

# Multi-Output Selective Ensemble Identification of Nonlinear and Nonstationary Industrial Processes

Tong Liu, Sheng Chen, *Fellow, IEEE*, Shan Liang, *Member, IEEE*, Shaojun Gan, Chris J. Harris

**Abstract**—A key characteristic of biological systems is the ability to update the memory by learning new knowledge and removing out-of-date knowledge so that intelligent decision can be made based on the relevant knowledge acquired in the memory. Inspired by this fundamental biological principle, this paper proposes a multi-output selective ensemble regression (SER) for online identification of multi-output nonlinear time-varying industrial processes. Specifically, an adaptive local learning approach is developed to automatically identify and encode newly emerging process state by fitting a local multi-output linear model based on the multi-output hypothesis testing. This growing strategy ensures a highly diverse and independent local model set. The online modeling is constructed as a multi-output SER predictor by optimizing the combining weights of the selected local multi-output models based on a probability metric. An effective pruning strategy is also developed to remove the unwanted out-of-date local multi-output linear models in order to achieve low online computational complexity without sacrificing the prediction accuracy. A simulated two-output process and two real-world identification problems are used to demonstrate the effectiveness of the proposed multi-output SER over a range of benchmark schemes for real-time identification of multi-output nonlinear and nonstationary processes, in terms of both online identification accuracy and computational complexity.

**Index Terms**—Multi-output nonlinear time-varying industrial processes, adaptive local learning, multivariate statistic hypothesis testing, selective ensemble, pruning

## I. INTRODUCTION

Generally, a biological system has two types of evolution or adaptation. The first type of evolution or long-term adaptation is over many generations of a biological system, based on the famous survival of the fittest principle. Many artificial learning algorithms mimic this long-term evolution principle, including genetic algorithm [1] and differential evolutionary algorithm [2]. A biological system also experiences the second type of adaptation or short-term evolving during its ‘daily

life’. In its existence, a biological system must evolve or adapt to fast changing environment and, therefore, it constantly updates its memory or ‘brain’ by learning new knowledge and removing out-of-date knowledge in order to make intelligent decision based on the relevant knowledge stored in its memory. Inspired by this fundamental short-term biological adaptation principle, in this work, we design a multi-output selective ensemble regression (SER) for online identification of multi-output nonlinear time-varying industrial processes.

Many real-life processes exhibit inherently nonlinear and nonstationary dynamic behaviours [3]–[9], where the underlying data generating mechanisms are fast changing over time. Under such circumstances, a predictive model must have the flexibility of adapting its structure and parameters from the fast-arriving nonstationary data stream, in order to maintain its performance in the changing environment.

Construction of neural network models can be interpreted as learning and encoding the nonlinear dynamic information of the system presented in the training data in their hidden-layer nodes. For example, each hidden node of the radial basis function (RBF) network encodes a process state. When the system is nonstationary, it is necessary for a model to acquire the newly emerging process state by updating its structure. A classical approach for estimating and tracking the temporal variations of a nonlinear process is to employ adaptive algorithms such as recursive least square (RLS) [10], [11] and its more recent variant multi-innovation RLS (MRLS) [12]. In particular, if the process’s input space is known a priori, by covering the input space with sufficiently dense nodes, the online sequential extreme learning machine (OS-ELM) only needs to update model weights online using the RLS algorithm [13]–[15]. Because the size of OS-ELM has to be very large to cover the training data space, online adaptation of the model weights is computationally costly and, moreover, there is no guarantee that fixed hidden nodes, no matter how dense they are, will also cover the changing nonstationary data space well. Hence the OS-ELM is unsuitable for fast time-varying data.

An alternative to single global nonlinear modeling approach is the ensemble learning that employs multiple models to separately model the data space. The online ensemble learning (OEL) has attracted considerable attentions recently, but most of the researches focus on online classification [16]–[23] and the regression problem is rarely discussed. The development of OEL for nonstationary systems is a challenging task, as the ensemble learner must not only maintain highly diverse models to cover a wide range of the observed data subspaces, but also learn the new concept as fast as possible to timely capture the changing dynamics. The ensemble of OS-ELM (EOS-ELM) [24], however, does not have this capability, as all the base extreme learning machine (ELM) models are trained on the

T. Liu and S. Liang are with Key Laboratory of Dependable Service Computing in Cyber Physical Society (Ministry of Education) and School of Automation, Chongqing University, Chongqing 400044, China (E-mails: liutong42@cqu.edu.cn, lightsun@cqu.edu.cn).

S. Chen and C.J. Harris are with School of Electronics and Computer Science, University of Southampton, Southampton SO17 1BJ, UK (E-mails: sqc@ecs.soton.ac.uk, chrisharris57@msn.com). S. Chen is also with King Abdulaziz University, Jeddah 21589, Saudi Arabia.

S. Gan is with Beijing Key Laboratory of Traffic Engineering, College of Metropolitan Transportation, Beijing University of Technology, Beijing 100124, China (E-mail: s.gan@bjut.edu.cn).

T. Liu would like to thank the sponsorship of Chinese Scholarship Council for funding his research at School of Electronics and Computer Science, University of Southampton, UK. This work was partly supported by the National Natural Science Foundation of China under grant 61771077, Key Research Program of Chongqing Science & Technology under grant CSTC2017jcyjBX0025, and Graduate Scientific Research and Innovation Foundation of Chongqing, China, under grant CYB19072.

This manuscript is to appear in the special issue on Biologically inspired methodologies for sensing, control and decision making.

same dataset and they are not updated online.

To design better OEL models, we take the inspiration from the aforementioned fundamental learning principle of biological systems, namely, the ability to update the memory by learning new knowledge and removing out-of-date knowledge so that intelligent decision can be made based on the most relevant knowledge acquired in the memory. For single-output nonlinear and nonstationary processes, the recently proposed selective ensemble based multiple local model (SEMLM) learning [25] enables automatically identifying newly emerging process states online and combining the most up-to-date local linear models to make an accurate SER based prediction. The SEMLM, however, does not have the ability to remove out-of-date process states that are no longer needed. Hence, a potential problem associated with the SEMLM is that for fast time varying systems, over a long period of online adaptation, the local model set may grow to be very large, which may impose high online computational complexity. To overcome this drawback and in particular to fully consider the aforementioned short-term evolving principle of biological learning system, the work [26] further proposed a growing and pruning SER (GAP-SER) for nonlinear and nonstationary data modeling. The GAP-SER can not only learn the newly emerging concept but also forget the past accumulated old concepts that are no longer relevant, hence balancing well the so-called stability-plasticity dilemma of a learning system. The experimental results of [26] shows that the GAP-SER typically outperforms the SEMLM, in terms of both online prediction accuracy and computational complexity.

Although the GAP-SER achieves great success in online modeling of nonstationary data, it is suitable only for single-output nonlinear and nonstationary systems. This is because the key components of the GAP-SER, including its local model growing and pruning strategies as well as selective ensemble prediction, are restricted to single-output modeling and they are not applicable to multi-output modeling. Practical systems and processes often contain multiple manipulated variables and controlled variables that exhibit strong coupling. Although we may apply the GAP-SER to identify the multiple single-output models for a multi-output system, this will increase the modeling effort considerably and more importantly, it will lead to the degradation in achievable online prediction accuracy. The latter is owing to the fact that the multiple single-output models ignore the correlation or coupling effects of the different output variables.

Therefore, it is necessary to investigate efficient identification techniques for multi-output nonlinear and nonstationary systems. The online multi-output regression problem remains largely under-studied. To our best knowledge, most of the existing online modeling approaches are restricted to single-output systems. For example, like the GAP-SER, the methods of [27]–[29] are restricted to single-output systems. Although a few studies have developed predictive models for multi-output nonlinear systems, such as the support vector regression (SVR) [30], the ELM [31] and the orthogonal least squares (OLS) method [32]–[35], these models are designed specifically for stationary systems, and they cannot directly be applied to highly nonstationary systems. In the online soft sensor design,

developing multi-output soft sensor to predict multiple primary variables has been demonstrated to be vital to achieve better performance than building multiple single-output soft sensors [36]. This motivates our current work.

In this paper, we propose a multi-output SER-based evolving model for online identification of multi-output nonlinear and time-varying processes, which is inspired by the fundamental evolving principle of biological systems for coping with a fast changing environment. As aforementioned, the short-term evolving of a biological system involves three levels of adaptation: adaptively acquiring new knowledge, adaptively making intelligent decision based on the most relevant knowledge stored in the memory, and adaptively removing the out-of-date knowledge which are no longer valid from the memory. Our multi-output SER-based evolving system also involves these three levels of adaptation. At the level of acquiring new knowledge, the newly emerging process state is automatically identified and a multi-output local linear model is fitted to it. Note that the growing strategy of [26] cannot be applied to the multi-output system as it is based on the single-output statistical testing. We develop a new adaptive local learning approach with the appropriate statistics for multi-output hypothesis testing in order to construct a highly diverse and independent multi-output local model set. At the level of making intelligent decision or online prediction, we construct a multi-output SER predictor based on a probability metric by optimizing the corresponding combining weights of the selective ensemble of the local multi-output linear models. At the level of removing out-of-date knowledge, a pruning strategy is developed to remove reliably the old local multi-output linear models that are no longer needed in order to further reduce the online computational complexity without degrading the prediction accuracy.

