

Explainable Intelligent Fault Diagnosis for Nonlinear Dynamic Systems: From Unsupervised to Supervised Learning

Hongtian Chen^{ID}, *Member, IEEE*, Zhigang Liu^{ID}, *Senior Member, IEEE*, Cesare Alippi^{ID}, *Fellow, IEEE*,
Biao Huang^{ID}, *Fellow, IEEE*, and Derong Liu^{ID}, *Fellow, IEEE*

Abstract—The increased complexity and intelligence of automation systems require the development of intelligent fault diagnosis (IFD) methodologies. By relying on the concept of a suspected space, this study develops explainable data-driven IFD approaches for nonlinear dynamic systems. More specifically, we parameterize nonlinear systems through a generalized kernel representation for system modeling and the associated fault diagnosis. An important result obtained is a unified form of kernel representations, applicable to both unsupervised and supervised learning. More importantly, through a rigorous theoretical analysis, we discover the existence of a *bridge* (i.e., a bijective mapping) between some supervised and unsupervised learning-based entities. Notably, the designed IFD approaches achieve the same performance with the use of this bridge. In order to have a better understanding of the results obtained, both unsupervised and supervised neural networks are chosen as the learning tools to identify the generalized kernel representations and design the IFD schemes; an invertible neural network is then employed to build the bridge between them. This article is a perspective article, whose contribution lies in proposing and formalizing the fundamental concepts for explainable intelligent learning methods, contributing to system modeling and data-driven IFD designs for nonlinear dynamic systems.

Index Terms—Intelligent fault diagnosis (IFD), neural networks, nonlinear dynamic systems, supervised learning, bridge, unsupervised learning.

I. INTRODUCTION

OVER the past three decades, fault diagnosis has undergone tremendous development [1], [2], [3], [4],

Manuscript received 27 January 2022; revised 18 July 2022; accepted 22 August 2022. This work was supported in part by the Natural Sciences and Engineering Research Council of Canada. (*Corresponding authors: Zhigang Liu; Biao Huang.*)

Hongtian Chen and Biao Huang are with the Department of Chemical and Materials Engineering, University of Alberta, Edmonton, AB T6G 1H9, Canada (e-mail: hongtian.chen@ieee.org; biao.huang@ualberta.ca).

Zhigang Liu is with the Institute of Rail Transit, Tongji University, Shanghai 201804, China, and also with the School of Electrical Engineering, Southwest Jiaotong University, Chengdu 611756, China (e-mail: liuzg_cd@126.com).

Cesare Alippi is with the Faculty of Informatics, Università della Svizzera italiana, 69000 Lugano, Switzerland, and also with the Department of Electronics, Information and Bioengineering, Politecnico di Milano, 20133 Milan, Italy (e-mail: alippi@elet.polimi.it).

Derong Liu is with the School of Automation, Guangdong University of Technology, Guangzhou 510006, China (e-mail: derong@gdut.edu.cn).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TNNLS.2022.3201511>.

Digital Object Identifier 10.1109/TNNLS.2022.3201511

[5], [6], [7]. Due to the increasing demands of safe and economic operations, it is playing an essential role in system performance evaluation and maintenance [1]. Fault diagnosis has a broad spectrum of applications, for instance, ranging from chemical processes [8], electrical systems [9], intelligent transportation [10], and aerospace engineering [11] to medical imaging [12].

Fault diagnosis techniques have undergone a dramatic evolution with the development of control theory, computer science, big data, sensor technology, machine learning, and so on [13]. They are becoming more intelligent and diversified. A unified description of these developments is so-called intelligent fault diagnosis (IFD), which can be model-based [14], signal processing-based [15], and data-driven [16].

The success of deep learning has further promoted fault diagnosis through neural networks. These IFD approaches, as mentioned in [1], can accomplish tasks similar to what a human can do. Deep neural networks have the ability to extract helpful features so that distinguishing faulty conditions from normal conditions becomes simpler. The accompanying problem is an additional demand for sufficient labeled data [13]. Furthermore, neural networks have also been a popular tool in adaptive dynamic programming that is a novel approximate optimal control strategy for nonlinear systems [17], [18], [19]. Following the idea of value and policy iterations [20], Song *et al.* [21] proposed an off-policy reinforcement learning algorithm for nonlinear systems with completely unknown dynamics. Recently, these strategies were introduced into fault-tolerant control [22] to enhance the system operation performance in faulty conditions through online optimization.

Depending on the operating range, neural networks-based IFD approaches can be divided into static and dynamic methods. In [23], a recurrent neural network was used to detect and identify faults of railway track circuits. In order to diagnose incipient interturn faults, a probabilistic neural network was adopted in [24] with consideration of the network size. By using a deep convolutional neural network, Liu *et al.* [25] proposed a diagnosis method for the loose strands of the isoelectric line. Taking the topological structure of system data into account, a graph convolutional network was developed in [26] to diagnose machine faults. Most recently, Chen *et al.* [27] proposed two fault detection schemes, where the first design is based on the finite impulse response filter

using a fully connected neural network and the second one constructs a recursive residual generator using a recurrent neural network.

Despite the aforementioned achievements, much of the inherent difficulty of the IFD designs arises from the following.

- 1) The presence of system nonlinearity and dynamics.
- 2) Fault diagnosis cannot be simply treated as a classification problem.
- 3) Lack of interpretability in many IFD approaches since a good representation should possess explainable posterior knowledge.

These challenges motivate the study in this perspective article with the following fourfold contributions.

- 1) Data-based modeling and construction of residual generators via the composite operators, enabling the learning procedures to be described quantitatively.
- 2) A *bridge* is built for the construction of the generalized kernel representations, based on which unsupervised and supervised learning-based IFD approaches can be transferred between each other. An additional value of the bridge representation is the evaluation of the performance of nonlinear IFD approaches.
- 3) Both unsupervised and supervised neural networks are employed in designing the two specific IFD algorithms, whose purpose is to help us understand their fundamentals.
- 4) An invertible neural network is proposed to build a bridge that allows unsupervised and supervised neural network-based IFD approaches to be bijectively connected.

The remainder of this study is organized as follows. Section II introduces some metrics, machine learning, and nonlinear dynamic systems. With the aid of composite operators, Section III is dedicated to constructing the bridge between unsupervised and supervised learning-based IFD approaches. Section IV details two specific IFD algorithms, respectively, using unsupervised and supervised neural networks, followed by an invertible neural network-based bridge. Section V concludes this study and delineates research opportunities for IFD.

Notations: All notations in the article are standard. \mathcal{M} denotes a metric, whose subscript signifies a specific form; \mathcal{R}^κ represents the space of real κ -dimensional vectors; $\text{Tr}(\cdot)$ is the trace operator; $|\cdot|$ refers to the absolute value; \circ is the cascade connection of multiple operators; superscript “f” refers to faulty conditions; $\hat{\kappa}$ is the estimate of κ ; $\text{Pr}(\cdot)$ is the probability; and $\begin{pmatrix} \psi_1 & \psi_2 \\ \psi_3 & \psi_4 \end{pmatrix}$ is a composite operator (a matrix if and only if ψ_i is a linear operator) obtained by stacking the operator ψ_i , $i = 1, \dots, 4$, in a suitable manner.

II. PRELIMINARIES AND PROBLEM FORMULATION

A metric, also called a measure, is essential for obtaining the objective function of the machine learning approaches. Therefore, four representative metrics are introduced in this section, followed by the description of the nonlinear dynamic systems of interest.

A. Metrics

Metrics are the measures that quantitatively assess, compare, and track performance [28]. In the following, several metrics are introduced for the purposes such as constructing residual signals, defining test statistics, and quantifying the influences of the unknown faults.

Consider two variables $\psi_1 \in \mathcal{R}^{k_{\psi_1}}$ and $\psi_2 \in \mathcal{R}^{k_{\psi_2}}$. If a metric \mathcal{M} is chosen as the Euclidean distance, one has

$$\mathcal{M}_{\text{euc}} = \|\psi_1 - \psi_2\|_2 \quad (1)$$

provided $k_{\psi_1} = k_{\psi_2}$. By introducing a covariance matrix Σ_ψ , the Mahalanobis distance can be defined as

$$\mathcal{M}_{\text{mah}}^2 = (\psi_1 - \psi_2)^T \Sigma_\psi^{-1} (\psi_1 - \psi_2) \quad (2)$$

where ψ_1 and ψ_2 are drawn from the same distribution of covariance Σ_ψ . By setting $\Sigma_\psi^{-1} = \mathbf{I}$, (3) becomes

$$\mathcal{M}_{\text{spe}}^2 = (\psi_1 - \psi_2)^T (\psi_1 - \psi_2) \quad (3)$$

which is called squared prediction error (SPE). A metric can also be defined by the correlation such as [29]

$$\mathcal{M}_{\text{cor}}(\psi_1, \psi_2) = k_{\psi_1} - \text{Tr}\left(\Sigma_{\psi_1}^{-1/2} \Sigma_{\psi_1, \psi_2} \Sigma_{\psi_2}^{-1/2}\right) \quad (4a)$$

$$\mathcal{M}_{\text{cor}}(\psi_2, \psi_1) = k_{\psi_2} - \text{Tr}\left(\Sigma_{\psi_2}^{-1/2} \Sigma_{\psi_2, \psi_1} \Sigma_{\psi_1}^{-1/2}\right) \quad (4b)$$

where Σ_{ψ_1, ψ_2} is the covariance matrix between ψ_1 and ψ_2 , and k_{ψ_1} may not equal to k_{ψ_2} .

B. Unsupervised and Supervised Machine Learning

Depending upon linear or nonlinear mappings, machine learning is generally divided into linear and nonlinear approaches. An important step is to evaluate performance through the defined matrices (such as \mathcal{M}_{euc} , \mathcal{M}_{mah} , \mathcal{M}_{spe} , and \mathcal{M}_{cor}), and the subsequent step will be devoted to revealing the relationship between unsupervised and supervised machine learning approaches.

Define $\psi_2 = \mathcal{S}\psi_3$, where \mathcal{S} represents the mapping of a supervised learning method. By minimizing (1), the following relation holds:

$$\begin{aligned} \min \mathcal{M}_{\text{euc}} &= \min \left\| \begin{bmatrix} \mathbf{I} & -\mathcal{S} \end{bmatrix} \begin{bmatrix} \psi_1 \\ \psi_3 \end{bmatrix} \right\|_2 \\ &= \min \left\| \underbrace{\mathcal{P} \circ \begin{bmatrix} \mathbf{I} & -\mathcal{S} \end{bmatrix}}_{\text{unsupervised learning}} \begin{bmatrix} \psi_1 \\ \psi_3 \end{bmatrix} \right\|_2 \end{aligned} \quad (5)$$

w.r.t. \mathcal{S} , where \mathcal{P} is any (linear or nonlinear) operator that does not cause information loss or change the minimization of (5). In (1) and (5), ψ_1 is a reference signal; ψ_2 and ψ_3 are the output and input of supervised learning models, respectively. By introducing \mathcal{P} , $\mathcal{P} \circ [\mathbf{I} \quad -\mathcal{S}]$ is a composite operator whose input is $[\psi_1^T \quad \psi_3^T]^T$. Instead of optimizing w.r.t. \mathcal{S} , we can treat $\mathcal{P} \circ [\mathbf{I} \quad -\mathcal{S}]$ as a whole mapping to be learned. In this way, the objective function given in (5) can be minimized through an unsupervised learning approach. Similarly, the objective functions defined by other metrics (including \mathcal{M}_{mah} , \mathcal{M}_{spe} , and \mathcal{M}_{cor}) can also be considered.

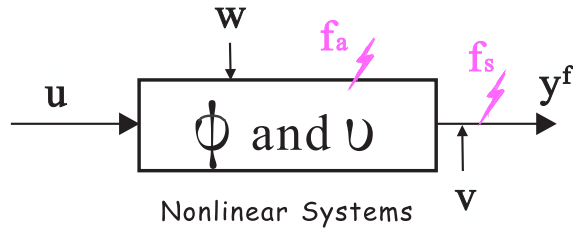


Fig. 1. Nonlinear dynamic systems with actuator and sensor faults.

It is seen that, through simple mathematical derivations, an objective function can be achieved by using both the supervised and unsupervised learning approaches.

Remark 1: The concept behind the aforementioned transformations is a bridge linking the unsupervised to supervised machine learning methods. This fact motivates this study to focus on the development of a unified framework of the IFD and related parameter-identification approaches for nonlinear systems. ▽

C. Nonlinear Dynamic Systems and Objectives

Consider a nonlinear dynamic system driven by

$$\begin{aligned} \mathbf{x}(k+1) &= \phi(\mathbf{x}(k), \mathbf{u}(k), \mathbf{w}(k)) \\ \mathbf{y}(k) &= v(\mathbf{x}(k), \mathbf{u}(k)) + \mathbf{v}(k) \end{aligned} \quad (6)$$

where k is the discrete time index; $\mathbf{x}(k) \in \mathcal{R}^{k_x}$, $\mathbf{u}(k) \in \mathcal{R}^{k_u}$, and $\mathbf{y}(k) \in \mathcal{R}^{k_y}$ are the state, input, and output vectors, respectively; $\mathbf{w}(k)$ and $\mathbf{v}(k)$ are random noise variables; and $\phi(\cdot)$ and $v(\cdot)$ are continuous nonlinear mappings.

Both the actuator and sensor faults affect (6). Mathematically, the nonlinear system with faults becomes

$$\begin{aligned} \mathbf{x}^f(k+1) &= \phi(\mathbf{x}(k), \mathbf{u}(k), \mathbf{w}(k), \mathbf{f}_a(k)) \\ \mathbf{y}^f(k) &= v(\mathbf{x}(k), \mathbf{u}(k), \mathbf{f}_a(k)) + \mathbf{f}_s(k) + \mathbf{v}(k) \end{aligned} \quad (7)$$

where \mathbf{f}_a and \mathbf{f}_s represent actuator and sensor faults, respectively. The schematic of (7) is given in Fig. 1. Now, we introduce a lemma (given in [30]) that constitutes another foundation of this study.

Lemma 1: An operator \mathcal{K} is called the generalized stable kernel representation for the nonlinear system (6) if, for $\mathbf{w} = \mathbf{0}$ and $\mathbf{v} = \mathbf{0}$, the following relationship holds:

$$\mathcal{K}(z) \begin{bmatrix} \mathbf{u}(z) \\ \mathbf{y}(z) \end{bmatrix} = \mathbf{0} \quad (8)$$

for an initial system state $\mathbf{x}(0)$ and a given \mathbf{u} .

