\mu_{\rm NN}$ in SIMs. We first used a combination of convolution and max-pooling layers, which follows the VGG architecture (Simonyan & Zisserman, 2015) and then applied a linear layer to fit the dimension of the lifted space.

Table B.11 lists the configuration of SIMs. In this task, we tuned the bandwidth of the RFF, the batch size, the number of epochs, the learning rate, and the learning rate schedule. We set the model architectures of SIMs as almost the same size as that of the compared monDEQs. We also applied the training algorithm of the monDEQs to SIMs: the learning rate schedules

Table B.8:	Configurations	of com	pared :	models i	in the	copy	memory	task.
	0 0 0 0					~~~~~./		

	DEQ	SIM (single-tier)	SIM (two-tier)
Model architectures			
n_head	8	-	-
d_head	5	-	-
d_model	40	-	-
d_inner	40	-	-
pre_lnorm	True	-	-
wnorm	True	-	-
f_thres	30	-	-
#pretraining steps	0	-	-
N	-	32	32
M	-	-	32
Bandwidth of RFF	-	-	[1e-3, 1e+3)
#hidden units of $ u_{\rm NN}$	-	32	32
Dropout rate	0.0	0.0	0.0
Training algorithms			
Batch size	$\{1, 2, 5, 10, 20\}$	$\{1, 2, 5, 10, 20\}$	$\{1, 2, 5, 10, 20\}$
#epochs	20	20	20
Optimizer	Adam	Adam (ϵ =1e-5)	Adam (ϵ =1e-5)
Learning rate	[1e-4, 0.05)	[1e-4, 0.05)	[1e-4, 0.05)
Learning rate schedule	$\cos(\text{step-wise})$	constant	constant
Gradient clipping value	0.25	1.0	1.0

Table B.9: Statistics of three datasets in the image classification task.

	#train	#valid	#test	Input dim.	Output dim.
CIFAR-10 SVHN MNIST	$\begin{array}{c} 45,000 \\ 65,931 \\ 54,000 \end{array}$	$5,000 \\ 7,326 \\ 6,000$	$10,000 \\ 26,032 \\ 10,000$	$\begin{array}{c} 32 \times 32 \times 3 \\ 32 \times 32 \times 3 \\ 28 \times 28 \times 1 \end{array}$	10×1 10×1 10×1

are also implemented in the same way as the monDEQs. As in the monDEQs, we minimized the cross-entropy loss for SIMs.

Appendix B.3. Image Regression

Table B.12 shows the statistics of a dataset in the image regression task. All the images have 512×512 pixels, and we used equally spaced 1/4 (65,536) pixels from an image as training data, another 1/4 pixels as validation data, and other 1/4 pixels as test data. It should be noted that all the images are the same size and thus have the same statistics.

Table B.13 shows the architecture of $\mu_{\rm NN}$ in SIMs. We used a neural network that consists of two fully connected layers and one additional linear layer to fit the dimension of the lifted space.

Table B.14 shows the configuration of SIMs. In this task, we tuned the number of hidden units of μ_{NN} , the first and second lift dimensions (i.e. N

Table B.10: Architecture of $\mu_{\rm NN}$ in SIMs for the image classification task.

	Details			
Convolution layer $\times 2$				
2D convolution	kernel size: 3, padding: 1, output channels: #channels			
Batch norm.				
ReLU				
Max pooling layer	kernel size: 2, stride: 2			
Convolution layer $\times 2$				
2D convolution	kernel size: 3, padding: 1, output channels: $\#$ channels			
Batch norm.				
ReLU				
Max pooling layer	kernel size: 2, stride: 2			
Linear layer				
Flatten				
Linear	out features: N			

Table B.11: Configurations of SIMs in the image classification task.

	SIM (single-tier)	SIM (two-tier)
Model architectures		
#channels of $\mu_{ m NN}$	24 (MNIST) or 38 (others)	21 (MNIST) or 36 (others)
N	50	50
M	-	104
Bandwidth of RFF	-	[1e-3, 1e+3)
#hidden units of $\nu_{\rm NN}$	32	32
Training algorithms		
Batch size	$\{64, 128, 256, 512\}$	$\{64, 128, 256, 512\}$
#epochs	$\{20, 40, 60, 80\}$	$\{20, 40, 60, 80\}$
Optimizer	Adam	Adam
Learning rate	[1e-4, 0.05)	[1e-4, 0.05)
Learning rate schedule	$\{1cycle, step, constant\}$	$\{1cycle, step, constant\}$
Step size for step schedule	$\{5i \mid i \in (1, 2, \dots 10)\}$	$\{5i \mid i \in (1, 2, \dots 10)\}$