The main contribution of this work therefore is to develop a highly efficient and accurate online evolving model for multi-output selective ensemble identification of nonlinear and time-varying processes, with the completed design of online adaptive growing and pruning strategies and SER based adaptive prediction modeling. This new multi-output SER evolving model fully complies with the aforementioned biological system learning principle. Three case studies, a simulated two-output nonlinear system, a non-isothermal continuous stirred-tank reactor (CSTR) process [37], [38] and a real-world microwave heating system [39], are used to demonstrate the superior performance of our proposed multi-output SER over a range of benchmark schemes, in terms of both real-time prediction accuracy and online computational complexity.

## II. PROPOSED METHOD

To achieve high online prediction accuracy while imposing low computational complexity, our proposed evolving model mimics the fundamental learning principle from biological systems in coping with a fast changing environment, namely, the ability to update memory by learning new knowledge, to make intelligent decision based on the most relevant knowledge in the memory, and to remove out-of-date knowledge from the memory. Consequently, it involves three levels of adaptation,

namely, adaptive local learning that grows the multi-output local linear model set by identifying the newly emerged process state, adaptive online prediction by selective ensemble of the most relevant local models from the memory, and adaptive removal of the most out-of-date local linear models in the memory to free space for acquiring new knowledge. We now detail these levels of adaptation or the three components of our proposed evolving model.

#### A. Adaptive local learning via multivariate statistic

At the first level of adaptation, the evolving model absorbs new information from the fast-arriving data stream, and this is achieved by an adaptive local learning strategy which automatically encodes the newly emerging process states. Consider a nonlinear time-varying process with  $m$ -dimensional input  $\mathbf{x}(t) \in \mathbb{R}^m$  and  $p$ -dimensional output  $\mathbf{y}(t) \in \mathbb{R}^p$ . The task of local learning is to establish the local experts  $\{\mathbf{f}_l\}_{l=1}^L$  that accurately model the process's  $L$  local states or local regions, represented by the sub-datasets  $\{\mathbf{X}_l, \mathbf{Y}_l\}_{l=1}^L$ , where  $\mathbf{f}_l$  are linear models. The local linear model set  $\{\mathbf{f}_l\}_{l=1}^L$  correspond to the process states observed over time and, therefore, they represent the knowledge of the process acquired and stored in the evolving model's memory.

The basic idea of adaptive local learning is as follow. Let a local window  $\mathcal{W}_{\text{ini}} = \{\mathbf{X}_{\text{ini}} \in \mathbb{R}^{W_G \times m}, \mathbf{Y}_{\text{ini}} \in \mathbb{R}^{W_G \times p}\}$  with  $W_G$  consecutive samples  $\{\mathbf{x}(t), \mathbf{y}(t)\}_{t=t_{\text{ini}}+1}^{t_{\text{ini}}+W_G}$  be initially set. A multi-output local linear model  $\mathbf{f}_{\text{ini}}$  is built on it as

$$\hat{\mathbf{Y}}_{\text{ini}} = \mathbf{f}_{\text{ini}}(\mathbf{X}_{\text{ini}}) = \Phi \Omega, \quad (1)$$

where  $\Phi = [\mathbf{1}_{W_G} \ \mathbf{X}_{\text{ini}}] \in \mathbb{R}^{W_G \times (1+m)}$ ,  $\mathbf{1}_{W_G}$  is the  $W_G$ -dimensional vector whose elements are all one, and the model parameter matrix  $\Omega \in \mathbb{R}^{(1+m) \times p}$  is given by the least squares (LS) estimate as

$$\Omega = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{Y}_{\text{ini}}. \quad (2)$$

In the challenging situation of the singular input observation matrix  $\Phi$  and multicollinearity among the inputs, we may change the LS estimate to the regularized LS estimate

$$\Omega = (\Phi^T \Phi + \lambda \mathbf{I}_{m+1})^{-1} \Phi^T \mathbf{Y}_{\text{ini}},$$

where the regularization parameter  $\lambda$  is a very small positive number, e.g.,  $\lambda = 10^{-6}$  and  $\mathbf{I}_n$  is the  $n \times n$  identity matrix. The predicted error or residual matrix of this local model is

$$\mathbf{Y}_{\text{ini}} - \mathbf{f}_{\text{ini}}(\mathbf{X}_{\text{ini}}) \in \mathbb{R}^{W_G \times p}. \quad (3)$$

By shifting the data window one sample ahead, a new window  $\mathcal{W}_{\text{sft}} = \{\mathbf{X}_{\text{sft}}, \mathbf{Y}_{\text{sft}}\}$  is obtained, which contains the samples  $\{\mathbf{x}(t), \mathbf{y}(t)\}_{t=t_{\text{ini}}+1}^{t_{\text{ini}}+1+W_G}$ . If the two local data regions  $\mathcal{W}_{\text{ini}}$  and  $\mathcal{W}_{\text{sft}}$  are not significantly different, it can be considered that the data within  $\mathcal{W}_{\text{sft}}$  follow the same distribution as in  $\mathcal{W}_{\text{ini}}$  and the window continues to be shifted forward. Otherwise,  $\mathcal{W}_{\text{sft}}$  is considered to represent a new process state different from the one for  $\mathcal{W}_{\text{ini}}$ , and a new local linear model  $\mathbf{f}_{\text{new}}$  should be developed based on  $\mathcal{W}_{\text{sft}}$ . Let the estimation error matrix produced by  $\mathbf{f}_{\text{ini}}$  on  $\mathcal{W}_{\text{sft}}$  be denoted as

$$\mathbf{Y}_{\text{sft}} - \mathbf{f}_{\text{ini}}(\mathbf{X}_{\text{sft}}) \in \mathbb{R}^{W_G \times p}. \quad (4)$$

Whether the two local data regions  $\mathcal{W}_{\text{ini}}$  and  $\mathcal{W}_{\text{sft}}$  are similar or not can then be turned into the equivalent testing that tests whether  $\mathbf{Y}_{\text{ini}}$  and  $\mathbf{Y}_{\text{sft}}$  are significantly different or not. Since  $\mathbf{f}_{\text{ini}}$  is a  $p$ -output linear model,  $\mathbf{Y}_{\text{ini}}$  and  $\mathbf{Y}_{\text{sft}}$  are considered not significantly different when both their mean vectors and covariance matrices are the same. Accordingly, two null hypotheses  $H_0^\mu$  and  $H_0^\Sigma$  are set as

$$H_0^\mu : \mu_{\text{ini}} = \mu_{\text{sft}}, \quad (5)$$

$$H_0^\Sigma : \Sigma_{\text{ini}} = \Sigma_{\text{sft}}, \quad (6)$$

where the mean vectors of  $\mu_{\text{ini}}$  and  $\mu_{\text{sft}}$  as well as the covariance matrices of  $\Sigma_{\text{ini}}$  and  $\Sigma_{\text{sft}}$  come from the populations of  $\mathbf{Y}_{\text{ini}}$  and  $\mathbf{Y}_{\text{sft}}$ , respectively. Since  $\mathbf{f}_{\text{ini}}$  is an unbiased estimator,  $\mu_{\text{ini}} = \mathbf{0}_p$  and  $\Sigma_{\text{ini}} = \frac{1}{W_G} \mathbf{Y}_{\text{ini}}^T \mathbf{Y}_{\text{ini}}$ , where  $\mathbf{0}_p$  denotes the  $p$ -dimensional zero vector. In order to determine whether or not to accept  $H_0^\mu$  and  $H_0^\Sigma$ , two statistics  $F_\mu$  and  $F_\Sigma$  are constructed as [40]

$$F_\mu = \frac{(W_G - p)W_G}{(W_G - 1)p} (\hat{\mu}_{\text{sft}} - \mu_{\text{ini}})^T \hat{\Sigma}_{\text{sft}}^{-1} (\hat{\mu}_{\text{sft}} - \mu_{\text{ini}}), \quad (7)$$

$$F_\Sigma = (W_G - 1) \left( \ln |\Sigma_{\text{ini}}| - p - \ln |\hat{\Sigma}_{\text{sft}}| + \text{tr}(\hat{\Sigma}_{\text{sft}} \Sigma_{\text{ini}}^{-1}) \right), \quad (8)$$

where  $\hat{\mu}_{\text{sft}}$  and  $\hat{\Sigma}_{\text{sft}}$  are the mean vector and covariance matrix of  $\mathbf{Y}_{\text{sft}}$ , respectively, estimated with  $\hat{\mu}_{\text{sft}} = \frac{1}{W_G} \sum_{i=1}^{W_G} \mathbf{Y}_{\text{sft}}(i, :)$  and  $\hat{\Sigma}_{\text{sft}} = \frac{1}{W_G - 1} \sum_{i=1}^{W_G} (\mathbf{Y}_{\text{sft}}(i, :)) (\mathbf{Y}_{\text{sft}}(i, :)) - \hat{\mu}_{\text{sft}} \hat{\mu}_{\text{sft}}^T$  in which  $\mathbf{Y}_{\text{sft}}(i, :)$  denotes the  $i$ th row of  $\mathbf{Y}_{\text{sft}}$ , while  $|\bullet|$  and  $\text{tr}(\bullet)$  are the determinant and trace operators, respectively.