In (8), z refers to discrete variables in the z domain. Motivated by [16], [27], [29], and [31], this study focuses on IFD and its related parameter identification for nonlinear dynamic systems. The three objectives are cast as follows:

- 1) to construct a unified IFD framework that can include both unsupervised and supervised learning-based approaches;
- 2) to design two IFD approaches using the unsupervised and supervised neural networks, respectively;
- 3) to quantify the fault influences in each situation, whose purpose is to mine the interpretability of IFD methods.

D. Revisiting Stable Kernel Representations

As mentioned in [16] for linear systems and [31] for nonlinear systems, $\mathcal{K}(z)$ plays an essential role in constructing an observer and completing IFD tasks. In order to have an intuitive understanding of $\mathcal{K}(z)$, we introduce the form of $\mathcal{K}(z)$ and how to derive it for both linear and nonlinear dynamic systems, as given in the following.

Linear Example: Consider a linear time-invariant system described by

$$\begin{aligned} \mathbf{x}(k+1) &= \mathbf{A}\mathbf{x}(k) + \mathbf{B}\mathbf{u}(k) \\ \mathbf{y}(k) &= \mathbf{C}\mathbf{x}(k) + \mathbf{D}\mathbf{u}(k) \end{aligned} \quad (9)$$

where \mathbf{A} , \mathbf{B} , \mathbf{C} , and \mathbf{D} are the real matrices with appropriate dimensions; other variables are defined in (6). Given a gain matrix \mathbf{L} , a full-order observer has the following dynamic equations:

$$\begin{aligned} \hat{\mathbf{x}}(k+1) &= (\mathbf{A} - \mathbf{L}\mathbf{C})\hat{\mathbf{x}}(k) + [\mathbf{B} - \mathbf{L}\mathbf{D} \quad \mathbf{L}] \begin{bmatrix} \mathbf{u}(k) \\ \mathbf{y}(k) \end{bmatrix} \\ \hat{\mathbf{y}}(k) &= \mathbf{C}\hat{\mathbf{x}}(k) + \mathbf{D}\mathbf{u}(k) \\ \mathbf{r}(k) &= \mathbf{y}(k) - \hat{\mathbf{y}}(k), \\ \mathbf{r}(k) &= -\mathbf{C}\hat{\mathbf{x}}(k) + [-\mathbf{D} \quad \mathbf{I}] \begin{bmatrix} \mathbf{u}(k) \\ \mathbf{y}(k) \end{bmatrix} \rightarrow \mathbf{0} \end{aligned} \quad (10)$$

where $\mathbf{r}(k)$ is the residual signal. Based on the state-space representation of (10), we can define the transfer function from $[\mathbf{u}^T \mathbf{y}^T]^T$ to \mathbf{r} as $\mathcal{K}(z)$ of (9), i.e.,

$$\mathbf{r}(z) = \underbrace{[-\hat{\mathbf{N}}(z) \quad \hat{\mathbf{M}}(z)]}_{\mathcal{K}(z)} \begin{bmatrix} \mathbf{u}(z) \\ \mathbf{y}(z) \end{bmatrix} = \mathbf{0} \quad (11)$$

where $\hat{\mathbf{M}}(z)$ and $\hat{\mathbf{N}}(z)$ are

$$\begin{aligned} \hat{\mathbf{M}}(z) &= \mathbf{I} - \mathbf{C}(z\mathbf{I} - \mathbf{A} + \mathbf{L}\mathbf{C})^{-1}\mathbf{L}, \\ \hat{\mathbf{N}}(z) &= \mathbf{D} + \mathbf{C}(z\mathbf{I} - \mathbf{A} + \mathbf{L}\mathbf{C})^{-1}(\mathbf{B} - \mathbf{L}\mathbf{D}). \end{aligned} \quad (12)$$

Nonlinear Example: In noise-free cases, (6) is redefined by

$$\mathbf{y}(z) = \Pi(z) \circ \mathbf{u}(z) \quad (13)$$

in order to simplify the analysis. Based on the assumption that the nonlinear system $\Pi(z)$ meets the stable condition as in [32], $\Pi(z)$ can be rewritten as

$$\begin{aligned} \Pi(z) &= \Pi_{\mathbf{M}}^{-1}(z) \circ \Pi_{\mathbf{N}}(z) \implies \\ \mathbf{y}(z) &= \Pi_{\mathbf{M}}^{-1}(z) \circ \Pi_{\mathbf{N}}(z)\mathbf{u}(z) \end{aligned} \quad (14)$$

where $\Pi_{\mathbf{M}}$ is an invertible operator and $\Pi_{\mathbf{N}}$ is stable. Then, a generalized $\mathcal{K}(z)$ can define the following residual generator:

$$\mathbf{r}(z) = \underbrace{[-\Pi_{\mathbf{N}}(z) \quad \Pi_{\mathbf{M}}(z)]}_{\mathcal{K}(z)} \begin{bmatrix} \mathbf{u}(z) \\ \mathbf{y}(z) \end{bmatrix} = \mathbf{0}. \quad (15)$$

Note that $\mathcal{K}(z)$ can be obtained based on both system information and input-output data. In addition, the co-inner-outer factorization, having a similar form to (14), can be used to estimate unknown signals, such as external disturbances and unexpected faults [13].

III. UNIFIED IFD FRAMEWORK: FROM SUPERVISED TO UNSUPERVISED LEARNING

By using the composite operators and suspected spaces (see [29]), this section will present a data-driven implementation of residual generators using both unsupervised and supervised learning, with a focus on the construction of the bridge between unsupervised and supervised learning-based approaches.

A. Data-Based Modeling

As both $\mathbf{x}(k)$ and the innovation from $\mathbf{x}(k-1)$ to $\mathbf{x}(k)$ are unknown, the extended form of (6) is usually adopted for system identification and fault diagnosis. For obtaining the extended form, several notations, referring to data and operators, are introduced

$$\begin{aligned} \phi_{\mathbf{xu}}\mathbf{u}(k) &: \mathbf{u}(k) \in \mathcal{R}^{k_u} \rightarrow \mathbf{x}(k+1) \in \mathcal{R}^{k_x} \\ v_{\mathbf{xu}}\mathbf{u}(k) &: \mathbf{u}(k) \in \mathcal{R}^{k_u} \rightarrow \mathbf{x}(k) \in \mathcal{R}^{k_x}, \\ \phi_{\mathbf{xx}}^k &= \phi_{\mathbf{xx}}^{k-1} \circ \phi_{\mathbf{xx}} = \phi_{\mathbf{xx}}^{k-2} \circ \phi_{\mathbf{xx}} \circ \phi_{\mathbf{xx}} \end{aligned} \quad (16)$$

and

$$\mathbf{u}_s(k) = [\mathbf{u}^T(k) \ \cdots \ \mathbf{u}^T(k+s)]^T \in \mathcal{R}^{(s+1)k_u}. \quad (17)$$

In (16), ϕ_{\cdot} and v_{\cdot} are nonlinear operators and can be replaced by any other operator; ϕ_{\cdot}^k represents a high-order composite operator. In (17), s is the stack length, and \mathbf{u} can be replaced by any variable in (6) and (7).

Remark 2: In order to parameterize nonlinear dynamic systems, the composite operators defined in (16), similar to Koopman operators [33], simplify the treatment processes for obtaining an equivalent system representation [34]. ∇

By using the notation given in (16), we rewrite (6) as

$$\mathbf{x}(k+1) = \underbrace{(\phi_{\mathbf{xx}} \ \phi_{\mathbf{xu}} \ \phi_{\mathbf{xw}})}_{\phi} \begin{bmatrix} \mathbf{x}(k) \\ \mathbf{u}(k) \\ \mathbf{w}(k) \end{bmatrix} \quad (18a)$$

$$\mathbf{y}(k) = \underbrace{(v_{\mathbf{yx}} \ v_{\mathbf{yu}} \ v_{\mathbf{yv}})}_v \begin{bmatrix} \mathbf{x}(k) \\ \mathbf{u}(k) \\ \mathbf{v}(k) \end{bmatrix} \quad (18b)$$

where $v_{\mathbf{yv}} = \mathbf{I}$. Similar to the parity space in [2] and [11], the data-based system model can be described as

$$\mathbf{y}_s(k) = (\Upsilon_{\mathbf{x}} \ \Upsilon_{\mathbf{u}}) \begin{bmatrix} \mathbf{x}(k) \\ \mathbf{u}_s(k) \end{bmatrix} + \Upsilon_{\mathbf{w}}\mathbf{w}_s(k) + \mathbf{v}_s(k) \quad (19)$$

in which the composite operators $\Upsilon_{\mathbf{x}}$, $\Upsilon_{\mathbf{u}}$, and $\Upsilon_{\mathbf{w}}$ are

$$\Upsilon_{\mathbf{x}} := \begin{pmatrix} v_{\mathbf{yx}} \\ \vdots \\ v_{\mathbf{yx}}\phi_{\mathbf{xx}}^s \end{pmatrix} : \mathcal{R}^{k_x} \rightarrow \mathcal{R}^{(s+1)k_y} \quad (20)$$

$$\Upsilon_{\mathbf{u}} := \begin{pmatrix} v_{\mathbf{yu}} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ v_{\mathbf{yx}}\phi_{\mathbf{xx}}^{s-1}\phi_{\mathbf{xu}} & \cdots & v_{\mathbf{yu}} \end{pmatrix} : \mathcal{R}^{(s+1)k_u} \rightarrow \mathcal{R}^{(s+1)k_y} \quad (21)$$

$$\Upsilon_{\mathbf{w}} := \begin{pmatrix} \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ v_{\mathbf{yx}}\phi_{\mathbf{xx}}^{s-1}\phi_{\mathbf{xw}} & \cdots & \mathbf{0} \end{pmatrix} : \mathcal{R}^{(s+1)k_x} \rightarrow \mathcal{R}^{(s+1)k_y}. \quad (22)$$

In the data-based model (19), the state $\mathbf{x}(k)$ is generally unknown and needs to be estimated by using past data. The recent studies [27], [35], [36] can bypass estimation of $\mathbf{x}(k)$ and suggest to directly estimate \mathbf{y} through data in the past moving window, i.e.,

$$\mathbf{e}(k) = \mathbf{y}(k) - \hat{\mathbf{y}}(k), \hat{\mathbf{y}}(k) = v(\hat{\mathbf{x}}(k), \mathbf{u}(k)) \quad (23a)$$

$$\hat{\mathbf{x}}(k+1) = \phi(\hat{\mathbf{x}}(k), \mathbf{u}(k)) + \ell(\mathbf{y}(k) - \hat{\mathbf{y}}(k)) \quad (23b)$$

in which $\ell : \mathcal{R}^{k_y} \rightarrow \mathcal{R}^{k_x}$ signifies the (unknown but existing) projection from \mathbf{e} to $\hat{\mathbf{x}}$. Therefore, one can obtain that

$$\begin{aligned} \hat{\mathbf{x}}(k) &= \phi_{\hat{\mathbf{x}}\hat{\mathbf{x}}}^{s_p+1}\hat{\mathbf{x}}(k-s_p-1) + [\mathcal{L}_{\mathbf{p,u}} \ \mathcal{L}_{\mathbf{p,y}}]\mathbf{z}_{\mathbf{p}}(k) \\ &\quad \times \phi_{\hat{\mathbf{x}}\hat{\mathbf{x}}}^{s_p+1}\hat{\mathbf{x}}(k-s_p-1) \approx \mathbf{0} \\ \implies \hat{\mathbf{x}}(k) &\approx \underbrace{[\mathcal{L}_{\mathbf{p,u}} \ \mathcal{L}_{\mathbf{p,y}}]}_{\mathcal{L}_{\mathbf{p}}} \mathbf{z}_{\mathbf{p}}(k) \end{aligned} \quad (24)$$

where $\mathcal{L}_{\mathbf{p,u}}$ and $\mathcal{L}_{\mathbf{p,y}}$ are the composite operators similar to (21) and (22); the past stacked vector $\mathbf{z}_{\mathbf{p}}(k)$ is defined by

$$\begin{aligned} \mathbf{z}_{\mathbf{p}}(k) &= \begin{bmatrix} \mathbf{u}_{\mathbf{p}}(k) \\ \mathbf{y}_{\mathbf{p}}(k) \end{bmatrix}, \mathbf{u}_{\mathbf{p}}(k) = \mathbf{u}_{s_p}(k-s_p-1) \\ \mathbf{y}_{\mathbf{p}}(k) &= \mathbf{y}_{s_p}(k-s_p-1). \end{aligned} \quad (25)$$

Combining (19) with (24) yields an equivalent model

$$\begin{aligned} \mathbf{y}_s(k) &= \Upsilon_{\mathbf{x}}\mathcal{L}_{\mathbf{p}}\mathbf{z}_{\mathbf{p}}(k) + \Upsilon_{\mathbf{u}}\mathbf{u}_s(k) + \Upsilon_{\mathbf{e}}\mathbf{e}_s(k) \\ &= (\Upsilon_{\mathbf{x}}\mathcal{L}_{\mathbf{p}} \ \Upsilon_{\mathbf{u}}) \begin{bmatrix} \mathbf{z}_{\mathbf{p}}(k) \\ \mathbf{u}_s(k) \end{bmatrix} + \Upsilon_{\mathbf{e}}\mathbf{e}_s(k) \end{aligned} \quad (26)$$

where \mathbf{e}_s contains the influences caused by \mathbf{w} and \mathbf{v} , and $\Upsilon_{\mathbf{e}}$ has the following form:

$$\Upsilon_{\mathbf{e}} = \begin{pmatrix} \mathbf{I} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ v_{\hat{\mathbf{y}}\hat{\mathbf{x}}}\phi_{\hat{\mathbf{x}}\hat{\mathbf{x}}}^{s-1}\ell_{\hat{\mathbf{x}}\mathbf{e}} & \cdots & \mathbf{I} \end{pmatrix}. \quad (27)$$

Remark 3: By the use of the composite operator, (26) details a generalized nonlinear predictor from $\mathbf{z}_{\mathbf{p}}$ and \mathbf{u}_s to the current system output \mathbf{y}_s while resembling a linear model. It makes system modeling possible and explainable, especially achieving a consensus among different IFD approaches using both unsupervised and supervised learning. ∇

B. Data-Driven Implementation of Residual Generators

On the basis of Lemma 1, the following corollary is obtained, whose embryonic development credits to [27].