and M), the bandwidth of the RFF, and the learning rate. Following (Tancik et al., 2020), we applied Adam with the constant learning rate and performed the full batch gradient descent which uses all data points (pixels) for the update of parameters. #iterations denotes the number of steps in the gradient descent and was set to 2000 as in (Tancik et al., 2020). We minimized the mean squared loss for all models. The difference from (Tancik et al., 2020) is that we did not apply the sigmoid function before the output for each model: we observed such a setting improved the performance.

References

Almeida, L. B. (1987). A learning rule for asynchronous perceptrons with feedback in a combinatorial environment. In *IEEE First International*

Table B.12: Statistics of each dataset in the image regression task. It should be noted that all 32 datasets have the same statistics.

#train	#valid	#test	Input dim.	Output dim.
$65,\!536$	$65,\!536$	$65,\!536$	2×1	3×1

Table B.13: Architecture of $\mu_{\rm NN}$ in SIMs for the image regression task.

	Details	
Fully connected layer × 2 Linear ReLU	output features: #hidden units	
Linear layer Linear	output features: N	

Conference on Neural Networks (pp. 609–618).

- Baevski, A., & Auli, M. (2019). Adaptive input representations for neural language modeling. In 7th International Conference on Learning Representations.
- Bai, S., Kolter, J. Z., & Koltun, V. (2018). An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. arXiv preprint arXiv:1803.01271.
- Bai, S., Kolter, J. Z., & Koltun, V. (2019a). Deep equilibrium models. In Advances in Neural Information Processing Systems 32 (pp. 688–699).
- Bai, S., Kolter, J. Z., & Koltun, V. (2019b). Trellis networks for sequence modeling. In 7th International Conference on Learning Representations.
- Bai, S., Koltun, V., & Kolter, J. Z. (2020). Multiscale deep equilibrium models. In Advances in Neural Information Processing Systems 33 (pp. 5238–5250).
- Bai, S., Koltun, V., & Kolter, J. Z. (2021). Stabilizing equilibrium models by Jacobian regularization. In *Proceedings of the 38th International Conference on Machine Learning* (pp. 554–565).
- Billings, L., & Schwartz, I. B. (2008). Identifying almost invariant sets in stochastic dynamical systems. *Chaos*, 18.

Table B.14: Configurations of SIMs in the image regression task.

	SIM (single-tier)	SIM (two-tier)	SIM (RFF only)
Model architectures			
#hidden units of $\mu_{ m NN}$	$\{32, 64, 128, 256, 512\}$	$\{32, 64, 128, 256, 512\}$	-
N	$\{32, 64, 128, 256, 512\}$	$\{32, 64, 128, 256, 512\}$	$\{32, 64, 128, 256, 512\}$
M	-	$\{32, 64, 128, 256, 512\}$	-
Bandwidth of RFF	-	[1e-3, 1e+3)	[1e-3, 1e+3)
#hidden units of $ u_{\rm NN}$	256	256	256
Training algorithms			
#iterations	2000	2000	2000
Optimizer	Adam	Adam	Adam
Learning rate	[1e-4, 0.05)	[1e-4, 0.05)	[1e-4, 0.05)
Learning rate schedule	constant	constant	constant

- Chang, B., Chen, M., Haber, E., & Chi, E. H. (2019). AntisymmetricRNN: A dynamical system view on recurrent neural networks. In 7th International Conference on Learning Representations.
- Chen, T. Q., Rubanova, Y., Bettencourt, J., & Duvenaud, D. (2018). Neural ordinary differential equations. In Advances in Neural Information Processing Systems 31 (pp. 6572–6583).
- Cvitanović, P., Artuso, R., Mainieri, R., Tanner, G., & Vattay, G. (2020). *Chaos: Classical and Quantum.* ChaosBook.org. Niels Bohr Institute, Copenhagen.
- Dabre, R., & Fujita, A. (2019). Recurrent stacking of layers for compact neural machine translation models. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence* (pp. 6292–6299).
- Dehghani, M., Gouws, S., Vinyals, O., Uszkoreit, J., & Kaiser, L. (2019). Universal transformers. In 7th International Conference on Learning Representations.
- Dogra, A. S., & Redman, W. T. (2020). Optimizing neural networks via Koopman operator theory. In Advances in Neural Information Processing Systems 33 (pp. 2087–2097).
- Froyland, G., & Padberg, K. (2009). Almost-invariant sets and invariant manifolds—connecting probabilistic and geometric descriptions of coherent structures in flows. *Physica D: Nonlinear Phenomena*, 238, 1507–1523.
- Gaspard, P. (1998). Chaos, Scattering and Statistical Mechanics volume 9 of Cambridge Nonlinear Science Series. Cambridge University Press.