Assuming that both  $\mathbf{Y}_{\text{ini}}$  and  $\mathbf{Y}_{\text{sft}}$  follow the multivariate normal distributions, then  $F_\mu$  follows the  $F$  distribution with the degrees of freedom  $p$  and  $W_G - p$ , denoted as  $F_\mu \sim \mathcal{F}(p, W_G - p)$ , when  $H_0^\mu$  holds, while  $F_\Sigma$  follows the  $F$  distribution with the degrees of freedom  $f_1$  and  $f_2$ , denoted as  $F_\Sigma \sim b\mathcal{F}(f_1, f_2)$ , when  $H_0^\Sigma$  holds, where  $b = \frac{f_1}{1 - D_1 - f_1/f_2}$ ,  $f_1 = \frac{p(p+1)}{2}$ ,  $f_2 = \frac{f_1+2}{D_2 - D_1^2}$ ,  $D_1 = \frac{2p+1-2/(p+1)}{6(W_G-1)}$ , and  $D_2 = \frac{(p-1)(p+2)}{6(W_G-1)^2}$ . See for example [40]. Therefore, the condition of accepting both  $H_0^\mu$  and  $H_0^\Sigma$  are

$$F_\mu < \lambda_\mu \ \& \ F_\Sigma < \lambda_\Sigma, \quad (9)$$

where  $\lambda_\mu$  is the threshold value given the significance level  $\alpha_\mu$  which satisfies  $\Pr\{F_\mu < \lambda_\mu\} = 1 - \alpha_\mu$ , while  $\lambda_\Sigma$  is the threshold value given the significance level  $\alpha_\Sigma$  which satisfies  $\Pr\{F_\Sigma < \lambda_\Sigma\} = 1 - \alpha_\Sigma$ .

Let the local model set contains  $L > 1$  independent local linear models  $\{\mathbf{f}_l\}_{l=1}^L$ , and  $\mathbf{f}_{\text{ini}} = \mathbf{f}_L$ . When condition (9) is violated,  $\mathcal{W}_{\text{ini}}$  and  $\mathcal{W}_{\text{sft}}$  are significantly different, and the new local linear model  $\mathbf{f}_{\text{new}} = \mathbf{f}_{\text{sft}}$  is different from  $\mathbf{f}_L$ . We still need to test whether  $\mathbf{f}_{\text{new}}$  differs from  $\{\mathbf{f}_l\}_{l=1}^{L-1}$ . This task is also fulfilled based on the hypothesis testing. Let the prediction error matrices of  $\mathcal{W}_{\text{new}} = \{\mathbf{X}_{\text{sft}} \in \mathbb{R}^{W_G \times m}, \mathbf{Y}_{\text{sft}} \in \mathbb{R}^{W_G \times p}\}$  based on  $\mathbf{f}_{\text{new}}$  and  $\mathbf{f}_l$  be defined respectively by

$$\mathbf{Y}_{\text{new}} - \mathbf{f}_{\text{new}}(\mathbf{X}_{\text{sft}}) \in \mathbb{R}^{W_G \times p}, \quad (10)$$

$$\mathbf{Y}_l - \mathbf{f}_l(\mathbf{X}_{\text{sft}}) \in \mathbb{R}^{W_G \times p}, \ 1 \leq l \leq L - 1. \quad (11)$$

To test whether  $\mathbf{Y}_{\text{new}}$  and  $\mathbf{Y}_l$  are significantly different or not,

two null hypotheses  $H_0^{\mu_l}$  and  $H_0^{\Sigma_l}$  are set as

$$H_0^{\mu_l} : \mu_l = \mu_{\text{new}}, \quad (12)$$

$$H_0^{\Sigma_l} : \Sigma_l = \Sigma_{\text{new}}, \quad (13)$$

where the mean vectors  $\mu_l$  and  $\mu_{\text{new}}$  as well as the covariance matrices  $\Sigma_l$  and  $\Sigma_{\text{new}}$  come from the populations of  $\Upsilon_l$  and  $\Upsilon_{\text{new}}$ , respectively. Here,  $\mu_{\text{new}}$  and  $\Sigma_{\text{new}}$  are estimated based on  $\Upsilon_{\text{new}}$  as  $\mu_{\text{new}} = \mathbf{0}_p$  and  $\Sigma_{\text{new}} = \frac{1}{W_G} \Upsilon_{\text{new}}^T \Upsilon_{\text{new}}$ .

Again, to determine whether or not to accept  $H_0^{\mu_l}$  and  $H_0^{\Sigma_l}$ , two statistics are constructed as

$$F_{\mu}^{(l)} = \frac{(W_G - p)W_G}{(W_G - 1)p} (\hat{\mu}_l - \mu_{\text{new}})^T \hat{\Sigma}_l^{-1} (\hat{\mu}_l - \mu_{\text{new}}), \quad (14)$$

$$F_{\Sigma}^{(l)} = (W_G - 1) \left( \ln |\Sigma_{\text{new}}| - p \ln |\hat{\Sigma}_l| + \text{tr}(\hat{\Sigma}_l \Sigma_{\text{new}}^{-1}) \right), \quad (15)$$

where  $\hat{\mu}_l$  and  $\hat{\Sigma}_l$  are the mean vector and covariance matrix of  $\Upsilon_l$ , estimated with  $\hat{\mu}_l = \frac{1}{W_G} \sum_{i=1}^{W_G} \Upsilon_l(i, :)$  and  $\hat{\Sigma}_l = \frac{1}{W_G - 1} \sum_{i=1}^{W_G} (\Upsilon_l(i, :) - \hat{\mu}_l)^T (\Upsilon_l(i, :) - \hat{\mu}_l)$ . Under the assumption that  $\Upsilon_l$  and  $\Upsilon_{\text{new}}$  follow the multivariate normal distributions,  $F_{\mu}^{(l)} \sim \mathcal{F}(p, W_G - p)$  when  $H_0^{\mu_l}$  holds, and  $F_{\Sigma}^{(l)} \sim b\mathcal{F}(f_1, f_2)$  when  $H_0^{\Sigma_l}$  holds. Hence,  $f_l$  and  $f_{\text{new}}$  are regarded to be identical if the following condition is met

$$F_{\mu}^{(l)} < \lambda_{\mu} \ \& \ F_{\Sigma}^{(l)} < \lambda_{\Sigma}. \quad (16)$$

---

#### Algorithm 1 Adaptive local learning

---

- 1: **Initialization**
  - 2: Collect  $\mathcal{W}_{\text{ini}}$  with  $W_G$  consecutive samples from historical data, and construct multi-output LS linear model  $f_{\text{ini}}$  on  $\mathcal{W}_{\text{ini}}$ .
  - 3: Calculate  $\Upsilon_{\text{ini}}$ , and estimate  $\mu_{\text{ini}}$  and  $\Sigma_{\text{ini}}$ .
  - 4: Set  $L = 1$ ,  $\{\mathcal{W}_L, f_L\} = \{\mathcal{W}_{\text{ini}}, f_{\text{ini}}\}$  and  $\mathcal{W}_{\text{sft}} = \mathcal{W}_L$ .
  - 5: **Step 1: New local model detection**
  - 6: When a new data sample is available, shift  $\mathcal{W}_{\text{sft}}$  one sample ahead.
  - 7: Calculate  $\Upsilon_{\text{sft}}$ , and estimate  $\hat{\mu}_{\text{sft}}$  and  $\hat{\Sigma}_{\text{sft}}$ .
  - 8: Construct  $F_{\mu}$  and  $F_{\Sigma}$  statistics using (7) and (8).
  - 9: **If** condition (9) is satisfied
  - 10:   Go to **Step 1**.
  - 11: **End if**
  - 12: Construct multi-output LS linear model  $f_{\text{sft}}$  on  $\mathcal{W}_{\text{sft}}$ .
  - 13: Set  $\mathcal{W}_{\text{new}} = \mathcal{W}_{\text{sft}}$  and  $f_{\text{new}} = f_{\text{sft}}$ .
  - 14: Calculate  $\Upsilon_{\text{new}}$ , and estimate  $\mu_{\text{new}}$  and  $\Sigma_{\text{new}}$ .
  - 15: **Step 2: Redundant local model deletion**
  - 16: **For**  $l = 1, 2, \dots, L - 1$
  - 17:   Compute  $\Upsilon_l$ , and estimate  $\hat{\mu}_l$  and  $\hat{\Sigma}_l$ .
  - 18:   Construct  $F_{\mu}^{(l)}$  and  $F_{\Sigma}^{(l)}$  statistics using (14) and (15).
  - 19:   **If** condition (16) is satisfied
  - 20:     Delete  $f_l$ , set  $f_i = f_{i+1}$  for  $i = l, l + 1, \dots, L - 1$ , set  $L = L - 1$ , then go to **Step 3**.
  - 21:   **End if**
  - 22: **End for**
  - 23: **Step 3: Add new local model**
  - 24: Set  $L = L + 1$ ,  $\mathcal{W}_L = \mathcal{W}_{\text{new}}$  and  $f_L = f_{\text{new}}$ .
  - 25: Return to **Step 1**.
- 

Under this circumstance, either  $f_l$  or  $f_{\text{new}}$  is redundant and one of them should be removed. Since  $f_l$  is ‘older’ than the  $f_{\text{new}}$ ,  $f_{\text{new}}$  is kept and  $f_l$  is removed. On the other hand, if condition (16) is violated  $\forall l \in \{1, 2, \dots, L - 1\}$ ,  $f_{\text{new}}$  is different from  $f_l$  for  $1 \leq l \leq L$ . Thus, we have identified a new process state, and we add  $f_{\text{new}}$  to the local model set by setting  $L = L + 1$  and  $f_L = f_{\text{new}}$ .