Corollary 1: An operator \mathcal{K}_{s_p+s} is called the data-driven generalized kernel representation for (6) if, for $\mathbf{w} = \mathbf{0}$ and $\mathbf{v} = \mathbf{0}$, the following relationship holds:

$$\underbrace{(\mathcal{K}_{\mathbf{z}}^{\text{unsp}} \ \mathcal{K}_{\mathbf{u}}^{\text{unsp}} \ \mathcal{K}_{\mathbf{y}}^{\text{unsp}})}_{\mathcal{K}_{s_p+s}^{\text{unsp}}} \begin{bmatrix} \mathbf{z}_{\mathbf{p}} \\ \mathbf{u}_s \\ \mathbf{y}_s \end{bmatrix} = \mathbf{0} \quad (28a)$$

$$\text{or } \underbrace{(\mathcal{K}_{s_p+s}^{\text{sp}} - \mathbf{I})}_{\mathcal{K}_{s_p+s}^{\text{sp}}} \begin{bmatrix} \mathbf{z}_{\mathbf{p}} \\ \mathbf{u}_s \\ \mathbf{y}_s \end{bmatrix} = \mathbf{0} \quad (28b)$$

for an initial system state $\mathbf{x}(0)$ and an arbitrary \mathbf{u}_s , where the superscripts “unsp” and “sp” signify the unsupervised and supervised learning scenarios, respectively.

The obtained \mathcal{K}_{s_p+s} , by using the unsupervised or supervised learning methods, can be directly applied in constructing a unified residual generator such that

$$\mathbf{r}_s(k) = \mathcal{K}_{s_p+s} \begin{bmatrix} \mathbf{z}_p(k) \\ \mathbf{u}_s(k) \\ \mathbf{y}_s(k) \end{bmatrix} \in \mathcal{R}^{(s+1)k_y} \quad (29)$$

in fault-free conditions; \mathcal{K}_{s_p+s} can be replaced by $\mathcal{K}_{s_p+s}^{\text{unsp}}$ and $\mathcal{K}_{s_p+s}^{\text{sp}}$ provided in Corollary 1. When (6) is affected by faults, $\mathbf{r}_s(k)$ given in (29) becomes

$$\mathbf{r}_s^f(k) = \mathbf{r}_s(k) + \mathbf{f}_{a,\text{term}} + \mathbf{f}_{s,\text{term}} \quad (30)$$

where $\mathbf{f}_{a,\text{term}}$ and $\mathbf{f}_{s,\text{term}}$ are the actuator and sensor fault-related terms, respectively; their forms can be obtained via composite operators. Then, the test statistics [1], such as T^2 , can be defined on the residual signal \mathbf{r}_s , i.e.,

$$T^2(\mathbf{r}_s(k)) = \mathbf{r}_s^T(k) \Sigma_{\mathbf{r}_s}^{-1} \mathbf{r}_s(k), \quad \Sigma_{\mathbf{r}_s} = \mathbb{E}(\mathbf{r}_s(k) \mathbf{r}_s^T(k)). \quad (31)$$

Correspondingly, the threshold is determined by the chosen confidence level α

$$J_{th} \leftarrow \Pr(T^2(\mathbf{r}_s(k)) > J_{th}) = \alpha. \quad (32)$$

The following theorem presents the essence of this study, sketching *the bridge* (i.e., *one-to-one mapping*) between unsupervised and supervised learning-based parameter identification, together with the IFD approaches, for nonlinear systems.

Theorem 1: Consider the two residual generators defined by

$$\mathbf{r}_s^{\text{unsp}}(k) = \mathcal{K}_{s_p+s}^{\text{unsp}} \begin{bmatrix} \mathbf{z}_p(k) \\ \mathbf{u}_s(k) \\ \mathbf{y}_s(k) \end{bmatrix}, \quad \mathbf{r}_s^{\text{sp}}(k) = \mathcal{K}_{s_p+s}^{\text{sp}} \begin{bmatrix} \mathbf{z}_p(k) \\ \mathbf{u}_s(k) \\ \mathbf{y}_s(k) \end{bmatrix} \quad (33)$$

where $\mathcal{K}_{s_p+s}^{\text{unsp}}$ and $\mathcal{K}_{s_p+s}^{\text{sp}}$ are defined by

$$\mathcal{K}_{s_p+s}^{\text{unsp}} = [\mathbf{0} \quad \mathbf{0} \quad \mathbf{I}] - \mathcal{U}_y, \quad (34a)$$

$$\mathcal{K}_{s_p+s}^{\text{sp}} = (-\Upsilon_x \mathcal{L}_p \quad -\Upsilon_u \quad \mathbf{I}) \quad (34b)$$

in which \mathcal{U}_y is a segment of unsupervised learning \mathcal{U}

$$\begin{bmatrix} \hat{\mathbf{z}}_p(k) \\ \hat{\mathbf{u}}_s(k) \\ \hat{\mathbf{y}}_s(k) \end{bmatrix} = \underbrace{(\mathcal{U}_z \quad \mathcal{U}_u \quad \mathcal{U}_y)}_{\mathcal{U}} \begin{bmatrix} \mathbf{z}_p(k) \\ \mathbf{u}_s(k) \\ \mathbf{y}_s(k) \end{bmatrix}. \quad (35)$$

There exists a nonlinear operator $\mathcal{P}_{\text{sp/unsp}}$ such that

$$\mathbf{r}_s^{\text{sp}}(k) = \mathcal{P}_{\text{sp/unsp}} \mathbf{r}_s^{\text{unsp}}(k) \quad (36)$$

and the inverse function of $\mathcal{P}_{\text{sp/unsp}}$, denoted as $\mathcal{P}_{\text{unsp/sp}} = \mathcal{P}_{\text{sp/unsp}}^{-1}$, must exist such that

$$\mathbf{r}_s^{\text{unsp}}(k) = \mathcal{P}_{\text{unsp/sp}} \mathbf{r}_s^{\text{sp}}(k). \quad (37)$$

Proof: The complete proof is given in the Appendix. ■

The following remark is made to set forth contributions of Theorem 1.

Remark 4: With the help of the bridge, Theorem 1 provides us with an elegant way to evaluate the performance of both

unsupervised and supervised learning-based IFD approaches in the sense of fault-detection capacity. ▽

Based on the metric $\mathcal{M}_{\text{euc}}^2$, the suspected space is used to design IFD approaches, whose purpose is to gain a more in-depth understanding of Theorem 1.

C. IFD Using Unsupervised Learning

In order to develop the unsupervised learning-based IFD approaches, a suspected space, denoted as $\mathcal{Y}^{\text{unsp}}$, is defined according to \mathcal{U}_y in the following definition.

Definition 1: Given any \mathcal{U} in (35), $\mathcal{Y}^{\text{unsp}}$ obtained by

$$\mathcal{Y}^{\text{unsp}} = \mathcal{U}_y [\mathbf{z}_p^T \quad \mathbf{u}_s^T \quad \mathbf{y}_s^T]^T \quad (38)$$

spans a space $\mathcal{Y}^{\text{unsp}}$ that is called the suspected space of (6).

Corresponding to $\mathcal{Y}^{\text{unsp}}$, \mathcal{Y} is called the measurement space, where $\mathbf{y}_s \in \mathcal{Y} \subset \mathcal{R}^{(s+1)k_y}$. Based on the metric defined in (1), the objective function of IFD approaches using unsupervised learning can be formulated as

$$\min \mathcal{M}_{\text{euc}}(\mathbf{y}_s, \mathcal{Y}^{\text{unsp}}) \quad (39a)$$

$$\Rightarrow \mathcal{U}_y^* := \arg \min_{\mathcal{U}_y} \mathcal{M}_{\text{euc}}(\mathbf{y}_s, \mathcal{Y}^{\text{unsp}}) \quad (39b)$$

where the superscript “*” signifies the best solution. It is of interest to find that the residual space, denoted as $\mathcal{E}_s^{\text{unsp}}$, can be obtained based on \mathcal{U}_y^*

$$\begin{aligned} \mathbf{e}_s &= \mathbf{y}_s - \mathcal{Y}^{\text{unsp},*}, \quad \mathbf{e}_s \in \mathcal{E}_s^{\text{unsp}} \subset \mathcal{R}^{(s+1)k_y} \\ \mathcal{Y}^{\text{unsp},*} &= \mathcal{U}_y^* [\mathbf{z}_p^T \quad \mathbf{u}_s^T \quad \mathbf{y}_s^T]^T \in \mathcal{Y}^{\text{unsp},*} \end{aligned} \quad (40)$$

which satisfies

$$\mathcal{E}_s^{\text{unsp}} \perp \mathcal{Y}^{\text{unsp},*}, \quad \mathcal{E}_s^{\text{unsp}} (\approx) \perp \mathcal{Y}. \quad (41)$$

Theorem 2: Considering a nonlinear system (6), its generalized kernel representation is defined by $\mathcal{K}_{s_p+s}^{\text{unsp}} = [\mathbf{0} \quad \mathbf{0} \quad \mathbf{I}] - \mathcal{U}_y^$, where \mathcal{U}_y^* is obtained through (39b). For the faults $\mathbf{f}_a(k)$ and $\mathbf{f}_s(k)$ occurring from the k th time instant, $\mathbf{r}_s^{\text{unsp}}(k)$ given in (33) becomes*

$$\begin{aligned} \mathbf{r}_s^{\text{unsp},f}(k) &= \mathbf{e}_s(k) + \mathbf{f}_{a,\text{term}}^{\text{unsp}}(k) + \mathbf{f}_{s,\text{term}}^{\text{unsp}}(k) \\ \mathbf{f}_{a,\text{term}}^{\text{unsp}}(k) &= \Upsilon_{f_a} \mathbf{f}_{a,s}(k) \\ \mathbf{f}_{s,\text{term}}^{\text{unsp}}(k) &= \Upsilon_{f_s} \mathbf{f}_{s,s}(k) \end{aligned} \quad (42)$$

where $\mathbf{f}_{a,\text{term}}^{\text{unsp}}$ and $\mathbf{f}_{s,\text{term}}^{\text{unsp}}$ have the following forms:

$$\Upsilon_{f_a} = \Upsilon_{f_u}, \quad \Upsilon_{f_s} = \begin{pmatrix} \mathbf{I} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ -v_{\hat{\mathbf{y}}\hat{\mathbf{x}}} \mathcal{A}_\ell^{s-1} \ell_{\hat{\mathbf{x}}\mathbf{e}} & \cdots & \mathbf{I} \end{pmatrix}. \quad (43)$$

Then, T^2 defined on $\mathbf{r}_s^{\text{unsp},f}(k)$ according to

$$T^2(\mathbf{r}_s^{\text{unsp},f}(k)) = \mathbf{r}_s^{\text{unsp},f,T}(k) \Sigma_{\mathbf{r}_s}^{-1} \mathbf{r}_s^{\text{unsp},f}(k) \quad (44)$$

has the optimal fault-detection power.

Proof: When (6) is fault-free, we have

$$\mathbf{r}_s^{\text{unsp}}(k) = \mathbf{e}_s(k) \quad \Sigma_{\mathbf{r}_s^{\text{unsp}}} = \mathbb{E}(\mathbf{e}_s(k) \mathbf{e}_s^T(k)). \quad (45)$$

Since \mathbf{f}_a and \mathbf{f}_s , respectively, represent the actuator and sensor faults, it is reasonable to adopt

$$\phi_{x\mathbf{f}_a} = \phi_{x\mathbf{u}}, \quad v_{y\mathbf{f}_a} = v_{y\mathbf{u}}, \quad \text{and} \quad v_{y\mathbf{f}_s} = \mathbf{I}. \quad (46)$$

When (6) is affected by \mathbf{f}_a and \mathbf{f}_s , $\hat{\mathbf{x}}^f(k)$ and $\mathbf{y}^f(k)$ in (23a) and (23b) can be expressed by

$$\begin{aligned} \hat{\mathbf{x}}^f(k+1) &= \phi_{\hat{x}\hat{x}}\hat{\mathbf{x}}(k) + \phi_{\hat{x}\mathbf{u}}\mathbf{u}(k) + \phi_{\hat{x}\mathbf{f}_a}\mathbf{f}_a(k) \\ &\quad + \ell(\mathbf{y}^f(k) - \hat{\mathbf{y}}(k) - \mathbf{f}_s)\mathbf{y}^f(k) \\ &= v_{\hat{y}\hat{x}}\hat{\mathbf{x}}(k) + v_{\hat{y}\mathbf{u}}\mathbf{u}(k) + v_{\hat{y}\mathbf{f}_a}\mathbf{f}_a(k) \\ &\quad + v_{y\mathbf{f}_s}\mathbf{f}_s(k) + \mathbf{e}(k). \end{aligned} \quad (47)$$

Combining it with (A.2) yields (42) with $\underline{\Upsilon}_{\mathbf{f}_a}$ and $\underline{\Upsilon}_{\mathbf{f}_s}$ described by (43). Furthermore, taking expectation of (44) obtains

$$\begin{aligned} \mathbb{E}(T^2(\mathbf{r}_s^{\text{unsp},f}(k))) &= T^2(\mathbf{e}_s) + \mathbf{f}_{a,\text{term}}^T \Sigma_{\mathbf{r}_s^{\text{unsp}}}^{-1} \mathbf{f}_{a,\text{term}} \\ &\quad + \mathbf{f}_{f,\text{term}}^T \Sigma_{\mathbf{r}_s^{\text{unsp}}}^{-1} \mathbf{f}_{f,\text{term}} \end{aligned} \quad (48)$$

because, in general, $\mathbf{e}_s(k)$, $\mathbf{f}_{a,\text{term}}(k)$, and $\mathbf{f}_{f,\text{term}}(k)$ are independent. As proved in [29], \mathcal{U}_y^* minimizes $\Sigma_{\mathbf{r}_s^{\text{unsp}}}$ via (39b); thus, the last two terms in (48) are maximized to the greatest extent possible. It means that \mathcal{U}_y^* , a segment of \mathcal{U}^* , can achieve optimal performance in detecting \mathbf{f}_a and \mathbf{f}_s .

The theorem is proven. \blacksquare

D. IFD Using Supervised Learning

In parallel to $\mathcal{Y}^{\text{unsp}}$, another suspected space \mathcal{Y}^{sp} using supervised learning can be defined according to \mathcal{S} as follows.

Definition 2: Given any \mathcal{S} in (A.8), $\underline{\mathbf{y}}_s^{\text{sp}}$ obtained by

$$\underline{\mathbf{y}}_s^{\text{sp}} = \mathcal{S}[\mathbf{z}_p^T \quad \mathbf{u}_s^T]^T \quad (49)$$

spans a space \mathcal{Y}^{sp} that is also the suspected space of (6), where \mathcal{S} is defined by

$$\mathcal{S} = (\mathcal{S}_z \quad \mathcal{S}_u) := (\Upsilon_x \mathcal{L}_p \quad \Upsilon_u). \quad (50)$$

Similar to (39a) and (39b), the optimal \mathcal{S}^* can be obtained via

$$\min \mathcal{M}_{\text{euc}}(\mathbf{y}_s, \underline{\mathbf{y}}_s^{\text{sp}}) \quad (51a)$$

$$\implies \mathcal{S}^* := \arg \min_{\mathcal{S}} \mathcal{M}_{\text{euc}}(\mathbf{y}_s, \underline{\mathbf{y}}_s^{\text{sp}}) \quad (51b)$$

based on which $\mathcal{Y}^{\text{sp},*}$ can be obtained.