- Gori, M., Monfardini, G., & Scarselli, F. (2005). A new model for learning in graph domains. In *Proceedings of 2005 IEEE International Joint Conference on Neural Networks, vol. 2* (pp. 729–734).
- Grave, E., Joulin, A., Cissé, M., Grangier, D., & Jégou, H. (2017). Efficient softmax approximation for GPUs. In *Proceedings of the 34th International Conference on Machine Learning* (pp. 1302–1310).
- He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer* Vision and Pattern Recognition (pp. 770–778).
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. Neural Computation, 9, 1735–1780.
- Kag, A., Zhang, Z., & Saligrama, V. (2020). RNNs incrementally evolving on an equilibrium manifold: A panacea for vanishing and exploding gradients? In 8th International Conference on Learning Representations.
- Kawaguchi, K. (2021). On the theory of implicit deep learning: Global convergence with implicit layers. In 9th International Conference on Learning Representations.
- Kawahara, Y. (2016). Dynamic mode decomposition with reproducing kernels for Koopman spectral analysis. In Advances in Neural Information Processing Systems 29 (pp. 911–919).
- Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations.
- Koopman, B. O. (1931). Hamiltonian systems and transformation in Hilbert space. Proceedings of the National Academy of Sciences of the United States of America, 17, 315–318.
- Krantz, S. G., & Parks, H. R. (2013). The Implicit Function Theorem: History, Theory, and Applications. Modern Birkhäusr Classics. Birkhäusr.
- Krizhevsky, A. (2009). Learning multiple layers of features from tiny images.
- Lasota, A., & Mackey, M. C. (1994). Chaos, Fractals, and Noise: Stochastic Aspects of Dynamics volume 97 of Applied Mathematical Sciences. Springer New York, NY.

- LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86, 2278–2324.
- Li, L., Jamieson, K. G., Rostamizadeh, A., Gonina, E., Ben-tzur, J., Hardt, M., Recht, B., & Talwalkar, A. (2020). A system for massively parallel hyperparameter tuning. In *Proceedings of Machine Learning and Systems* (pp. 230–246). volume 2.
- Liao, R., Xiong, Y., Fetaya, E., Zhang, L., Yoon, K., Pitkow, X., Urtasun, R., & Zemel, R. (2018). Reviving and improving recurrent back-propagation. In *Proceedings of the 35th International Conference on Machine Learning* (pp. 3082–3091).
- Liaw, R., Liang, E., Nishihara, R., Moritz, P., Gonzalez, J. E., & Stoica, I. (2018). Tune: A research platform for distributed model selection and training. arXiv preprint arXiv:1807.05118.
- Linsley, D., Ashok, A. K., Govindarajan, L. N., Liu, R., & Serre, T. (2020). Stable and expressive recurrent vision models. In Advances in Neural Information Processing Systems 33 (pp. 10456–10467).
- Manek, G., & Kolter, J. Z. (2019). Learning stable deep dynamics models. In Advances in Neural Information Processing Systems 32 (pp. 10718–10728).
- Manojlović, I., Fonoberova, M., Mohr, R., Andrejcuk, A., Drmăc, Z., Kevrekidis, Y., & Mezić, I. (2020). Applications of Koopman mode analysis to neural networks. arXiv preprint arXiv:2006.11765.
- Mauroy, A., Mezić, I., & Susuki, Y. (2020). The Koopman Operator in Systems and Control. Springer International Publishing.
- Mezić, I. (2005). Spectral properties of dynamical systems, model reduction and decompositions. *Nonlinear Dynamics*, 41, 309–325.
- Nair, V., & Hinton, G. E. (2010). Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference* on Machine Learning (pp. 807–814).
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., & Ng, A. Y. (2011). Reading digits in natural images with unsupervised feature learning. In NIPS Workshop on Deep Learning and Unsupervised Feature Learning.