The significance level  $\alpha_{\mu}$  and  $\alpha_{\Sigma}$  are usually set to small values, e.g., 0.05, 0.01, and they can be different according to the process data characteristics. Similar to our previous work [25], [26], the selection of window size  $W_G$  is a trade-off between the adaptive ability of capturing the local characteristics and the accuracy of local model. The proposed local learning procedure is summarized in Algorithm 1.

*Remark 1:* This local learning procedure can operate both offline and online. During online operation, when the newest data sample  $\{\mathbf{x}(t_{\text{next}}), \mathbf{y}(t_{\text{next}})\}$  is available, the data window shift one sample ahead, and the corresponding learning procedure can then be carried out. Unlike our previous work for single-output modeling [26], we consider multi-output modeling via multivariate statistics. This local learning procedure automatically encodes a newly emerging process state in the memory as a new local linear model. The local models in the memory are independent and represent different states of the process. Hence, our proposed learner is capable of achieving the desired maximum diversity and it is capable of acquiring all the different process states that have appeared.

#### B. Multi-output selective ensemble based online prediction

To mimic the adaptive intelligent decision making procedure of biological systems, at the second level of adaptation, our evolving model constructs a selective ensemble of the most relevant subset local linear models to produce an accurate prediction at each sample. Specifically, after the online operations at sample  $t$ , the local linear model set  $\{f_l\}_{l=1}^L$  has been produced by Algorithm 1, which represents the knowledge that the evolving system has acquired and stored in its memory. At the next sample of  $t_{\text{next}} = t + 1$ , the task of online modeling is to produce the model prediction  $\hat{\mathbf{y}}(t_{\text{next}})$  for the process’s true output  $\mathbf{y}(t_{\text{next}})$ , given the process input  $\mathbf{x}(t_{\text{next}})$  and the available local model set  $\{f_l\}_{l=1}^L$ . Our evolving model adopts an ensemble of the selected  $M$  local linear models from the model library  $\{f_l\}_{l=1}^L$  based on the  $q$  latest labeled data  $\{\mathbf{x}(t - i), \mathbf{y}(t - i)\}_{i=0}^{q-1}$ . Let the modeling error matrix of the  $l$ th multi-output local linear model  $f_l$  over the prediction window  $\{\mathbf{x}(t - i), \mathbf{y}(t - i)\}_{i=0}^{q-1}$  be defined as

$$\Upsilon_l(t) = [e_l(t) \ e_l(t - 1) \ \dots \ e_l(t - q + 1)]^T \in \mathbb{R}^{q \times p}, \quad (17)$$

where for  $0 \leq i \leq q - 1$ ,

$$\begin{aligned} e_l(t - i) &= \mathbf{y}(t - i) - f_l(\mathbf{x}(t - i)) \\ &= [e_{l,1}(t - i) \ e_{l,2}(t - i) \ \dots \ e_{l,p}(t - i)]^T. \end{aligned} \quad (18)$$

The performance metric of the  $l$ th local model is defined by

$$J_l(t) = \text{tr}(\Upsilon_l^T(t) \Upsilon_l(t)). \quad (19)$$

The sum of squared errors  $J_l(t)$  is further transformed into a probability metric. Specifically,  $J_l(t)$  is first converted to a

similarity measure [41] ranging from 0 to 1 as follows

$$\text{Sm}_l(t) = \frac{1}{1 + J_l(t)}. \quad (20)$$

The probability metric  $\text{Pr}_l(t)$  of the  $l$ th model is computed as the normalized similarity measure according to

$$\text{Pr}_l(t) = \frac{\text{Sm}_l(t)}{\sum_{i=1}^L \text{Sm}_i(t)}. \quad (21)$$

$\text{Pr}_l(t)$  quantifies the contribution of the  $l$ th local model to the ensemble, since a large  $\text{Pr}_l(t)$  indicates that the  $l$ th model is a good identifier for the ensemble model and vice versa.

Arrange all the  $L$  local models according to their probability values in descending order as

$$\text{Pr}_{l_1}(t) \geq \dots \geq \text{Pr}_{l_M}(t) \geq \text{Pr}_{l_{M+1}}(t) \geq \dots \geq \text{Pr}_{l_L}(t). \quad (22)$$

We select the first  $M$  best local models for constructing the ensemble model when the termination criterion

$$1 - \sum_{m=1}^M \text{Pr}_{l_m}(t) < \xi, \quad (23)$$

is met, where  $0 < \xi < 1$  is a desired tolerance. The selected models yield the  $M$  model outputs

$$\hat{\mathbf{y}}_{l_m}(t-i) = \mathbf{f}_{l_m}(\mathbf{x}(t-i)), \quad 1 \leq m \leq M, \quad (24)$$

for  $0 \leq i \leq q-1$ . Denote  $\mathbf{y}(t-i) = [y_1(t-i) \dots y_p(t-i)]^T$  and  $\hat{\mathbf{y}}_{l_m}(t) = [\hat{y}_{l_m,1}(t-i) \dots \hat{y}_{l_m,p}(t-i)]^T$ . The estimate  $\hat{y}_j(t-i)$  of the  $j$ th system output  $y_j(t-i)$  is given as the weighted sum of the  $M$  selected subset models, which is computed by

$$\hat{y}_j(t-i) = \sum_{m=1}^M \theta_{m,j}(t) \hat{y}_{l_m,j}(t-i), \quad 0 \leq i \leq q-1, \quad (25)$$

for  $1 \leq j \leq p$ , where the nonnegative  $\theta_{m,j}(t)$  is the combining coefficient for the  $m$ th selected local model of the  $j$ th output, and the combining coefficients must satisfy the constraint

$$\sum_{m=1}^M \theta_{m,j}(t) = 1, \quad 1 \leq j \leq p. \quad (26)$$

The  $j$ th estimation errors for  $1 \leq j \leq p$

$$\varepsilon_j(t-i) = y_j(t-i) - \hat{y}_j(t-i), \quad 0 \leq i \leq q-1, \quad (27)$$

are utilized to determine the combining coefficients of  $j$ th output. Specifically, the optimal combining coefficients can be obtained by minimizing the following cost function

$$V(t) = \frac{1}{2} \sum_{j=1}^p V_j(t) = \sum_{j=1}^p \sum_{i=0}^{q-1} \varepsilon_j^2(t-i). \quad (28)$$

Minimizing  $V(t)$  is equivalent to minimizing each  $V_j(t)$ ,  $1 \leq j \leq p$ , separately. Because of the constrain  $\sum_{m=1}^M \theta_{m,j}(t) = 1$ ,

$$\begin{aligned} V_j(t) &= \frac{1}{2} \sum_{i=0}^{q-1} \left( y_j(t-i) - \sum_{m=1}^M \theta_{m,j}(t) \hat{y}_{l_m,j}(t-i) \right)^2 \\ &= \frac{1}{2} \sum_{i=0}^{q-1} \left( \sum_{m=1}^M \theta_{m,j}(t) y_j(t-i) - \sum_{m=1}^M \theta_{m,j}(t) \hat{y}_{l_m,j}(t-i) \right)^2 \\ &= \frac{1}{2} \sum_{i=0}^{q-1} \left( \sum_{m=1}^M \theta_{m,j}(t) e_{l_m,j}(t-i) \right)^2 = \frac{1}{2} \boldsymbol{\theta}_j^T(t) \bar{\mathbf{E}}_j(t) \boldsymbol{\theta}_j(t), \end{aligned} \quad (29)$$

for  $1 \leq j \leq p$ , where  $\boldsymbol{\theta}_j(t) = [\theta_{1,j}(t) \dots \theta_{M,j}(t)]^T$  and  $\bar{\mathbf{E}}_j(t)$  is the estimated error covariance matrix of the  $j$ th output, which is given by

$$\bar{\mathbf{E}}_j(t) = \sum_{i=0}^{q-1} \begin{bmatrix} e_{l_1,j}^2(t-i) & \dots & e_{l_1,j}(t-i)e_{l_M,j}(t-i) \\ \vdots & \ddots & \vdots \\ e_{l_M,j}(t-i)e_{l_1,j}(t-i) & \dots & e_{l_M,j}^2(t-i) \end{bmatrix}. \quad (30)$$