By the use of \mathcal{S}^* , the following theorem is derived for IFD using supervised learning.

Theorem 3: Considering a nonlinear system (6), its generalized kernel representation is defined by $\mathcal{K}_{s_p+s}^{\text{sp}} = (\mathcal{S}^* \quad -\mathbf{I})$, where \mathcal{S}^* is obtained through (51b). For the faults $\mathbf{f}_a(k)$ and $\mathbf{f}_s(k)$ occurring from the k th time instant, $\mathbf{r}_s^{\text{sp}}(k)$ given in (33) becomes

$$\begin{aligned} \mathbf{r}_s^{\text{sp},f}(k) &= \Upsilon_e \mathbf{e}_s(k) + \mathbf{f}_{a,\text{term}}^{\text{sp}}(k) + \mathbf{f}_{s,\text{term}}^{\text{sp}}(k) \\ \mathbf{f}_{a,\text{term}}^{\text{sp}}(k) &= \Upsilon_{\mathbf{f}_a} \mathbf{f}_{a,s}(k) \\ \mathbf{f}_{s,\text{term}}^{\text{sp}}(k) &= \Upsilon_{\mathbf{f}_s} \mathbf{f}_{s,s} \end{aligned} \quad (52)$$

where $\mathbf{f}_{a,\text{term}}^{\text{sp}}$ and $\mathbf{f}_{s,\text{term}}^{\text{sp}}$ have the following forms:

$$\Upsilon_{\mathbf{f}_a} = \Upsilon_u \Upsilon_{\mathbf{f}_s} = \mathbf{I}. \quad (53)$$

Then, T^2 defined on $\mathbf{r}_s^{\text{sp},f}(k)$ according to

$$T^2(\mathbf{r}_s^{\text{sp},f}(k)) = \mathbf{r}_s^{\text{sp},f,T}(k) \Sigma_{\mathbf{r}_s^{\text{sp},f}}^{-1} \mathbf{r}_s^{\text{sp},f}(k) \quad (54)$$

has the optimal fault-detection power.

Proof: The proof is similar to Theorem 2. Owing to space constraints, the detail is omitted here.

In addition, another sketch of the proof for the theorem is presented as follows. By using the bridge $\mathcal{P}_{\text{sp}/\text{unsp}}$ given in Theorem 1, simple mathematical manipulations can yield

$$T^2(\mathbf{r}_s^{\text{sp}}(k)) = T^2(\mathbf{r}_s^{\text{unsp}}(k)) \quad (55)$$

and (53). Furthermore, it can also be verified

$$T^2(\mathbf{r}_s^{\text{sp},f}(k)) = T^2(\mathbf{r}_s^{\text{unsp},f}(k)) \quad (56)$$

which has the same (i.e., optimal) performance of fault detection as $T^2(\mathbf{r}_s^{\text{unsp},f}(k))$ and, thus, completes this proof. \blacksquare

Remark 5: As presented in Theorems 2 and 3, performance evaluation of the two proposed IFD frameworks for nonlinear systems becomes possible. The main reason is that the difference between $\mathbf{y}_s(k)$ and its estimation in $\mathcal{Y}^{\text{unsp},*}$ and $\mathcal{Y}^{\text{sp},*}$ can be measured through quantitative metrics. ∇

E. Notes on the Bridge

In what follows, several remarks (including fault features and geometric interpretation) and perspectives (i.e., a more general version) are made to set forth contributions and essences of the bridge provided in Theorem 1.

1) *Fault Features:* Along with Theorems 2 and 3, the fault features, corresponding to $\mathcal{Y}^{\text{unsp},*}$ and $\mathcal{Y}^{\text{sp},*}$, are

$$\mathbf{r}_s^{\text{unsp},f} = \mathbf{e}_s + \mathbf{f}_{a,\text{term}}^{\text{unsp}} + \mathbf{f}_{s,\text{term}}^{\text{unsp}}, \quad (57a)$$

$$\mathbf{r}_s^{\text{sp},f} = \mathcal{P}_{\text{sp}/\text{unsp}} \mathbf{e}_s + \mathbf{f}_{a,\text{term}}^{\text{sp}} + \mathbf{f}_{s,\text{term}}^{\text{sp}}. \quad (57b)$$

It is interesting to verify

$$\begin{aligned} \mathcal{P}_{\text{sp}/\text{unsp}} \mathbf{f}_{a,\text{term}}^{\text{unsp}} &= \mathbf{f}_{a,\text{term}}^{\text{sp}} \\ \mathcal{P}_{\text{unsp}/\text{sp}} \mathbf{f}_{a,\text{term}}^{\text{sp}} &= \mathcal{P}_{\text{sp}/\text{unsp}}^{-1} \mathbf{f}_{a,\text{term}}^{\text{sp}} = \mathbf{f}_{a,\text{term}}^{\text{unsp}} \end{aligned} \quad (58)$$

for $\mathbf{f}_{a,\text{term}}$, and

$$\mathcal{P}_{\text{sp}/\text{unsp}} \mathbf{f}_{s,\text{term}}^{\text{unsp}} = \mathbf{f}_{s,\text{term}}^{\text{sp}}, \quad \mathcal{P}_{\text{unsp}/\text{sp}} \mathbf{f}_{s,\text{term}}^{\text{sp}} = \mathbf{f}_{s,\text{term}}^{\text{unsp}} \quad (59)$$

for $\mathbf{f}_{s,\text{term}}$. Combining (57a) and (57b) with (58) and (59) yields

$$\mathcal{P}_{\text{sp}/\text{unsp}} \mathbf{r}_s^{\text{unsp},f} = \mathbf{r}_s^{\text{sp},f} = \mathcal{P}_{\text{sp}/\text{unsp}} \circ \mathcal{P}_{\text{sp}/\text{unsp}}^{-1} \mathbf{r}_s^{\text{sp},f} \quad (60)$$

which indicates that, based on the bridge given in Theorem 1, the fault features can be transformed between each other. Also, (60) can be used as an auxiliary evidential statement to illustrate the validity of Theorems 2 and 3 because of

$$\begin{aligned} \Sigma_{\mathbf{r}_s^{\text{unsp}}} &= \mathbb{E}(\mathbf{e}_s(k) \mathbf{e}_s^T(k)), \\ \Sigma_{\mathbf{r}_s^{\text{sp}}} &= \mathcal{P}_{\text{sp}/\text{unsp}} \Sigma_{\mathbf{r}_s^{\text{unsp}}} \mathcal{P}_{\text{sp}/\text{unsp}}^T \end{aligned} \quad (61)$$

such that

$$\begin{aligned} \mathbb{E}(T^2(\mathbf{r}_s^{\text{unsp}}(k))) &= \mathbb{E}(T^2(\mathbf{r}_s^{\text{sp}}(k))) \\ \mathbb{E}(T^2(\mathbf{r}_s^{\text{unsp},f}(k))) &= \mathbb{E}(T^2(\mathbf{r}_s^{\text{sp},f}(k))). \end{aligned} \quad (62)$$

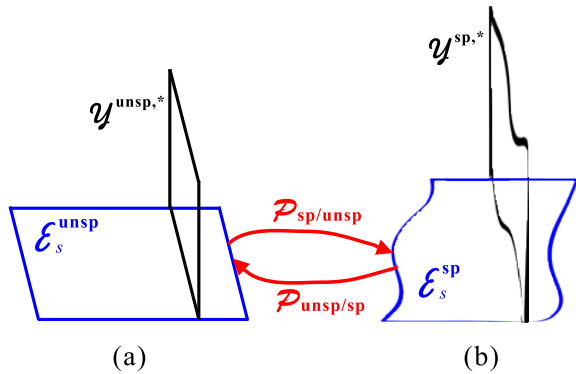


Fig. 2. Geometric descriptions of the suspected and residual spaces in system identification. (a) Unsupervised learning. (b) Supervised learning.

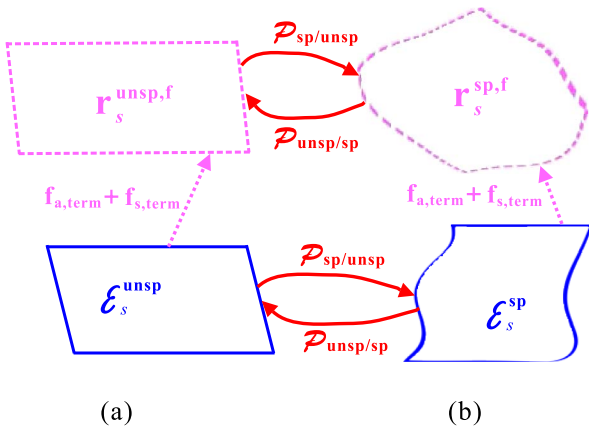


Fig. 3. Geometric descriptions of the residual spaces in both normal and faulty scenarios. (a) Unsupervised learning. (b) Supervised learning.

2) *Geometric Interpretation*: Fig. 2 provides the geometric descriptions of both unsupervised and supervised learning-aided parameter identification, where the bridge, $\mathcal{P}_{\text{unsp/sp}}$ and $\mathcal{P}_{\text{sp/unsp}}$, is highlighted by the two red curves. It is worth mentioning that $\mathbf{y}^{\text{unsp},*}$ and $\mathbf{y}^{\text{sp},*}$ reconstruct the dynamic behaviors of nonlinear systems by carrying the maximum variation information (i.e., the largest uncertainty). They are suitable for nonlinear parameter identification. At the same time, $\mathcal{E}^{\text{unsp},*}$ and $\mathcal{E}^{\text{sp},*}$ with the minimum uncertainty are the best choices for fault diagnosis [16].

Fig. 3 sketches the changes caused by \mathbf{f}_a and \mathbf{f}_s in residual spaces, together with the same bridge in both normal and faulty scenarios.

3) *Generalized Version*: In Theorems 1–3, \mathcal{M}_{euc} is chosen as the metric to measure the difference not only between \mathbf{y}_s and $\mathbf{y}_s^{\text{unsp}}$ when utilizing unsupervised learning to identify $\mathcal{K}_{s_p+s}^{\text{unsp}}$ but also between \mathbf{y}_s and \mathbf{y}_s^{sp} when utilizing supervised learning to estimate $\mathcal{K}_{s_p+s}^{\text{sp}}$. For the fault diagnosis purpose, the same performance can be obtained by using \mathcal{M}_{cor} defined in (4a) and (4b). As pointed in [29] and [31], an IFD design relying on the \mathcal{M}_{cor} metric is a generalized version that can provide multiple optimal solutions to both $\mathcal{K}_{s_p+s}^{\text{unsp}}$ and $\mathcal{K}_{s_p+s}^{\text{sp}}$.

In order to distinguish it from the solution based on \mathcal{M}_{euc} , the notations in the generalized version are marked by “ $\underline{\cdot}$ ”. Therefore, the objective functions of two generalized versions,

respectively, corresponding to unsupervised and supervised learning can be formulated as follows:

$$\begin{aligned} & \min \mathcal{M}_{\text{cor}}(\mathbf{y}_s, \underline{\mathbf{y}}_s^{\text{unsp}}) \\ & = \min k_{y_s} - \text{Tr} \left(\left| \Sigma_{y_s}^{-1/2} \Sigma_{\underline{\mathbf{y}}_s^{\text{unsp}}} \Sigma_{\underline{\mathbf{y}}_s^{\text{unsp}}}^{-1/2} \right| \right) \\ & \iff \min k_{y_s} - \text{Tr}(\Sigma^{\text{unsp},T} \Sigma^{\text{unsp}}) \end{aligned} \quad (63)$$

and

$$\begin{aligned} & \min \mathcal{M}_{\text{cor}}(\mathbf{y}_s, \underline{\mathbf{y}}_s^{\text{sp}}) \\ & = \min k_{y_s} - \text{Tr} \left(\left| \Sigma_{y_s}^{-1/2} \Sigma_{\underline{\mathbf{y}}_s^{\text{sp}}} \Sigma_{\underline{\mathbf{y}}_s^{\text{sp}}}^{-1/2} \right| \right) \\ & \iff \min k_{y_s} - \text{Tr}(\Sigma^{\text{sp},T} \Sigma^{\text{sp}}) \end{aligned} \quad (64)$$

where Σ^{unsp} and Σ^{sp} are obtained via the following singular value decompositions:

$$\begin{aligned} \Sigma_{y_s}^{-1/2} \Sigma_{\underline{\mathbf{y}}_s^{\text{unsp}}} \Sigma_{\underline{\mathbf{y}}_s^{\text{unsp}}}^{-1/2} &= \underline{\Gamma}^{\text{unsp}} \underline{\Sigma}^{\text{unsp}} \underline{\Upsilon}^{\text{unsp},T} \\ \Sigma_{y_s}^{-1/2} \Sigma_{\underline{\mathbf{y}}_s^{\text{sp}}} \Sigma_{\underline{\mathbf{y}}_s^{\text{sp}}}^{-1/2} &= \underline{\Gamma}^{\text{sp}} \underline{\Sigma}^{\text{sp}} \underline{\Upsilon}^{\text{sp},T}. \end{aligned} \quad (65)$$

Now, we define two sets of linear mappings as follows:

$$\underline{\mathbf{G}}_{y_s}^{\text{unsp}} = \underline{\Gamma}^{\text{unsp},T} \Sigma_{y_s}^{-1/2}, \quad \underline{\mathbf{G}}_{\underline{\mathbf{y}}_s^{\text{unsp}}}^{\text{unsp}} = \underline{\Upsilon}^{\text{unsp},T} \Sigma_{\underline{\mathbf{y}}_s^{\text{unsp}}}^{-1/2}, \quad (66a)$$

$$\underline{\mathbf{G}}_{y_s}^{\text{sp}} = \underline{\Gamma}^{\text{sp},T} \Sigma_{y_s}^{-1/2}, \quad \underline{\mathbf{G}}_{\underline{\mathbf{y}}_s^{\text{sp}}}^{\text{sp}} = \underline{\Upsilon}^{\text{sp},T} \Sigma_{\underline{\mathbf{y}}_s^{\text{sp}}}^{-1/2}. \quad (66b)$$