- Pabbaraju, C., Winston, E., & Kolter, J. Z. (2021). Estimating Lipschitz constants of monotone deep equilibrium models. In 9th International Conference on Learning Representations.
- Papamakarios, G., Nalisnick, E. T., Rezende, D. J., Mohamed, S., & Lakshminarayanan, B. (2021). Normalizing flows for probabilistic modeling and inference. *Journal of Machine Learning Research*, 22, 1–64.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., & Chintala, S. (2019). PyTorch: An imperative style, highperformance deep learning library. In Advances in Neural Information Processing Systems 32 (pp. 8024–8035).
- Pineda, F. J. (1987). Generalization of back-propagation to recurrent neural networks. *Physical review letters*, 59, 2229–2232.
- Proctor, J. L., Brunton, S. L., & Kutz, J. N. (2018). Generalizing Koopman theory to allow for inputs and control. SIAM J. Appl. Dyn. Syst., 17, 909–930.
- Rahimi, A., & Recht, B. (2007). Random features for large-scale kernel machines. In Advances in Neural Information Processing Systems 20 (pp. 1177–1184).
- Rowley, C. W., Mezić, I., Bagheri, S., Schlatter, P., & Henningson, D. S. (2009). Spectral analysis of nonlinear flows. *Journal of Fluid Mechanics*, 641, 115–127.
- Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., & Monfardini, G. (2008). The graph neural network model. *IEEE Transactions on Neural Networks*, 20, 61–80.
- Schmid, P. J. (2010). Dynamic mode decomposition of numerical and experimental data. Journal of Fluid Mechanics, 656, 5–28.
- Selverston, A. I., & Moulins, M. (1985). Oscillatory neural networks. Annual Review of Physiology, 47, 29–48.

- Simonyan, K., & Zisserman, A. (2015). Very deep convolutional networks for large-scale image recognition. In 3rd International Conference on Learning Representations.
- Sitzmann, V., Martel, J. N. P., Bergman, A. W., Lindell, D. B., & Wetzstein, G. (2020). Implicit neural representations with periodic activation functions. In Advances in Neural Information Processing Systems 33 (pp. 7462–7473).
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15, 1929–1958.
- Strogatz, S. H. (2015). Nonlinear Dynamics and Chaos: With Applications to Physics, Biology, Chemistry, and Engineering, Second Edition. CRC press.
- Takeishi, N., & Kawahara, Y. (2021). Learning dynamics models with stable invariant sets. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence* (pp. 9782–9790).
- Takeishi, N., Kawahara, Y., Tabei, Y., & Yairi, T. (2017a). Bayesian dynamic mode decomposition. In Proceedings of the 26th International Joint Conference on Artificial Intelligence (pp. 2814–2821).
- Takeishi, N., Kawahara, Y., & Yairi, T. (2017b). Learning Koopman invariant subspaces for dynamic mode decomposition. In Advances in Neural Information Processing Systems 30 (pp. 1130–1140).
- Tancik, M., Srinivasan, P. P., Mildenhall, B., Fridovich-Keil, S., Raghavan, N., Singhal, U., Ramamoorthi, R., Barron, J. T., & Ng, R. (2020). Fourier features let networks learn high frequency functions in low dimensional domains. In Advances in Neural Information Processing Systems 33 (pp. 7537–7547).
- Townley, S., Ilchmann, A., Weiß, M. G., McClements, W., Ruiz, A. C., Owens, D. H., & Pratzel-Wolters, D. (2000). Existence and learning of oscillations in recurrent neural networks. *IEEE Transactions on Neural Networks*, 11, 205–214.

- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention is all you need. In Advances in Neural Information Processing Systems 30 (pp. 5998–6008).
- Williams, M. O., Kevrekidis, I. G., & Rowley, C. W. (2015). A data–driven approximation of the Koopman operator: Extending dynamic mode decomposition. *Journal of Nonlinear Science*, 25, 1307–1346.
- Winston, E., & Kolter, J. Z. (2020). Monotone operator equilibrium networks. In Advances in Neural Information Processing Systems 33 (pp. 10718–10728).
- Xie, X., Wang, Q., Ling, Z., Li, X., Wang, Y., Liu, G., & Lin, Z. (2021). Optimization induced equilibrium networks. arXiv preprint arXiv:2105.13228.