Hence, we can form the following optimization problem to determine the optimal  $\boldsymbol{\theta}_j(t)$  for the  $j$ th output

$$\begin{aligned} \min_{\boldsymbol{\theta}_j} \quad & \frac{1}{2} \boldsymbol{\theta}_j^T(t) \bar{\mathbf{E}}_j(t) \boldsymbol{\theta}_j(t), \\ \text{s.t.} \quad & \sum_{m=1}^M \theta_{m,j}(t) = 1, \end{aligned} \quad (31)$$

where  $1 \leq j \leq p$ . The Lagrangian function for the single-output optimization (31) is given by

$$L(\boldsymbol{\theta}_j(t); \gamma) = \frac{1}{2} \boldsymbol{\theta}_j^T(t) \bar{\mathbf{E}}_j(t) \boldsymbol{\theta}_j(t) + \gamma (\mathbf{1}_M^T \boldsymbol{\theta}_j(t) - 1), \quad (32)$$

where  $\gamma > 0$  is a Lagrange multiplier. Letting  $\frac{\partial L}{\partial \boldsymbol{\theta}_j(t)} = \mathbf{0}_M$  yields

$$\bar{\mathbf{E}}_j(t) \boldsymbol{\theta}_j(t) + \gamma \mathbf{1}_M = \mathbf{0}_M, \quad 1 \leq j \leq p. \quad (33)$$

This suggests that the optimal combining vector  $\hat{\boldsymbol{\theta}}_j$  can be obtained as follows. First, calculate

$$\tilde{\boldsymbol{\theta}}_j(t) = [\tilde{\theta}_{1,j}(t) \dots \tilde{\theta}_{M,j}(t)]^T = \bar{\mathbf{E}}_j^{-1}(t) \mathbf{1}_M, \quad (34)$$

which is followed by the normalization

$$\hat{\theta}_{m,j}(t) = \frac{1}{\sum_{i=1}^M \tilde{\theta}_{i,j}(t)} \tilde{\theta}_{m,j}(t), \quad 1 \leq m \leq M. \quad (35)$$

The  $j$ th prediction  $\hat{y}_j(t_{\text{next}})$  for the  $j$ th system's true output

---

#### Algorithm 2 Multi-output selective ensemble prediction

---

- 1: **Initialization**
  - 2: Give  $W_G$ ,  $q$  and  $\xi$ .
  - 3: At beginning of online operation with sampling instance  $t$ , local model set  $\{\mathbf{f}_l\}_{l=1}^L$  has been constructed.
  - 4: **Step 1: Online prediction**
  - 5: Give input  $\mathbf{x}(t_{\text{next}})$  at new sample time  $t_{\text{next}} = t + 1$ .
  - 6: Calculate probability  $\text{Pr}_l(t)$  of each local model using (21) for  $1 \leq l \leq L$ .
  - 7: Select  $M$  subset models with termination criterion (23).
  - 8: Calculate error covariance matrix  $\bar{\mathbf{E}}_j(t)$ ,  $1 \leq j \leq p$ , for each output using (30).
  - 9: Calculate optimal combining coefficients  $\hat{\boldsymbol{\theta}}_j(t)$ ,  $1 \leq j \leq p$ , using (34) and (35).
  - 10: Predict true system outputs  $y_j(t_{\text{next}})$ ,  $1 \leq j \leq p$ , with selective ensemble prediction (36).
  - 11: **Step 2: Model adaptation**
  - 12: When  $\mathbf{y}(t_{\text{next}})$  is available, add  $\{\mathbf{x}(t_{\text{next}}), \mathbf{y}(t_{\text{next}})\}$  to dataset with  $t = t + 1$ .
  - 13: Carry out relevant operations to adapt local model set.
  - 14: Set  $t_{\text{next}} = t_{\text{next}} + 1$ , and go to **Step 1**.
-

$y_j(t_{\text{next}})$  is produced as the selected ensemble

$$\hat{y}_j(t_{\text{next}}) = \sum_{m=1}^M \hat{\theta}_{m,j}(t) f_{l_m,j}(\mathbf{x}(t_{\text{next}})), \quad 1 \leq j \leq p, \quad (36)$$

where  $f_{l_m,j}(\mathbf{x}(t))$  denotes the  $j$ th element of  $\mathbf{f}_{l_m}(\mathbf{x}(t))$ .

Algorithm 2 summarizes the multi-output selective ensemble based online prediction. In line 13 of Algorithm 2, the relevant operations may include lines 6 to 24 of adaptive local learning in Algorithm 1.

*Remark 2:* The two algorithmic parameters for selective ensemble prediction are the desired tolerance  $\xi$  and the number of latest labeled data  $q$ . A large  $q$  benefits the accuracy of the online SER prediction but imposes high computational complexity. A small  $q$  on the other hand offers high adaptability and is suitable for highly time-varying data. The threshold  $\xi$  trades off the accuracy of the SER prediction and the computational complexity. From all the  $L$  acquired independent local models that include the newest process state information, this SER procedure selects the most relevant subset of  $M$  local models to form an accurate prediction of the current process output. Obviously, the size of the selected ensemble  $M$  is different for different prediction samples. A very small positive regularization term  $\lambda_E$  can be added to the diagonal elements of  $\bar{\mathbf{E}}_j(t)$  for  $1 \leq j \leq p$  to ensure invertibility.

### C. Local model set pruning

In a fast changing environment, a biological system must be capable of removing out-of-date knowledge that are no longer relevant from its memory to free up memory space for fast acquiring new knowledge. This characteristic is also vital for our multi-output SER evolving model, because for highly nonstationary processes, the base local model set is likely to become very large over a long period of online adaptation, and this imposes high online computational complexity in constructing the SER predictor. Therefore, at the third level of adaptation, a pruning strategy is adopted to remove out-of-date local models from the memory and hence to alleviate the online computational burden. The essence of local model pruning is to remove those base local models that are far from the current data dynamics, since these local models are not needed in modeling the current process dynamics.

Recalling (17) to (23),  $M$  local linear models are selected to produce the ensemble prediction, which also suggests a way of local model set pruning. For example, the  $l_{M+1}$ th to  $l_L$ th models are not selected for the current prediction and thus they may be removed. However, pruning a model based on its ‘one-sample’ prediction horizon may not be sufficiently reliable, as these unselected models at the current sample may be important at  $t_{\text{next}} + 1$ . To make the pruning more reliable, the work [26] introduced the concept of ‘memory depth’ for an ensemble learner. In the local learning procedure, a local linear model is constructed based on a data window with the window size  $W_G$ . Within this data window, the process’s dynamics are assumed to be stationary. Similarly, we introduce a data window for model pruning with the window size  $W_P$ . Specifically, if a local model is never selected for

---

### Algorithm 3 Multi-output local model pruning

---

- 1: **Initialization**
  - 2: Give  $W_P$ , set counters of all local models  $\text{count}_l = 0$  for  $1 \leq l \leq L$ , set  $t = t_{\text{ini}}$  and  $\text{index} = 0$ .
  - 3: **Step 1: Pruning in pruning model window**
  - 4: Perform multi-output selective ensemble prediction.
  - 5: **If**  $(t - t_{\text{ini}} \leq W_P)$
  - 6:   **For**  $l = 1, 2, \dots, L$
  - 7:     **If**  $f_l$  is not selected at current sample  $t$
  - 8:        $\text{count}_l = \text{count}_l + 1$ .
  - 9:     **End if**
  - 10:   **End for**
  - 11:   Set  $t = t + 1$  and go to **Step 1**.
  - 12: **Else**
  - 13:   **For**  $l = 1, 2, \dots, L$
  - 14:     **If**  $\text{count}_l = W_G$
  - 15:       Add  $l$  to pruning model index set  $\Gamma$ , and set  $\text{index} = \text{index} + 1$ .
  - 16:     **End if**
  - 17:   **End for**
  - 18:   Delete  $f_l$  for all  $l \in \Gamma$ , and set  $L = L - \text{index}$ .
  - 19: **End if**
  - 20: **Step 2: Pruning model window update**
  - 21: Clear counters for all local models, set  $t_{\text{ini}} = t$  and  $\text{index} = 0$ , and go to **Step 1**.
- 

the consecutive  $W_P$  samples of the window, then it can be removed with high confidence.

The multi-output local model pruning strategy is listed in Algorithm 3. Line 4 of Algorithm 3 corresponds to the SER prediction operations of lines 6 to 10 in Algorithm 2. The rest pruning operations in Algorithm 3 occur at line 13 of Algorithm 2. Since we have a window size  $W_G$  for adaptive local learning, we can conveniently set  $W_P = W_G$ .