Then, the two generalized residual generators can be constructed according to

$$\underline{\mathbf{r}}_s^{\text{unsp}} = \underline{\mathbf{G}}_{y_s}^{\text{unsp}} \mathbf{y}_s - \underline{\Sigma}^{\text{unsp}} \underline{\mathbf{G}}_{\underline{\mathbf{y}}_s^{\text{unsp}}}^{\text{unsp}} \underline{\mathbf{y}}_s^{\text{unsp}}, \quad (67a)$$

$$\underline{\mathbf{r}}_s^{\text{sp}} = \underline{\mathbf{G}}_{y_s}^{\text{sp}} \mathbf{y}_s - \underline{\Sigma}^{\text{sp}} \underline{\mathbf{G}}_{\underline{\mathbf{y}}_s^{\text{sp}}}^{\text{sp}} \underline{\mathbf{y}}_s^{\text{sp}}. \quad (67b)$$

It can be verified that

$$\Sigma_{\underline{\mathbf{r}}_s^{\text{unsp}}} = \mathbf{I} - \underline{\Sigma}^{\text{unsp}} \underline{\Sigma}^{\text{unsp},T}, \quad \Sigma_{\underline{\mathbf{r}}_s^{\text{sp}}} = \mathbf{I} - \underline{\Sigma}^{\text{sp}} \underline{\Sigma}^{\text{sp},T}. \quad (68)$$

Furthermore, by substituting (66a) into (67a), one can obtain

$$\begin{aligned} & \underline{\mathbf{G}}_{y_s}^{\text{unsp},-1} \underline{\mathbf{r}}_s^{\text{unsp}} \\ & = \mathbf{y}_s - \underline{\mathbf{G}}_{y_s}^{\text{unsp},-1} \underline{\Sigma}^{\text{unsp}} \underline{\mathbf{G}}_{\underline{\mathbf{y}}_s^{\text{unsp}}}^{\text{unsp}} \underline{\mathbf{y}}_s^{\text{unsp}} \\ & = \mathbf{y}_s - \Sigma_{y_s}^{1/2} \left(\underline{\Gamma}^{\text{unsp}} \underline{\Sigma}^{\text{unsp}} \underline{\Upsilon}^{\text{unsp},T} \right) \Sigma_{\underline{\mathbf{y}}_s^{\text{unsp}}}^{-1/2} \underline{\mathbf{y}}_s^{\text{unsp}} \\ & = \mathbf{y}_s - \Sigma_{\underline{\mathbf{y}}_s^{\text{unsp}}} \Sigma_{\underline{\mathbf{y}}_s^{\text{unsp}}}^{-1} \underline{\mathbf{y}}_s^{\text{unsp}} \end{aligned} \quad (69)$$

because $\underline{\mathbf{G}}_{y_s}^{\text{unsp}}$ has full rank. Define

$$\underline{\mathbf{y}}_s^{\text{unsp}} = \underline{\mathcal{U}} \begin{bmatrix} \mathbf{z}_p \\ \mathbf{u}_s \\ \mathbf{y}_s \end{bmatrix} \quad (70)$$

where its optimal choices, $\underline{\mathbf{y}}_s^{\text{unsp},*}$ and $\underline{\mathcal{U}}^*$, are obtained according to (63). It must be pointed out that $\underline{\mathcal{U}}^*$ and its corresponding output $\underline{\mathbf{y}}_s^{\text{unsp},*}$ are not unique because the objective functions are defined by \mathcal{M}_{euc} . Then, (69) can be further

written as

$$\begin{aligned} \underline{\underline{\mathbf{G}}}_{y_s}^{\text{unsp},-1} \underline{\underline{\mathbf{r}}}_s^{\text{unsp}} &= \mathbf{y}_s - \underbrace{\Sigma_{y_s, y_s}^{\text{unsp}} \Sigma_{y_s}^{-1} \underline{\underline{\mathcal{L}}}_{y_s}^*}_{\mathcal{U}_y^*} \begin{bmatrix} \mathbf{z}_p \\ \mathbf{u}_s \\ \mathbf{y}_s \end{bmatrix} \\ &= \mathbf{e}_s = \mathbf{r}_s^{\text{unsp}} \\ \iff \underline{\underline{\mathbf{r}}}_s^{\text{unsp}} &= \underline{\underline{\mathbf{G}}}_{y_s}^{\text{unsp}} \mathbf{r}_s^{\text{unsp}}. \end{aligned} \quad (71)$$

Similarly, the following relationship holds for supervised learning-based IFD approaches:

$$\underline{\underline{\mathbf{r}}}_s^{\text{sp}} = \underline{\underline{\mathbf{G}}}_{y_s}^{\text{sp}} \mathbf{r}_s^{\text{sp}}. \quad (72)$$

Eventually, the following two relationships hold:

$$\begin{aligned} T^2(\underline{\underline{\mathbf{r}}}_s^{\text{unsp}}(k)) &= \underline{\underline{\mathbf{r}}}_s^{\text{unsp},T}(k) \Sigma_{\underline{\underline{\mathbf{r}}}_s^{\text{unsp}}}^{-1} \underline{\underline{\mathbf{r}}}_s^{\text{unsp}}(k) \\ &= (\underline{\underline{\mathbf{G}}}_{y_s}^{\text{unsp}} \mathbf{r}_s^{\text{unsp}})^T \Sigma_{\underline{\underline{\mathbf{r}}}_s^{\text{unsp}}}^{-1} \underline{\underline{\mathbf{G}}}_{y_s}^{\text{unsp}} \mathbf{r}_s^{\text{unsp}} \\ &= T^2(\mathbf{r}_s^{\text{unsp}}(k)) \end{aligned} \quad (73)$$

and

$$T^2(\underline{\underline{\mathbf{r}}}_s^{\text{sp}}(k)) = T^2(\mathbf{r}_s^{\text{sp}}(k)). \quad (74)$$

Combining (74) with (55) obtains

$$\begin{aligned} T^2(\underline{\underline{\mathbf{r}}}_s^{\text{sp}}(k)) &= T^2(\underline{\underline{\mathbf{r}}}_s^{\text{unsp}}(k)) = \\ T^2(\mathbf{r}_s^{\text{sp}}(k)) &= T^2(\mathbf{r}_s^{\text{unsp}}(k)) \end{aligned} \quad (75)$$

for normal operations, and

$$\begin{aligned} T^2(\underline{\underline{\mathbf{r}}}_s^{\text{sp},f}(k)) &= T^2(\underline{\underline{\mathbf{r}}}_s^{\text{unsp},f}(k)) = \\ T^2(\mathbf{r}_s^{\text{sp},f}(k)) &= T^2(\mathbf{r}_s^{\text{unsp},f}(k)) \end{aligned} \quad (76)$$

for faulty operations.

Following Theorems 2 and 3, T^2 test statistic defined on the two generalized residual generators also has the optimal fault-detection power. The readers can refer to [29] for a more rigorous analysis of the generalized version.

Remark 6: We can find that the generalized versions of residual generators also deliver the least-squares estimation of \mathbf{y}_s because both $\underline{\underline{\mathbf{G}}}_{y_s}^{\text{unsp}}$ and $\underline{\underline{\mathbf{G}}}_{y_s}^{\text{sp}}$ play a role as normalization by limiting their covariance matrices to (68). ∇

In order to have an insightful observation, Fig. 4 depicts multiple solutions to the suspected spaces, where Fig. 4(a) and (b), respectively, corresponds to unsupervised and supervised learning strategies.

IV. IFD DESIGNS AND IMPLEMENTATIONS: FROM UNSUPERVISED TO SUPERVISED NEURAL NETWORKS

A study of the existing approaches reveals that both unsupervised and supervised machine learning techniques can be employed to enhance the IFD flexibility against system nonlinearity. For example, neural networks have received increasing attention for the designs of IFD. Directly motivated by Theorems 1–3, the fully connected neural networks are used in this section for the fault diagnosis purpose.

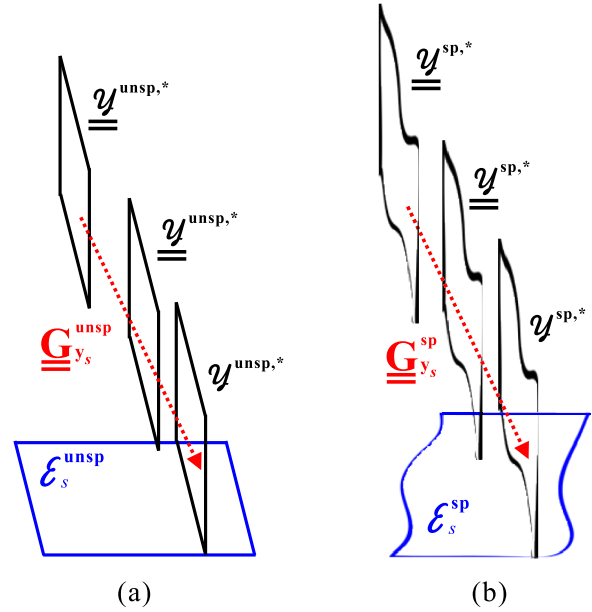


Fig. 4. More general versions of suspected and residual spaces. (a) Unsupervised learning. (b) Supervised learning.

A. Unit-Time Delay Operators

A recurrent neural network is capable of describing the dynamic behaviors of nonlinear systems, thanks to the unit-time delays attached to neurons [37]. As the system order increases, more delay units are necessary to fit the higher order dynamics, resulting in heavy computations together with problems associated with vanishing and exploding gradients [38]. To avoid such a situation and improve computation efficiency, the unit-time delay operator, denoted as z^{-1} , is used in the data preprocessing stage.

As shown in Fig. 5(a), $(s_p + s + 1)$ unit-time delay operators are defined on both $\mathbf{u}(k + s)$ and $\mathbf{y}(k + s)$ to obtain the inputs of an unsupervised neural network, i.e.,

$$\begin{aligned} \mathbf{u}(k + s - 1) &= z^{-1} \mathbf{u}(k + s) \quad \cdots, \\ \mathbf{u}(k - s_p - 1) &= z^{-(s_p + s + 1)} \mathbf{u}(k + s), \\ \mathbf{y}(k + s - 1) &= z^{-1} \mathbf{y}(k + s) \quad \cdots, \\ \mathbf{y}(k - s_p - 1) &= z^{-(s_p + s + 1)} \mathbf{y}(k + s). \end{aligned} \quad (77)$$

Similarly, $(s_p + s + 1)$ and $(s_p + 1)$ unit-time delay operators are defined on $\mathbf{u}(k + s)$ and $\mathbf{y}(k)$, respectively. The inputs of a supervised neural network given in Fig. 5(b) are $[\mathbf{u}_p^T \ \mathbf{u}_s^T]$ and

$$\begin{aligned} \mathbf{y}(k - 1) &= z^{-1} \mathbf{y}(k), \quad \cdots, \\ \mathbf{y}(k - s_p - 1) &= z^{-(s_p + 1)} \mathbf{y}(k). \end{aligned} \quad (78)$$

In addition, the unit-time delays are also used to obtain the reference outputs of both unsupervised and supervised neural networks. As observed from (77) and (78), and Fig. 5, these unit-time delay operators work independently of training neural networks, which reduces the requested layer number of neural networks from the inputs to outputs. In other words, according to the chain rule, the recommended location where delay operators are installed in this study can reduce the order of partial derivatives. Therefore, not only does it avoid the

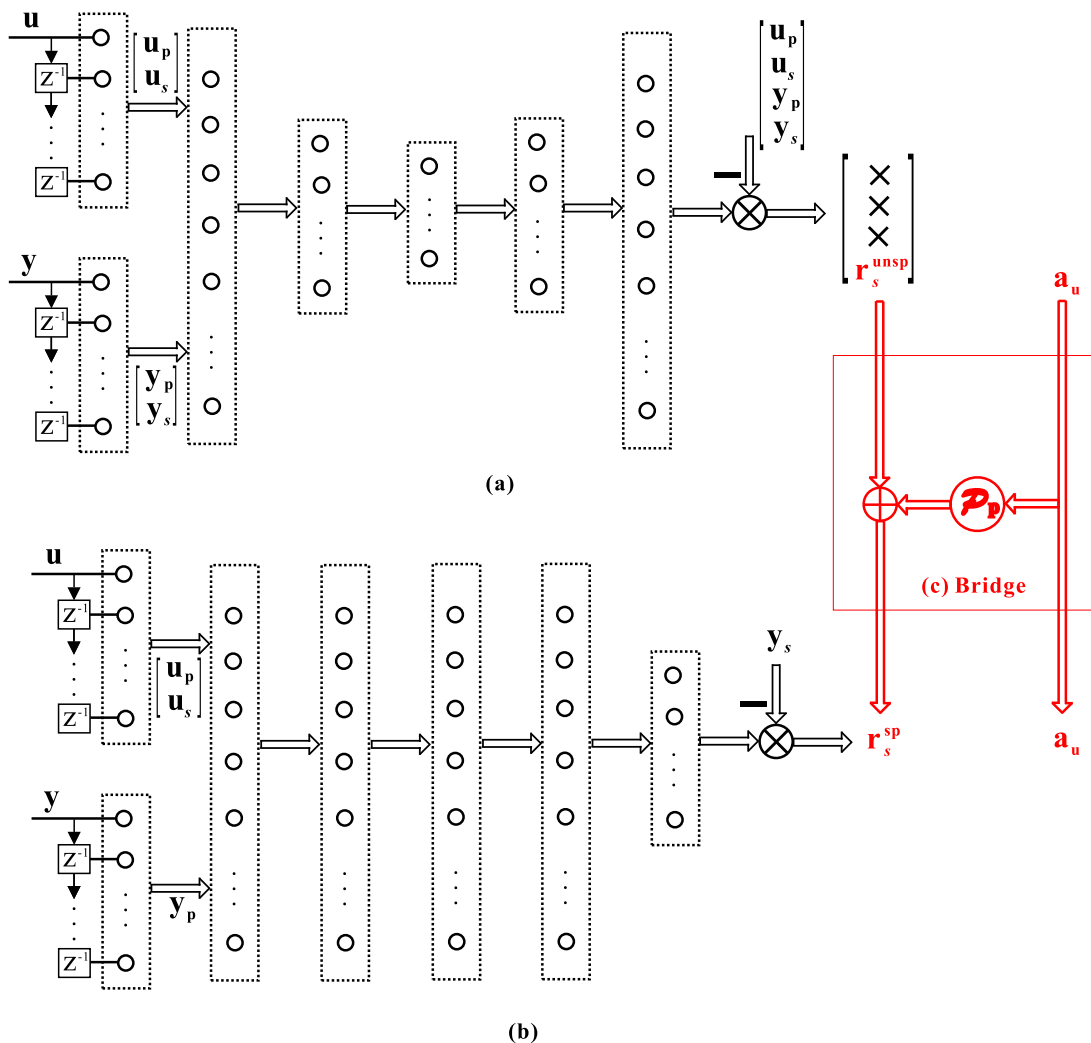


Fig. 5. Schematic of IFD frameworks with multiple unit-time delays, where “ \mathbf{a}_u ” represents an auxiliary variable. (a) is designed based on an unsupervised neural network (unsupervised learning-based IFD). (b) is implemented based on a supervised neural network (supervised learning-based IFD). (c) Highlighted in red is the bridge using an invertible neural network.

vanishing and exploding gradient problems but also improve the computation efficiency. It is worth mentioning that our recommended approach does not change the objective function of the neural networks.

B. Unsupervised IFD Approaches

Similar to (39a), the loss function L of the unsupervised neural network can be defined by

$$L^{\text{unsp}} = \left\| \begin{bmatrix} \mathbf{z}_p \\ \mathbf{u}_s \\ \mathbf{y}_s \end{bmatrix} - \underbrace{\mathcal{H}^{\text{unsp}}(\Theta_z \ \Theta_u \ \Theta_y)}_{\mathcal{H}^{\text{unsp}}(\Theta)} \begin{bmatrix} \mathbf{z}_p \\ \mathbf{u}_s \\ \mathbf{y}_s \end{bmatrix} \right\|_2^2 \quad (79)$$

where $\mathcal{H}^{\text{unsp}}(\Theta)$ is the architecture of the neural network with the hyperparameter Θ and the subscripts of Θ signify their corresponding outputs. Therefore, the optimal Θ can be obtained by minimizing L^{unsp} , i.e.,

$$\Theta^* = (\Theta_z^* \ \Theta_u^* \ \Theta_y^*) = \arg \min_{\Theta} L^{\text{unsp}} \quad (80)$$

based on which the residual signal using $\mathcal{H}^{\text{unsp}}(\Theta^*)$ is

$$\mathbf{r}_s^{\text{unsp}}(k) = \mathbf{y}_s - \mathcal{H}^{\text{unsp}}(\Theta_y^*) \circ \begin{bmatrix} \mathbf{z}_p \\ \mathbf{u}_s \\ \mathbf{y}_s \end{bmatrix} = (34a). \quad (81)$$

In Fig. 5(a), $\mathbf{r}_s^{\text{unsp}}(k)$ is highlighted in red, and “ \times ” refers to the reconstruction errors that may be neglected.

Combining with Fig. 5(a), the implementation procedures of an unsupervised neural network-based IFD approach are summarized in Algorithms 1 and 2.

C. Supervised IFD Approaches

Considering a supervised neural network $\mathcal{H}^{\text{sp}}(\Theta)$, its loss function is formulated as follows:

$$L^{\text{sp}} = \left\| \mathbf{y}_s - \mathcal{H}^{\text{sp}}(\Theta) \begin{bmatrix} \mathbf{z}_p \\ \mathbf{u}_s \end{bmatrix} \right\|_2^2 \quad (83)$$

when $\mathcal{M}_{\text{euc}}^2(\mathbf{y}_s, \mathbf{y}_s^{\text{sp}})$ is adopted. Then, the optimal hyperparameter can be obtained according to

$$\Theta^* = \arg \min_{\Theta} L^{\text{sp}} \quad (84)$$

Algorithm 1 Off-Line Learning: Unsupervised Neural Network-Based IFD Approaches

- 1: Collect system measurements \mathbf{u} and \mathbf{y} ;
- 2: Pre-process \mathbf{u} and \mathbf{y} by using the delay units according to (77) to obtain \mathbf{z}_p , \mathbf{u}_s and \mathbf{y}_s ;
- 3: Construct an unsupervised neural network whose architecture is $\mathcal{H}^{\text{unsp}}(\Theta)$ and loss function is defined in (79);
- 4: Update Θ according to (80) and obtain Θ_y^* ;
- 5: Generate the residual signal $\mathbf{r}_s^{\text{unsp}}$ via (81);
- 6: Determine the threshold J_{th} for T^2 based on (32);
- 7: (If necessary) identify the fault features based on system knowledge or using faulty data.

Algorithm 2 Online Fault Diagnosis: Unsupervised Neural Network-Based IFD Approaches

- 1: Read real-time data and employ unit-time delay operators;
- 2: Compute the online residual signal according to (81);
- 3: Compute the test statistic via (31);
- 4: Make an FD decision according to

$$T^2(\mathbf{r}_s^{\text{unsp}}(\text{online})) - J_{th} < 0 \implies \text{Fault-free}; \quad (82)$$

$$\text{Otherwise} \implies \text{Faulty};$$

- 5: (If necessary) diagnose the fault by using $\mathbf{r}_s^{\text{unsp}}(\text{online})$.

Algorithm 3 Off-Line Learning: Supervised Neural Network-Based IFD Approaches

- 1: Collect system measurements \mathbf{u} and \mathbf{y} ;
- 2: Pre-process \mathbf{u} and \mathbf{y} by using the delay units according to (77) and (78) to obtain \mathbf{z}_p and \mathbf{u}_s ;
- 3: Construct a supervised neural network whose architecture is $\mathcal{H}^{\text{sp}}(\Theta)$ and loss function is defined in (83);
- 4: Update Θ via (84) and obtain Θ^* ;
- 5: Generate the residual signal \mathbf{r}_s^{sp} via (85);
- 6: Determine the threshold J_{th} for T^2 based on (32);
- 7: (If necessary) identify the fault features based on system knowledge or using faulty data.

which allows for constructing the residual signal as follows:

$$\mathbf{r}_s^{\text{sp}}(k) = \mathbf{y}_s - \mathcal{H}^{\text{sp}}(\Theta^*) \begin{bmatrix} \mathbf{z}_p \\ \mathbf{u}_s \end{bmatrix} = (34b). \quad (85)$$

According to the analysis above and the schematic presented in Fig. 5(b), the complete designs and implementation procedures are given in Algorithms 3 and 4.

D. Invertible Neural Network-Aided Bridge

In order to build the bridge between unsupervised and supervised neural network-based IFD approaches, we introduce an auxiliary variable \mathbf{a}_u . The mapping generated by $\mathcal{H}(\Theta_{\text{ext}})$ in the red square of Fig. 5(c) is

$$\begin{bmatrix} \mathbf{r}_s^{\text{sp}} \\ \mathbf{a}_u \end{bmatrix} = \underbrace{\begin{pmatrix} \mathbf{I} & \mathcal{H}(\Theta_p) \\ \mathbf{0} & \mathbf{I} \end{pmatrix}}_{\mathcal{H}(\Theta_{\text{ext}})} \begin{bmatrix} \mathbf{r}_s^{\text{unsp}} \\ \mathbf{a}_u \end{bmatrix} \quad (87)$$

Algorithm 4 Online Fault Diagnosis: Supervised Neural Network-Based IFD Approaches

- 1: Read real-time data and employ unit-time delay operators;
- 2: Compute the online residual signal according to (85);
- 3: Compute the test statistic via (31);
- 4: Make an FD decision according to

$$T^2(\mathbf{r}_s^{\text{sp}}(\text{online})) - J_{th} < 0 \implies \text{Fault-free}; \quad (86)$$

$$\text{Otherwise} \implies \text{Faulty};$$

- 5: (If necessary) diagnose the fault by using $\mathbf{r}_s^{\text{sp}}(\text{online})$.

where the loss function of $\mathcal{H}(\Theta_{\text{ext}})$ is chosen as

$$L^{\text{ext}} = \left\| \begin{bmatrix} \mathbf{r}_s^{\text{sp}} \\ \mathbf{a}_u \end{bmatrix} - \mathcal{H}(\Theta_{\text{ext}}) \begin{bmatrix} \mathbf{r}_s^{\text{unsp}} \\ \mathbf{a}_u \end{bmatrix} \right\|_2^2. \quad (88)$$

Without loss of generality, we can choose $\mathbf{a}_u = \mathbf{r}_s^{\text{unsp}}$ to prove the existence of the bridge. Then, minimizing L^{ext} obtains

$$\Theta_{\text{ext}}^* = \arg \min_{\Theta_{\text{ext}}} L^{\text{ext}} \iff \mathcal{H}(\Theta_p^*) = \mathcal{P}_p \quad (89)$$

as shown in Fig. 5(c). It indicates that

$$\mathbf{r}_s^{\text{sp}} = (\mathcal{P}_p + \mathbf{I})\mathbf{r}_s^{\text{unsp}}, \mathbf{r}_s^{\text{unsp}} = \mathbf{r}_s^{\text{unsp}} \quad (90)$$

and

$$\begin{aligned} \mathbf{r}_s^{\text{unsp}} &= \mathbf{r}_s^{\text{sp}} - \mathcal{P}_p \\ \mathbf{r}_s^{\text{unsp}} &= \mathbf{r}_s^{\text{sp}} - \mathcal{H}(\Theta_p^*)\mathbf{r}_s^{\text{unsp}} \\ &= \mathbf{r}_s^{\text{sp}} - \mathcal{H}(\Theta_p^*)\mathbf{a}_u, \\ \mathbf{a}_u &= \mathbf{a}_u. \end{aligned} \quad (91)$$

It is interesting to see that, even though the inverse operation is not used, we can obtain the bridge by

$$\begin{bmatrix} \mathbf{r}_s^{\text{sp}} \\ \mathbf{a}_u \end{bmatrix} = \begin{pmatrix} \mathbf{I} & \mathcal{P}_p \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \begin{bmatrix} \mathbf{r}_s^{\text{unsp}} \\ \mathbf{a}_u \end{bmatrix} \quad (92)$$

and

$$\begin{bmatrix} \mathbf{r}_s^{\text{unsp}} \\ \mathbf{a}_u \end{bmatrix} = \begin{pmatrix} \mathbf{I} & -\mathcal{P}_p \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \begin{bmatrix} \mathbf{r}_s^{\text{sp}} \\ \mathbf{a}_u \end{bmatrix}. \quad (93)$$

In addition, the Jacobian matrix \mathbf{J} of $\mathcal{H}(\Theta_{\text{ext}})$ is

$$\mathbf{J} = \begin{pmatrix} \mathbf{I} & \mathcal{P}_p \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \quad (94)$$

whose determinant is always positive, i.e.,

$$|\mathbf{J}| = \left| \begin{pmatrix} \mathbf{I} & \mathcal{P}_p \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \right| = 1. \quad (95)$$

It means that $\mathcal{H}(\Theta_{\text{ext}})$ and its component $\mathcal{P}_{\text{sp/unsp}}$ are globally invertible.

Also, (92), (93), and (95) can be used as an auxiliary evidence to illustrate the validity of Theorem 1.

Remark 7: The solutions to the bridge depicted in Fig. 5(c) are not unique. For example, a "reversible residual network" given in [39] can also be directly employed to obtain a one-to-one mapping between $\mathbf{r}_s^{\text{unsp}}$ and \mathbf{r}_s^{sp} . ∇

V. CONCLUSION

Thanks to the advanced machine learning techniques, there appears to have room for further development of nonlinear IFD. System data can provide us with information to reveal physical principles, making more informative inferences and decisions possible. As the interpretability of these advanced methods is fundamental, the interplay between physical principles (such as a mathematical system description whose parameters may be unknown) and data becomes more critical.

This perspective article has developed three theorems related to IFD and related parameter identification for nonlinear dynamic systems. In order to obtain a unified form, Theorem 1 builds a bridge between unsupervised and supervised learning-based residual generators. Theorems 2 and 3, respectively, corresponding to unsupervised and supervised learning, develop the IFD structures and illustrate their optimal performance for fault detection. With the aid of three different kinds of neural networks, Section IV details the specific designs and implementations of Theorems 1–3. This work lays a foundation for the further development of explainable IFD.

The perspective article ends with discussions, including our opinions, expected challenges, and future research opportunities.

- 1) Not all IFD approaches are satisfactory from the viewpoint of fault-diagnosis performance. Our work provides researchers with some instructive guidance on designing more effective IFD algorithms, including both unsupervised and supervised learning-based schemes.
- 2) As shown in Fig. 5 and Theorem 2, not all results obtained through unsupervised learning are useful for fault diagnosis. The conclusion is also true and easier to explain in the case of linear dynamic systems. For example, the purpose of singular value decomposition [40] and QR [5] used in linear approaches is not for dimension reduction although they can do so. Keep in mind that the ultimate goal of unsupervised learning adopted in IFD approaches is always to minimize the reconstruction error of system outputs for both linear and nonlinear systems.
- 3) The initial must step to apply the results obtained through this work is to show the existence of the bridge $\mathcal{P}_{\text{sp/unsup}}$ and $\mathcal{P}_{\text{unsup/sp}}$ given in Theorem 1. Then, a series of invertible machine learning tools, such as invertible neural networks in [41], can be employed to build this bridge. The process should not be reversed.
- 4) An unmentioned condition in this study for nonlinear IFD approaches is that (6) is output reconstructible [42]. In fact, the condition is weak and easy to satisfy for our proposed IFD frameworks because a multilayer neural network can approximate a continuous function (such as $\phi(\cdot)$ and $v(\cdot)$ in Fig. 1) to an arbitrary accuracy [43].
- 5) The data-based model (19) is obtained from (6) with the assumption of additivity given in (18a) and (18b). In practical applications, many classes of nonlinear systems satisfy (18a) and (18b), such as Sector bounded systems, Hammerstein systems, and affine systems.

- 6) It is a fact that nonlinear IFD approaches and the associated thresholds show dependence on \mathbf{u} and $\mathbf{x}(0)$ [13]. In order to obtain a reasonable dataset for training, a uniformly distributed \mathbf{u} over different operation points is recommended so that a persistent excitation of the global system nonlinearities can be achieved.
- 7) In addition to the bridge developed in this work, other aspects of system identification and IFD approaches deserve more investigations, such as the following.
 - a) How many neurons should be used when utilizing a neural network?
 - b) How to determine the nonlinearity degree to approximate an unknown dynamic system without the problem of overfitting or underfitting?
 - c) What is the minimum size of training data for obtaining the desired performance?
- 8) Neural networks can generate a reproducing kernel Hilbert space [44], by which modeling nonlinear system dynamics is achievable. This comment has pushed us toward developing neural network-aided orthogonal projections to complete both the parameter (and system) identification and IFD tasks [45].
- 9) Zero-shot (or few-shot) learning-based IFD approaches do not mean that trustworthy results can be obtained without (or with less) data samples. On the contrary, sufficient data samples are necessary as learning prerequisites in zero-shot (and few-shot) learning [46].
- 10) The Vapnik–Chervonenkis dimension and the Rademacher complexity, respectively, corresponding to the data-independent and data-dependent measures [47], can be used to bound the generalization error of both unsupervised and supervised learning methods. Borrowing these theories from statistical learning will open up a new avenue for evaluating the generalization capability of nonlinear IFD approaches.