*Remark 3:* Since the newest local model  $f_L$  represents the newly emerged process state, it is critic to keep it in the current prediction horizon. Thus we always set  $\text{count}_L$  to zero. To maintain the local model set diversity and hence the identification accuracy, a minimal number of local model  $L_{\min}$  should be guaranteed. If the number of those unselected local models in  $\Gamma$  exceeds  $L - L_{\min}$ , only the oldest  $L - L_{\min}$  local models can be removed. Let  $L_{\text{ini}}$  be the number of local linear models obtained in the initial training. We can set  $L_{\min} = L_{\text{ini}}$ . This pruning strategy provides good plasticity to our multi-output SER learner and it can dramatically reduce online computational complexity in adaptive SER prediction.

## III. CASE STUDIES

The performance of the proposed multi-output SER learning method is evaluated using three case studies, which are: a simulated two-output nonlinear system, a non-isothermal CSTR process [37], [38] and an industrial microwave heating system [39]. The testing prediction error covariance  $\text{Cov}(\mathbf{E})$  is used to evaluate the multi-output online modeling performance, where  $\text{Cov}(\mathbf{E}) = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{E}_i - \bar{\mathbf{E}})(\mathbf{E}_i - \bar{\mathbf{E}})^T$ ,  $N$  is the number of test samples and  $\bar{\mathbf{E}}$  is the sample average of  $\mathbf{E}$ . The

online computational complexity is quantified by its averaged computation time per sample (ACTpS). The experiments are carried out on Matlab 2017a, running on a PC with i7-3770 3.40GHz processor of 4 cores and 16GB of RAM.

The performance of our proposed method is compared with existing typical online learning approaches, including the multi-output OLS (MOLS) [32]–[35], the OS-ELM [14], the EOS-ELM [24] and the recently proposed single-output GAP-SER [26], which is denoted as the S-SER. For the MOLS, nonlinear modeling is achieved by using the Gaussian RBF kernel and the initial training is carried out by constructing a compact multi-output RBF model using the OLS learning algorithm. During online operation, the weight adaptation is performed by the RLS algorithm with the fixed model structure. For the OS-ELM, a RBF network is initialized during training by randomly selecting a large number of input data points as its RBF centers, and the online adaptation of the OS-ELM involves the weight updating using the RLS algorithm. Similarly, for the EOS-ELM, a number of OS-ELM base models are trained on the same dataset, and each base model's centers are randomly selected from the training data points as in OS-ELM. The forgetting factor of the RLS algorithm is set to 0.98. For the S-SER, the multiple single-output SER models are constructed using the single-output GAP-SER of [26], one for modeling an output of the system.

#### A. Simulated two-output nonlinear system

The simulated two-output nonlinear system of [33], [34] is considered. The data set contained 1,000 noisy observations generated using the model

$$\begin{aligned} y_1(t) &= 0.5y_1(t-1) + u(t-1) + 0.4 \tanh(u(t-2)) \\ &\quad + 0.1 \sin(\pi y_1(t-2))y_2(t-1) + \epsilon_1(t), \\ y_2(t) &= 0.3y_2(t-1) + 0.1y_2(t-2)y_1(t-1) \\ &\quad + 0.4 \exp(-u^2(t-1))y_1(t-2) + \epsilon_2(t), \end{aligned} \quad (37)$$

where the system input  $u(t)$  is uniformly distributed in  $[-0.5, 0.5]$ , and the zero-mean Gaussian noise  $\epsilon(t) = [\epsilon_1(t) \ \epsilon_2(t)]^T$  has a covariance matrix  $0.04\mathbf{I}_2$ . Initial conditions are set as  $y_1(0) = y_2(0) = y_1(-1) = y_2(-1) = 0$  and  $u(0) = u(-1) = 0$ . The first 5,00 data generated from (37)

are used for training and the rest 5,00 samples for testing. Note that our proposed approach as well as the MOLS and the S-SER do not really need such a large number of training samples but the ELM-based models need a large number of training samples, as an ELM model must contain a large number of hidden nodes. In modeling this system, we use the model input vector as

$$\mathbf{x}(t) = [y_1(t-1) \ y_1(t-2) \ y_2(t-1) \ y_2(t-2) \ u(t-1) \ u(t-2)]^T \in \mathbb{R}^6. \quad (38)$$

In the simulation, 100 independent realizations are generated. The performance of each method are presented by its means and standard deviations (STDs) of the test  $\log(\det(\text{Cov}(\mathbf{E})))$  and ACTpS, respectively, over the 100 realizations.

For our proposed model, the window size, innovation length and decision threshold are empirically chosen to be  $W_G = 30$ ,  $q = 10$  and  $\xi = 0.9$ , respectively. For the S-SER, three algorithmic parameters are chosen to be  $W_G = 20$ ,  $q = 10$  and  $\xi = 0.1$  to achieve the best prediction performance. Note that we do not attempt to optimize the algorithmic parameters in training, as such an optimal algorithmic setting is only meaningful when the underlying system is stationary. For online identification of a nonstationary system, the optimal algorithmic setting is fast time-varying, and it is prohibitive to online learn this fast time-varying optimal algorithmic setting. Since each algorithmic parameter has clear physical interpretation, the appropriate value can be set empirically. The detailed analysis can be found in our previous work [25], [26].

The test results of all the five methods, including the model size, the ACTpS, the error covariance  $\log(\det(\text{Cov}(\mathbf{E})))$ , and the average ensemble size, are summarized in Table I. This system is not seriously time-varying, and all the methods are expected to perform well. This is confirmed by Table I. More specifically, the MOLS with 10 RBF nodes outperforms the OS-ELM with 100 nodes and the EOS-ELM of 5 base models with each base model having 100 nodes, in terms of both online modeling accuracy and computational complexity, and its performance is even slightly better than the S-SER. Among all the models, our proposed method attains the

TABLE I  
SIMULATED TWO-OUTPUT NONLINEAR SYSTEM: COMPARISON OF ONLINE PREDICTION AND ADAPTIVE MODELING PERFORMANCE (AVERAGE $\pm$ STD)  
FOR THE OS-ELM, EOS-ELM, MOLS, S-SER AND THE PROPOSED METHOD

Algorithm	$\log(\det(\text{Cov}(\boldsymbol{E})))$	ACTpS (ms)		Models/Nodes		Average ensemble size
				Initial	Final	
OS-ELM	-2.31±0.07	0.27±0.01		100	100	-
	-1.36±0.07	6.63±0.11		500	500	-
EOS-ELM	-2.32±0.04	2.30±0.05		100	100	5
	-1.46±0.16	58.71±2.37		500	500	5
MOLS	-2.60±0.04	0.04±0.01		10	10	-
	-2.59±0.05	0.06±0.01		20	20	-
	-2.54±0.05	0.08±0.01		30	30	-
S-SER	<b>-2.46±0.11</b>	SO1	0.15±0.02	2.13±1.01	3.02±1.39	2.71±1.19
		SO2	0.15±0.02	2.01±1.01	2.88±1.20	2.52±1.04
		MO	<b>0.30±0.04</b>	<b>4.14±1.35</b>	<b>5.90±1.87</b>	<b>5.23±1.61</b>
Proposed	<b>-2.72±0.06</b>	<b>0.47±0.04</b>		<b>4.99±1.57</b>	<b>5.12±1.57</b>	<b>1±0</b>



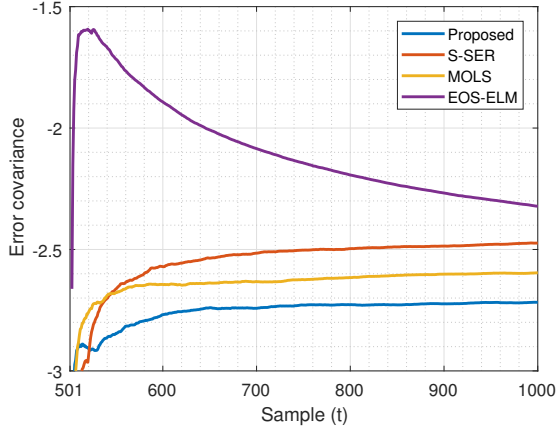


Fig. 1. Comparison of average  $\log(\det(\text{Cov}(\mathbf{E})))$  learning curves of various models for the simulated nonlinear system. The EOS-ELM has total of 500 hidden nodes, and the MOLS has 10 hidden nodes. The OS-ELM has a very close learning curve with the EOS-ELM and is omitted here.

lowest  $\log(\det(\text{Cov}(\mathbf{E})))$  and its average SER size is much smaller than that of the multiple ( $p=2$ ) S-SER models. The comparison of average error covariance learning curves for various models is depicted in Fig. 1.