APPENDIX

PROOF OF THEOREM 1

Given an arbitrary unsupervised learning approach whose notation is \mathcal{U} , it generates the nonlinear mapping (35) with the objective function, such as defined via $\mathcal{M}_{\text{cuc}}^2$

$$\begin{aligned} \min \left\| \begin{bmatrix} \mathbf{z}_p \\ \mathbf{u}_s \\ \mathbf{y}_s \end{bmatrix} - \begin{bmatrix} \hat{\mathbf{z}}_p \\ \hat{\mathbf{u}}_s \\ \hat{\mathbf{y}}_s \end{bmatrix} \right\|_2^2 \\ \implies \min \|\mathbf{y}_s - \hat{\mathbf{y}}_s\|_2^2 = \min \|\mathbf{e}_s\|_2^2. \end{aligned} \quad (\text{A.1})$$

In (35), the input of \mathcal{U}_z , \mathcal{U}_u , and \mathcal{U}_y is $[\mathbf{z}_p^T \ \mathbf{u}_s^T \ \mathbf{y}_s^T]^T$, and the subscripts correspond to the output variables.

Three steps are necessary to complete the proof.

Step 1 (Generation of Innovation Error \mathbf{e}_s Using Unsupervised Learning): Based on (23a) and (23b), one can obtain

$$\begin{aligned} \mathbf{y}(k) &= v_{\hat{\mathbf{y}}\hat{\mathbf{x}}}\hat{\mathbf{x}}(k) + v_{\hat{\mathbf{y}}\mathbf{u}}\mathbf{u}(k) + \mathbf{e}(k) \\ \mathbf{y}(k+1) &= v_{\hat{\mathbf{y}}\hat{\mathbf{x}}}\hat{\mathbf{x}}(k+1) \\ &\quad + v_{\hat{\mathbf{y}}\mathbf{u}}\mathbf{u}(k+1) + \mathbf{e}(k+1) \\ &= v_{\hat{\mathbf{y}}\mathbf{u}}\mathbf{u}(k+1) + \mathbf{e}(k+1) + v_{\hat{\mathbf{y}}\hat{\mathbf{x}}}\ell\mathbf{y}(k) \end{aligned}$$

$$\begin{aligned}
& + v_{\hat{y}\hat{x}} \underbrace{(\phi_{\hat{x}\hat{x}} - \ell v_{\hat{y}\hat{x}})}_{\mathcal{A}_\ell} \hat{\mathbf{x}}(k) \\
& + v_{\hat{y}\hat{x}} \underbrace{(\phi_{\hat{x}\mathbf{u}} - \ell v_{\hat{y}\mathbf{u}})}_{\mathcal{B}_\ell} \mathbf{u}(k) \\
& \dots
\end{aligned} \tag{A.2}$$

Then, $\mathbf{e}_s(k)$ can be described by

$$\mathbf{e}_s(k) = \mathbf{y}_s(k) - \underline{\Upsilon}_y \mathbf{y}_s(k) - \underline{\Upsilon}_x \hat{\mathbf{x}}(k) - \underline{\Upsilon}_u \mathbf{u}_s(k) \tag{A.3}$$

where $\underline{\Upsilon}_y$, $\underline{\Upsilon}_x$, and $\underline{\Upsilon}_u$ are the nonlinear composite operators

$$\underline{\Upsilon}_y = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} \\ v_{\hat{y}\hat{x}}\ell & \mathbf{0} & \dots & \mathbf{0} \\ \vdots & \ddots & \ddots & \vdots \\ v_{\hat{y}\hat{x}}\mathcal{A}_\ell^{s-1}\ell & \dots & v_{\hat{y}\hat{x}}\ell & \mathbf{0} \end{pmatrix}, \tag{A.4}$$

$$\underline{\Upsilon}_u = \begin{pmatrix} v_{\hat{y}\mathbf{u}} & \mathbf{0} & \dots & \mathbf{0} \\ v_{\hat{y}\hat{x}}\mathcal{B}_\ell & v_{\hat{y}\mathbf{u}} & \dots & \mathbf{0} \\ \vdots & \ddots & \ddots & \vdots \\ v_{\hat{y}\hat{x}}\mathcal{A}_\ell^{s-1}\mathcal{B}_\ell & \dots & v_{\hat{y}\hat{x}}\mathcal{B}_\ell & v_{\hat{y}\mathbf{u}} \end{pmatrix}, \tag{A.5}$$

$$\underline{\Upsilon}_x = \begin{pmatrix} v_{\hat{y}\hat{x}} \\ \vdots \\ v_{\hat{y}\hat{x}}\mathcal{A}_\ell^s \end{pmatrix}. \tag{A.6}$$

Through some mathematical manipulations, (26) can be rewritten as

$$\mathbf{y}_s(k) = \Upsilon_x \hat{\mathbf{x}}(k) + \Upsilon_u \mathbf{u}_s(k) + \Upsilon_e \mathbf{e}_s(k) \tag{A.7}$$

which can be achieved by using supervised learning approaches, as shown in (26), i.e.,

$$\hat{\mathbf{y}}_s(k) = \underbrace{(\Upsilon_x \mathcal{L}_p \quad \Upsilon_u)}_{\mathcal{S} := (\mathcal{S}_z \quad \mathcal{S}_u)} \begin{bmatrix} \mathbf{z}_p(k) \\ \mathbf{u}_s(k) \end{bmatrix} \tag{A.8}$$

where \mathcal{S} is the nonlinear composite operator generated by supervised learning.

Step 2 (A Unified Description of the Transformation Between $\mathbf{r}_s^{\text{unsp}}$ and \mathbf{r}_s^{sp}): Based on (A.3) and (A.8), two kinds of residual generators can be described by

$$\mathbf{r}_s^{\text{unsp}} = \mathbf{y}_s - \mathcal{U}_y \begin{bmatrix} \mathbf{z}_p \\ \mathbf{u}_s \\ \mathbf{y}_s \end{bmatrix} = \mathbf{e}_s, \tag{A.9a}$$

$$\mathbf{r}_s^{\text{sp}} = \mathbf{y}_s - (\mathcal{S}_z \quad \mathcal{S}_u) \begin{bmatrix} \mathbf{z}_p \\ \mathbf{u}_s \end{bmatrix} = \Upsilon_e \mathbf{e}_s. \tag{A.9b}$$

Step 3 (The Existence of $\mathcal{P}_{\text{sp/unsp}}^{-1}$): In fact, $\mathbf{r}_s^{\text{unsp}}$ in (A.9a) and \mathbf{r}_s^{sp} in (A.9b) have shown

$$\mathbf{r}_s^{\text{sp}} = \mathcal{P}_{\text{sp/unsp}} \mathbf{r}_s^{\text{unsp}}, \quad \mathcal{P}_{\text{sp/unsp}} := \Upsilon_e. \tag{A.10}$$

For the sake of simplicity, we rewrite (A.10) as

$$\mathbf{r}_s^{\text{sp}} = \mathbf{r}_s^{\text{unsp}} + \mathcal{P}_p \mathbf{r}_s^{\text{unsp}} \tag{A.11}$$

where \mathcal{P}_p is the component of $\mathcal{P}_{\text{sp/unsp}}$, i.e.,

$$\mathcal{P}_p = \begin{pmatrix} \mathbf{0} & \dots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ v_{\hat{y}\hat{x}}\phi_{\hat{x}\hat{x}}^{s-1}\ell_{\hat{x}\hat{e}} & \dots & \mathbf{0} \end{pmatrix}. \tag{A.12}$$

Consider a nonlinear operator $\mathcal{P}_{\text{unsp/sp}}$ and the Lipschitz constant of \mathcal{P}_p is $\mathbf{Lip} < 1$. By defining a sequence $\{\mathbf{e}_{s,j}\} \in \mathbf{e}_s$, one can obtain the following innovation:

$$\begin{aligned}
\mathcal{P}_{\text{unsp/sp}} : \mathbf{e}_{s,j+1} &= \Upsilon_e \mathbf{e}_s - \mathcal{P}_p \mathbf{e}_{s,j} \\
\implies \lim_{j \rightarrow \infty} \mathbf{e}_{s,j} &= \mathcal{P}_{\text{unsp/sp}} \circ \Upsilon_e \mathbf{e}_s.
\end{aligned} \tag{A.13}$$

Given a positive integer n , we have

$$\begin{aligned}
\|\mathbf{e}_{s,j+n} - \mathbf{e}_{s,j}\|_2 &\leq \|\mathbf{e}_{s,j+n} - \mathbf{e}_{s,j+n-1}\|_2 \\
&+ \dots \|\mathbf{e}_{s,j+1} - \mathbf{e}_{s,j}\|_2 \\
&= \|\mathcal{P}_p(\mathbf{e}_{s,j+n-1} - \mathbf{e}_{s,j+n-2})\|_2 + \dots \\
&+ \|\mathcal{P}_p(\mathbf{e}_{s,j} - \mathbf{e}_{s,j-1})\|_2 \\
&\leq (\mathbf{Lip}^{j+n-1} + \dots + \mathbf{Lip}^j) \|\mathbf{e}_{s,1} - \mathbf{e}_{s,0}\|_2 \\
&= \frac{1 - \mathbf{Lip}^{n-1}}{1 - \mathbf{Lip}} \mathbf{Lip}^j \|\mathbf{e}_{s,1} - \mathbf{e}_{s,0}\|_2 \\
&\leq \frac{\mathbf{Lip}^j}{1 - \mathbf{Lip}} \|\mathbf{e}_{s,1} - \mathbf{e}_{s,0}\|_2.
\end{aligned} \tag{A.14}$$

It indicates that, given an arbitrary small ε , there always exists an n such that

$$\|\mathbf{e}_{s,j+n} - \mathbf{e}_{s,j}\|_2 < \varepsilon \text{ or } \|\mathbf{e}_{s,j+n} - \mathbf{e}_{s,j}\|_2^2 < \varepsilon \tag{A.15}$$

when $\mathbf{Lip} < 1$, which guarantees the existence of $\mathcal{P}_{\text{sp/unsp}}^{-1} = \mathcal{P}_{\text{unsp/sp}}$. It is worth mentioning that $\{\mathbf{e}_{s,j}\}$ is a Cauchy sequence [48], and (A.15) holds for $\mathbf{Lip} < 1$ because of the Banach fixed-point theorem [49]. Furthermore, $\mathbf{Lip} < 1$ is a weak condition/requirement to satisfy when choosing \mathcal{P}_p .

Hence, the proof of Theorem 1 is completed.

REFERENCES

- [1] H. Chen, B. Jiang, S. X. Ding, and B. Huang, "Data-driven fault diagnosis for traction systems in high-speed trains: A survey, challenges, and perspectives," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 3, pp. 1700–1716, Mar. 2022.
- [2] S. X. Ding, *Model-Based Fault Diagnosis Techniques: Design Schemes, Algorithms and Tools*. Berlin, Germany: Springer-Verlag, 2008.
- [3] D. Zhou, Y. Zhao, Z. Wang, X. He, and M. Gao, "Review on diagnosis techniques for intermittent faults in dynamic systems," *IEEE Trans. Ind. Electron.*, vol. 67, no. 3, pp. 2337–2347, Mar. 2020.
- [4] S. J. Qin, "Survey on data-driven industrial process monitoring and diagnosis," *Annu. Rev. Control*, vol. 36, no. 2, pp. 220–234, Dec. 2012.
- [5] B. Huang and R. Kadali, *Dynamic Modeling, Predictive Control and Performance Monitoring: A Data-driven Subspace Approach*. London, U.K.: Springer-Verlag, 2008.
- [6] R. Isermann, "Model-based fault-detection and diagnosis—status and applications," *Annu. Rev. Control*, vol. 29, no. 1, pp. 71–85, Jan. 2005.
- [7] V. Venkatasubramanian, R. Rengaswamy, S. N. Kavuri, and K. Yin, "A review of process fault detection and diagnosis: Part III: Process history based methods," *Comput. Chem. Eng.*, vol. 27, no. 3, pp. 327–346, Mar. 2003.
- [8] S. Yin, S. X. Ding, X. Xie, and H. Luo, "A review on basic data-driven approaches for industrial process monitoring," *IEEE Trans. Ind. Electron.*, vol. 61, no. 11, pp. 6418–6428, Nov. 2014.
- [9] C. Dai, Z. Liu, K. Hu, and K. Huang, "Fault diagnosis approach of traction transformers in high-speed railway combining kernel principal component analysis with random forest," *IET Electr. Syst. Transp.*, vol. 6, no. 3, pp. 202–206, Sep. 2016.
- [10] H. Chen and B. Jiang, "A review of fault detection and diagnosis for the traction system in high-speed trains," *IEEE Trans. Intell. Transp. Syst.*, vol. 21, no. 2, pp. 450–465, Feb. 2020.
- [11] R. J. Patton and J. Chen, "Review of parity space approaches to fault diagnosis for aerospace systems," *J. Guid., Control Dyn.*, vol. 17, no. 2, pp. 278–285, Mar. 1994.