### B. Non-isothermal CSTR process

The non-isothermal CSTR with an irreversible reaction ( $A \rightarrow B$ ) is widely used as a benchmark for nonlinear and time-varying process modeling and identification [37], [38]. Based on the underlying physical-chemical laws, it is well known that this process is governed by the following two nonlinear ordinary differential equations [42]

$$\begin{aligned} \frac{dC_A}{dt} &= \frac{q_f}{V}(C_{Af} - C_A) - K_0 C_A \exp\left(-\frac{E}{RT}\right) \phi_c(t), \\ \frac{dT}{dt} &= \frac{q_f}{V}(T_f - T) + \frac{(-\Delta H)K_0 C_A}{\rho C_p} \exp\left(-\frac{E}{RT}\right) \phi_c(t) \\ &\quad + \frac{\rho_c C_{pc}}{\rho C_p V} q_c \left(1 - \exp\left(-\frac{h_A}{q_c \rho C_{pc}} \phi_h(t)\right)\right) (T_{cf} - T), \end{aligned} \quad (39)$$

where  $\phi_h(t)$  is the fouling coefficient,  $\phi_c(t)$  is the deactivation coefficient,  $C_A$  is the effluent concentration (controlled variable),  $T$  is the output temperature (controlled variable),  $q_c$  is the coolant flow rate (manipulated variable),  $q_f$  is the feed flow rate (manipulated variable),  $C_{Af}$  is the feed concentration,

TABLE II  
NOMINAL CSTR OPERATING CONDITION

$h_A = 7 \times 10^5 \text{ cal} \cdot \text{min}^{-1} \text{ K}^{-1}$	$T = 440.2 \text{ K}$
$K_0 = 7.2 \times 10^{10} \text{ min}^{-1}$	$T_f = 350 \text{ K}$
$E/R = 9.95 \times 10^3 \text{ K}$	$T_{cf} = 350 \text{ K}$
$-\Delta H = 2 \times 10^5 \text{ cal} \cdot \text{mol}^{-1}$	$C_{Af} = 1 \text{ mol} \cdot \text{l}^{-1}$
$C_A = 8.36 \times 10^{-2} \text{ mol} \cdot \text{l}^{-1}$	$V = 100 \text{ l}$
$C_p, C_{pc} = 1 \text{ cal} \cdot \text{g}^{-1} \text{ K}^{-1}$	$\rho_c, \rho = 1000 \text{ g} \cdot \text{l}^{-1}$
$q_c = 103.41 \text{ l} \cdot \text{min}^{-1}$	$q_f = 100 \text{ l} \cdot \text{min}^{-1}$

$T_f$  is the feed temperature and  $T_{cf}$  is the coolant inlet temperature. The other process parameters, together with the operating conditions, are given in Table II.

Accurate simulator of the non-isothermal CSTR is built by first-order differencing the ordinary differential equation (39) [37], [38]. To obtain the multiple operating conditions,  $\pm 5\%$  step changes are added to both coolant ( $q_c$ ) and feed ( $q_f$ ) flow rates. The step size is set to be 0.1. We follow the same approach to generate the process input-output data. The process inputs are  $q_f$  and  $q_c$ , and its outputs are  $C_A$  and  $T$ . For this process, the system input vector is chosen as

$$\mathbf{x}(t) = [C_A(t-1) \ T(t-1) \ q_f(t-1) \ q_c(t-1)]^T \in \mathbb{R}^4, \quad (40)$$

and the system output vector is

$$\mathbf{y}(t) = [C_A(t), T(t)]^T \in \mathbb{R}^2. \quad (41)$$

1) *Constant deactivation and fouling coefficients:* We first consider the case that the deactivation coefficient  $\phi_c(t)$  and the fouling coefficient  $\phi_h(t)$  are both constant, which are set to 1. Then 991 samples are generated with the first 5,00 samples for training and the remaining 491 data samples for testing. For our proposed method, the three algorithmic parameters are empirically chosen to be  $W_G = 100$ ,  $q = 20$  and  $\xi = 0.9$ , respectively. For each single-output model of the S-SER, the three algorithmic parameters are chosen to be  $W_G = 10$ ,  $q = 20$  and  $\xi = 0.9$  to achieve its best modeling accuracy.

The test performance of the five models are compared in Table III. The OS-ELM with 500 nodes and the EOS-ELM with 5 base models each having 500 nodes attain a similar performance of  $\log(\det(\text{Cov}(\mathbf{E}))) = -10.07$ , but they impose

TABLE III  
NON-ISOTHERMAL CSTR PROCESS WITH TWO CONSTANT PROCESS PARAMETERS: COMPARISON OF ONLINE PREDICTION AND ADAPTIVE MODELING PERFORMANCE FOR THE OS-ELM, EOS-ELM, MOLS, S-SER AND THE PROPOSED METHOD

Algorithm	$\log(\det(\text{Cov}(\mathbf{E})))$	ACTpS (ms)		Models/Nodes		Average ensemble size
				Initial	Final	
OS-ELM	-5.3252	0.3239		100	100	-
	-10.0731	6.5782		500	500	-
EOS-ELM	-5.2789	2.4132		100	100	5
	-10.0732	58.5402		500	500	5
MOLS	-7.3414	0.0716		10	10	-
	-9.1829	0.0942		20	20	-
	-7.1830	0.1069		30	30	-
S-SER	<b>-10.4252</b>	SO1	0.5562	65	65	5.66
		SO2	0.4801	62	62	6.01
		MO	<b>1.0372</b>	<b>127</b>	<b>127</b>	<b>11.67</b>
Proposed	<b>-12.5867</b>	<b>2.7334</b>		<b>45</b>	<b>53</b>	<b>5.47</b>



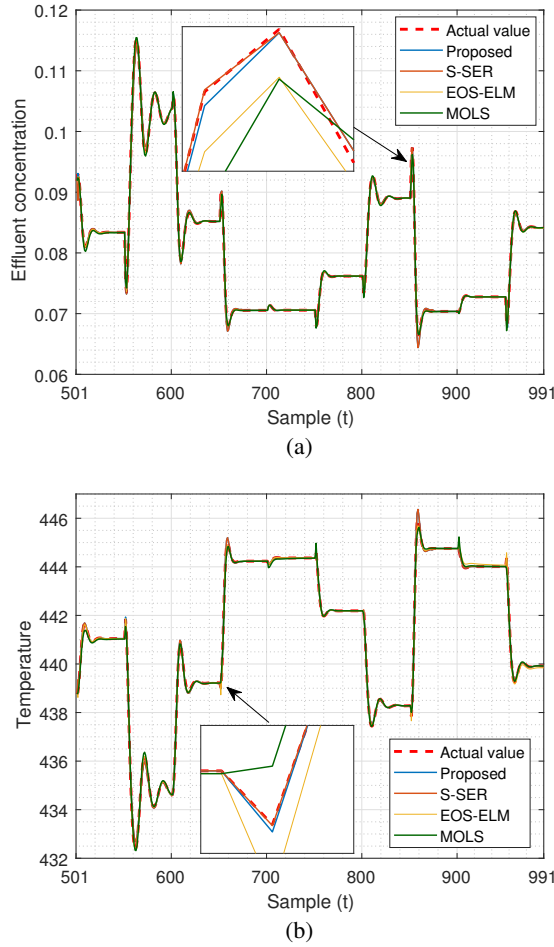


Fig. 2. Online identification of CSTR process with two constant process parameters: (a) effluent concentration, and (b) output temperature.

very high ACTpS of 6.58 ms and 58.54 ms, respectively. The MOLS with 20 nodes has a slightly worst accuracy but imposes the lowest ACTpS of 0.09 ms. The S-SER approach attains a better accuracy than the OS-ELM and EOS-ELM while imposing a significantly lower ACTpS of 1.04 ms. Our proposed method achieves the most accurate prediction with an ACTpS of 2.73 ms. The online prediction values of the MOLS with 20 nodes, the EOS-ELM with total of 2500 nodes, the

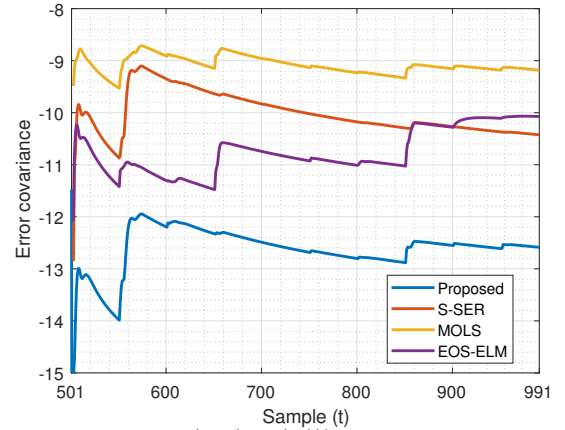


Fig. 3. Comparison of  $\log(\det(\text{Cov}(\mathbf{E})))$  learning curves of various models for the CSTR process with two constant process parameters. The EOS-ELM has total of 2500 hidden nodes, and the MOLS has 20 hidden nodes.

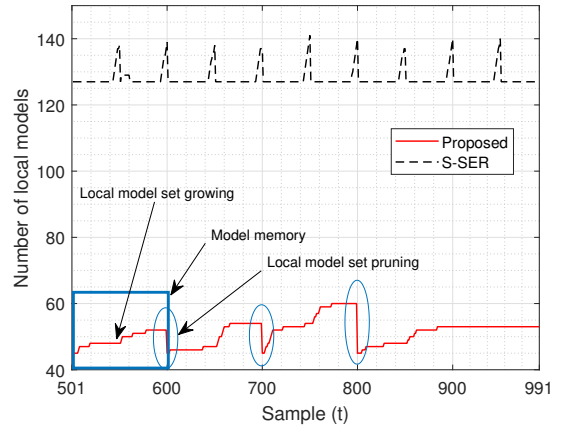


Fig. 4. Comparison of online local model set learning curves of the two SER methods for the CSTR process with two constant process parameters.