- [12] Z. Zhao, W. Jiang, and W. Gao, "Health evaluation and fault diagnosis of medical imaging equipment based on neural network algorithm," *Comput. Intell. Neurosci.*, vol. 2021, pp. 1–9, Sep. 2021, doi: [10.1155/2021/6092461](https://doi.org/10.1155/2021/6092461).
- [13] S. X. Ding, *Advanced Methods for Fault Diagnosis and Fault-tolerant Control*. Berlin, Germany: Springer, 2020.
- [14] J. Chen and R. J. Patton, *Robust Model-Based Fault Diagnosis for Dynamic System*. Norwell, MA, USA: Kluwer, 1998.
- [15] H. Wang, Z. Liu, A. Núñez, and R. Dollevoet, "Entropy-based local irregularity detection for high-speed railway catenaries with frequent inspections," *IEEE Trans. Instrum. Meas.*, vol. 68, no. 10, pp. 3536–3547, Oct. 2019.
- [16] S. X. Ding, *Data-Driven Design of Fault Diagnosis and Fault-tolerant Control Systems*. London, U.K.: Springer-Verlag, 2014.
- [17] H.-G. Zhang, X. Zhang, Y.-H. Luo, and J. Yang, "An overview of research on adaptive dynamic programming," *Acta Automat. Sinica*, vol. 39, no. 4, pp. 303–311, Apr. 2013.
- [18] J. Wang, Z. Zhang, B. Tian, and Q. Zong, "Event-based robust optimal consensus control for nonlinear multiagent system with local adaptive dynamic programming," *IEEE Trans. Neural Netw. Learn. Syst.*, early access, Jun. 22, 2022, doi: [10.1109/TNNLS.2022.3180054](https://doi.org/10.1109/TNNLS.2022.3180054).
- [19] H. Jiang, H. Zhang, K. Zhang, and X. Cui, "Data-driven adaptive dynamic programming schemes for non-zero-sum games of unknown discrete-time nonlinear systems," *Neurocomputing*, vol. 275, pp. 649–658, Jan. 2018.
- [20] D. P. Bertsekas, "Value and policy iterations in optimal control and adaptive dynamic programming," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 28, no. 3, pp. 500–509, Mar. 2017.
- [21] R. Song, Q. Wei, H. Zhang, and F. L. Lewis, "Discrete-time non-zero-sum games with completely unknown dynamics," *IEEE Trans. Cybern.*, vol. 51, no. 6, pp. 2929–2943, Jun. 2021.
- [22] C. Hua, S. X. Ding, and Y. A. W. Shardt, "A new method for fault tolerant control through Q -learning," *IFAC-PapersOnLine*, vol. 51, no. 24, pp. 38–45, 2018.
- [23] T. de Bruin, K. Verbert, and R. Babšška, "Railway track circuit fault diagnosis using recurrent neural networks," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 28, no. 3, pp. 523–533, Mar. 2017.
- [24] J. Seshadrinath, B. Singh, and B. K. Panigrahi, "Incipient interturn fault diagnosis in induction machines using an analytic wavelet-based optimized Bayesian inference," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 25, no. 5, pp. 990–1001, May 2014.
- [25] Z. Liu, L. Wang, C. Li, and Z. Han, "A high-precision loose strands diagnosis approach for isoelectric line in high-speed railway," *IEEE Trans. Ind. Informat.*, vol. 14, no. 3, pp. 1067–1077, Mar. 2018.
- [26] T. Li, Z. Zhao, C. Sun, R. Yan, and X. Chen, "Multireceptive field graph convolutional networks for machine fault diagnosis," *IEEE Trans. Ind. Electron.*, vol. 68, no. 12, pp. 12739–12749, Dec. 2021.
- [27] H. Chen, Z. Chai, O. Dogru, B. Jiang, and B. Huang, "Data-driven designs of fault detection systems via neural network-aided learning," *IEEE Trans. Neural Netw. Learn. Syst.*, early access, Apr. 14, 2021, doi: [10.1109/TNNLS.2021.3071292](https://doi.org/10.1109/TNNLS.2021.3071292).
- [28] A. Bellet, A. Habrard, and M. Sebban, "Metric learning," *Synth. Lectures Artif. Intell. Mach. Learn.*, vol. 9, no. 1, pp. 1–151, Jan. 2015.
- [29] H. Chen, Z. Chen, Z. Chai, B. Jiang, and B. Huang, "A single-side neural network-aided canonical correlation analysis with applications to fault diagnosis," *IEEE Trans. Cybern.*, vol. 52, no. 9, pp. 9454–9466, Sep. 2021, doi: [10.1109/TCYB.2021.3060766](https://doi.org/10.1109/TCYB.2021.3060766).
- [30] Y. Yang, S. X. Ding, and L. Li, "Parameterization of nonlinear observer-based fault detection systems," *IEEE Trans. Autom. Control*, vol. 61, no. 11, pp. 3687–3692, Nov. 2016.
- [31] H. Chen, L. Li, C. Shang, and B. Huang, "Fault detection for nonlinear dynamic systems with consideration of modeling errors: A data-driven approach," *IEEE Trans. Cybern.*, early access, Apr. 13, 2022, doi: [10.1109/TCYB.2022.3163301](https://doi.org/10.1109/TCYB.2022.3163301).
- [32] K. Fujimoto and T. Sugie, "Characterization of all nonlinear stabilizing controllers via observer-based kernel representations," *Automatica*, vol. 36, no. 8, pp. 1123–1135, Aug. 2000.
- [33] I. Mezić, "Koopman operator, geometry, and learning," 2020, *arXiv:2010.05377*.
- [34] A. L. Bruce, V. M. Zeidan, and D. S. Bernstein, "What is the Koopman operator? A simplified treatment for discrete-time systems," in *Proc. Amer. Control Conf. (ACC)*, Jul. 2019, pp. 1912–1917.
- [35] Y. Maki and K. A. Loparo, "A neural-network approach to fault detection and diagnosis in industrial processes," *IEEE Trans. Control Syst. Technol.*, vol. 5, no. 6, pp. 529–541, Nov. 1997.
- [36] I. Rivals and L. Personnaz, "Black-box modeling with state-space neural networks," *Neural Adapt. Control Technol.*, vol. 15, pp. 237–264, Apr. 1996.
- [37] S. Haykin, *Neural Networks and Learning Machines*. Hamilton, ON, Canada: Pearson, 2010.
- [38] H. Chen, Z. Chai, B. Jiang, and B. Huang, "Data-driven fault detection for dynamic systems with performance degradation: A unified transfer learning framework," *IEEE Trans. Instrum. Meas.*, vol. 70, pp. 1–12, 2021, doi: [10.1109/TIM.2020.3033943](https://doi.org/10.1109/TIM.2020.3033943).
- [39] A. N. Gomez, M. Ren, R. Urtasun, and R. B. Grosse, "The reversible residual network: Backpropagation without storing activations," in *Proc. 31st Int. Conf. Neural Inf. Process. Syst.*, Dec. 2017, pp. 2211–2221.
- [40] J. Wang and S. J. Qin, "A new subspace identification approach based on principal component analysis," *J. Process Control*, vol. 12, no. 8, pp. 841–855, Dec. 2002.
- [41] L. Ardizzone *et al.*, "Analyzing inverse problems with invertible neural networks," 2018, *arXiv:1808.04730*.
- [42] E. Sontag and Y. Wang, "Lyapunov characterizations of input to output stability," *SIAM J. Control Optim.*, vol. 39, no. 1, pp. 226–249, Jan. 2000.
- [43] K.-I. Funahashi, "On the approximate realization of continuous mappings by neural networks," *Neural Netw.*, vol. 2, no. 3, pp. 183–192, Mar. 1989.
- [44] D. Sejdinovic and A. Gretton, "What is an RKHS?" ResearchGate, Lect. Notes, Mar. 2012, pp. 1–24.
- [45] S. X. Ding, "A note on diagnosis and performance degradation detection in automatic control systems towards functional safety and cyber security," *Secur. Saf.*, vol. 1, p. 29, Aug. 2022, doi: [10.1051/sands/2022004](https://doi.org/10.1051/sands/2022004).
- [46] H. Chen, C. Cheng, O. Dogru, and B. Huang, "Performance evaluation of few-shot learning-based system identification," in *Proc. 5th Int. Conf. Robot. Control Autom. Eng.*, 2022.
- [47] D. Anguita, A. Ghio, L. Oneto, and S. Ridella, "A deep connection between the Vapnik–Chervonenkis entropy and the Rademacher complexity," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 25, no. 12, pp. 2202–2211, Dec. 2014.
- [48] S. Lang, *Algebraic Number Theory*. New York, NY, USA: Springer, 2013.
- [49] K. Ciesielski, "On Stefan Banach and some of his results," *Banach J. Math. Anal.*, vol. 1, no. 1, pp. 1–10, 2007.



Hongtian Chen (Member, IEEE) received the B.S. and M.S. degrees from the School of Electrical and Automation Engineering, Nanjing Normal University, Nanjing, China, in 2012 and 2015, respectively, and the Ph.D. degree from the College of Automation Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing, in 2019.

He had ever been a Visiting Scholar with the Institute for Automatic Control and Complex Systems, University of Duisburg-Essen, Duisburg, Germany, in 2018. He is currently a Post-Doctoral Fellow with the Department of Chemical and Materials Engineering, University of Alberta, Edmonton, AB, Canada. His research interests include process monitoring and fault diagnosis, data mining and analytics, machine learning, and quantum computation and their applications in high-speed trains, new energy systems, and industrial processes.

Dr. Chen was a recipient of the Grand Prize of Innovation Award of the Ministry of Industry and Information Technology of the People's Republic of China in 2019, the Excellent Ph.D. Thesis Award of Jiangsu Province in 2020, and the Excellent Doctoral Dissertation Award from the Chinese Association of Automation (CAA) in 2020. He also serves as an Associate Editor and a Guest Editor for a number of scholarly journals, such as IEEE TRANSACTIONS ON INSTRUMENTATION AND MEASUREMENT, IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS, and IEEE TRANSACTIONS ON ARTIFICIAL INTELLIGENCE.



Zhigang Liu (Senior Member, IEEE) received the Ph.D. degree in power system and its automation from Southwest Jiaotong University, Chengdu, China, in 2003.

He is currently a Full Professor with the School of Electrical Engineering, Southwest Jiaotong University. He is also a Guest Professor with Tongji University, Shanghai, China. He has authored three books and published more than 200 peer-reviewed journal articles and conference papers. His research interests include the electrical relationship of electric multiple unit (EMU) and traction, detection, and assessment of pantograph catenary in high-speed railways.

Dr. Liu won the Second Prize in the National Science and Technology Progress Award, the Sichuan Youth Science and Technology Award, the Special Prize of the China Railway Society, the Second Prize of the Science and Technology Progress Award of the Ministry of Education, and the China Electric Power Excellent Technologist Award. He was named the Outstanding Associate Editor of IEEE TRANSACTIONS ON INSTRUMENTATION AND MEASUREMENT for 2019 and 2020. He is also an Associate Editor-in-Chief of IEEE TRANSACTIONS ON INSTRUMENTATION AND MEASUREMENT and an Associate Editor of IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS and IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY.



Cesare Alippi (Fellow, IEEE) received the master's degree (*cum laude*) in electronic engineering and the Ph.D. degree from the Politecnico di Milano, Milan, Italy, in 1990 and 1995, respectively.

He has been a Visiting Researcher with University College London, London, U.K., the Massachusetts Institute of Technology, Cambridge, MA, USA, ESPCI, Paris, France, CASIA, Singapore, and the Agency for Science, Technology and Research (A*STAR), Singapore. He is currently a Full Professor of information processing systems with the Politecnico di Milano and a Full Professor of cyber-physical and embedded systems with the Università della Svizzera italiana, Lugano, Switzerland. He holds five patents. He published a monograph on *Intelligence for Embedded Systems* (Springer) in 2014 and (co)authored more than 200 papers in international journals and conference proceedings. His current research activity addresses adaptation and learning in nonstationary environments and intelligence for embedded systems.

Dr. Alippi is also a member and the chair of several IEEE committees. He received the IEEE Instrumentation and Measurement Society Young Engineer Award in 2004, the IBM Faculty Award in 2013, the 2016 IEEE TNNLS Outstanding Paper Award, and the 2016 International Neural Networks Society (INNS) Gabor Award. Among the others, he was the General Chair of the International Joint Conference on Neural Networks (IJCNN) in 2012 and its Program Chair and Co-Chair in 2014 and 2011, respectively. He was the General Chair of the IEEE Symposium Series on Computational Intelligence in 2014. He was an Associate Editor of the IEEE TRANSACTIONS ON INSTRUMENTATION AND MEASUREMENT and IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS. He is also a Distinguished Lecturer of the IEEE Computational Intelligence Society (CIS), a member of the Board of Governors of INNS, the Vice-President of IEEE CIS, and an Associate Editor of the *IEEE Computational Intelligence Magazine*.



Biao Huang (Fellow, IEEE) received the B.Sc. and M.Sc. degrees in automatic control from the Beijing University of Aeronautics and Astronautics, Beijing, China, in 1983 and 1986, respectively, and the Ph.D. degree in process control from the University of Alberta, Edmonton, AB, Canada, in 1997.

In 1997, he joined the Department of Chemical and Materials Engineering, University of Alberta, as an Assistant Professor, where he is currently a Professor. He held the positions of the Natural Sciences and Engineering Research Council of Canada's Industrial Research Chair in Control of Oil Sands Processes and the Alberta Innovates Technology Futures Industry Chair in Process Control. He has applied his expertise extensively in industrial practice, particularly in the oil sands industry. His research interests include process control, system identification, control performance assessment, Bayesian methods, and state estimation.

Dr. Huang is also a fellow of the Canadian Academy of Engineering and the Chemical Institute of Canada. He was a recipient of the Germany's Alexander von Humboldt Research Fellowship, the Canadian Chemical Engineer Society's Syncrude Canada Innovation and D. G. Fisher Awards, the APEGA Summit Research Excellence Award, the University of Alberta McCalla and Killam Professorship Awards, the Petro-Canada Young Innovator Award, and the Best Paper Award from the *Journal of Process Control*.



Derong Liu (Fellow, IEEE) received the Ph.D. degree in electrical engineering from the University of Notre Dame, Notre Dame, IN, USA, in 1994.

He was a Staff Fellow with the General Motors Research and Development Center, Warren, MI, USA, from 1993 to 1995. He was an Assistant Professor with the Department of Electrical and Computer Engineering, Stevens Institute of Technology, Hoboken, NJ, USA, from 1995 to 1999. He joined the University of Illinois at Chicago, Chicago, IL, USA, in 1999, where he became a Full Professor of electrical and computer engineering and of computer science in 2006. He was the Associate Director of the State Key Laboratory of Management and Control for Complex Systems, Institute of Automation, Chinese Academy of Sciences, Beijing, China, from 2010 to 2015. He has published 19 books.

Prof. Liu is also a fellow of the International Neural Network Society and the International Association of Pattern Recognition. He received the Faculty Early Career Development Award from the National Science Foundation in 1999, the University Scholar Award from the University of Illinois from 2006 to 2009, the Overseas Outstanding Young Scholar Award from the National Natural Science Foundation of China in 2008, the Outstanding Achievement Award from the Asia-Pacific Neural Network Assembly in 2014, the INNS Gabor Award in 2018, the IEEE TNNLS Outstanding Paper Award in 2018, the IEEE SMC Society Andrew P. Sage Best Transactions Paper Award in 2018, and the IEEE/CCA JAS Hsue-Shen Tsien Paper Award in 2019. He was selected for the 100 Talents Program by the Chinese Academy of Sciences in 2008. He was the Editor-in-Chief of the IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS from 2010 to 2015. He is also the Editor-in-Chief of *Artificial Intelligence Review* (Springer).