S-SER and our proposed method are compared with the actual process output observations in Fig. 2. Fig. 3 depicts the error covariance learning curves of various models.

From Table III, it can be seen that the online modeling efforts of our proposed multi-output SER method are lower than those of the S-SER method. Specifically, it imposes smaller model set size and smaller average ensemble, compared with the multiple S-SER models. This is confirmed by

TABLE IV  
NON-ISOTHERMAL CSTR PROCESS WITH TWO TIME-VARYING PARAMETERS: COMPARISON OF ONLINE PREDICTION AND ADAPTIVE MODELING PERFORMANCE FOR THE OS-ELM, EOS-ELM, MOLS, S-SER AND THE PROPOSED METHOD

Algorithm	$\log(\det(\text{Cov}(\boldsymbol{E})))$	ACTpS (ms)		Models/Nodes		Average ensemble size
				Initial	Final	
OS-ELM	-6.9808	0.3445		100	100	-
	-9.8185	6.5982		500	500	-
EOS-ELM	-6.9808	2.5011		100	100	5
	-9.9095	58.5402		500	500	5
MOLS	-8.0568	0.0796		10	10	-
	-8.4489	0.1036		20	20	-
	-8.4299	0.1088		30	30	-
S-SER	<b>-14.2614</b>	SO1	0.7722	40	40	14.83
		SO2	0.3797	27	27	10.45
		MO	<b>1.1519</b>	<b>67</b>	<b>67</b>	<b>25.29</b>
Proposed	<b>-21.0996</b>	<b>3.5615</b>		<b>47</b>	<b>57</b>	<b>5.56</b>

the online local model set learning curves of the two SER methods shown in Fig. 4. Observe the red curve in Fig. 4 that our proposed growing and pruning strategies can effectively add new local linear models and remove out-of-date ones, respectively, within the model memory  $W_G = 100$  (blue rectangle). However, the recorded ACTpS of the proposed method is higher than that of the S-SER. This is owing to Matlab software implementation. The high-dimensional matrix calculation in Matlab is complicated and time-consuming. In real-time operation with other programming platform, we expect that the online ACTpS of the proposed method will be significantly lower than that of the S-SER.

2) *Time-varying deactivation and fouling coefficients:* We next consider time-varying  $\phi_c(t)$  and  $\phi_h(t)$ . Fouling and catalyst deactivation phenomena present two main sources to produce nonstationary characteristics in the CSTR process during its normal operation. These time-varying dynamic characteristics are governed by the following equations

$$\phi_h(t) = 1 - 0.01t, \quad \phi_c(t) = \exp\left(-0.00067 \frac{E}{RT} t\right), \quad (42)$$

where the fouling effect is caused by depositing material on the heat transfer surface, and catalyst is deactivated due to poisoning. Again, 991 samples are generated with the first 5,00 data for training and the rest 491 data for testing.

Table IV compares the performance of different models.

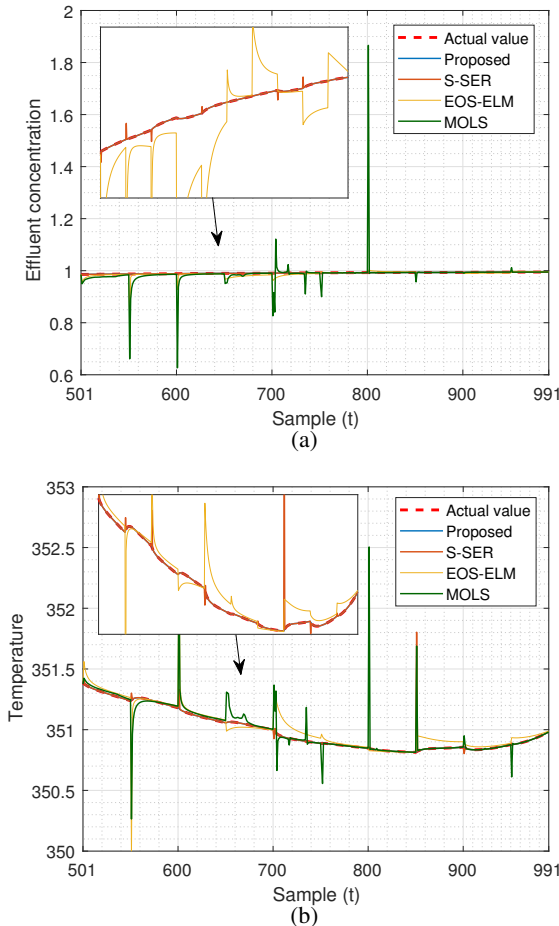


Fig. 5. Online identification of CSTR process with two time-varying process parameters: (a) effluent concentration, and (b) output temperature.

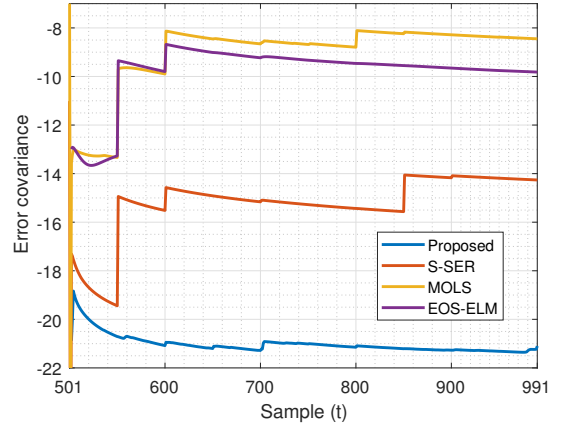


Fig. 6. Comparison of  $\log(\det(\text{Cov}(\mathbf{E})))$  learning curves of various models for the CSTR process with two time-varying process parameters. The EOS-ELM has total of 2500 hidden nodes, and the MOLS has 20 hidden nodes.

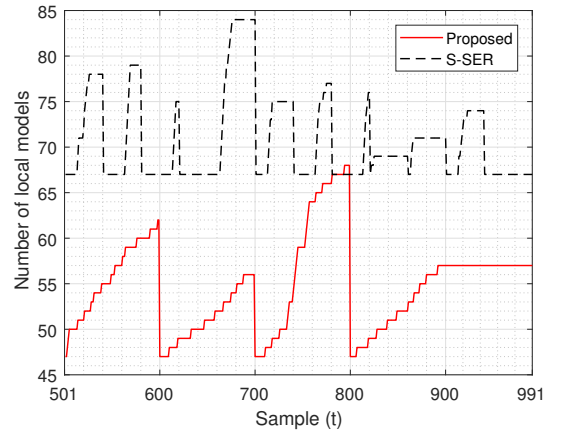


Fig. 7. Comparison of online local model set learning curves of the two SER methods for the CSTR process with two time-varying process parameters.

For this seriously nonstationary process, the S-SER significantly outperforms the OS-ELM and EOS-ELM, in terms of both online prediction accuracy and computational complexity. Furthermore, our proposed method considerably outperforms the S-SER in online prediction accuracy. The MOLS has the poorest accuracy but imposes the lowest ACTpS. Online prediction values by the various models are compared with the actual system output observations in Fig. 5, while Fig. 6 shows the error covariance learning curves of various models. Again although the recorded ACTpS of the proposed method is higher than that of the SER, it actually requires smaller model set size and much smaller average ensemble than the multiple S-SER models, as confirmed by the online local model set learning curves of the two SER methods given in Fig. 7.

### C. Microwave heating process

Microwave heating technology has found wide-ranging applications due to its many advantages over conventional heating methods, including selective and volumetric heating, rapid heat transfer and pollution-free environment. However, a major drawback of microwave heating is the temperature runaway, caused by properties of material and the inner electromagnetic field distribution, which may lead to unwanted combustion and destruction in industrial processes [39]. Therefore, in the

TABLE V  
MICROWAVE HEATING PROCESS: COMPARISON OF ONLINE PREDICTION AND ADAPTIVE MODELING PERFORMANCE FOR THE OS-ELM, EOS-ELM, MOLS, S-SER AND THE PROPOSED METHOD

Algorithm	$\log(\det(\text{Cov}(\boldsymbol{E})))$	ACTpS (ms)		Models/Nodes		Average ensemble size
				Initial	Final	
OS-ELM	-10.4544	0.3261		100	100	-
	-13.3419	6.5582		500	500	-
EOS-ELM	-10.4561	2.2521		100	100	5
	-13.3418	57.5702		500	500	5
MOLS	-10.3060	0.0499		5	5	-
	-12.3751	0.0507		10	10	-
	-3.8084	0.0571		15	15	-
S-SER	<b>-12.7403</b>	SO1	0.2629	11	11	5.69
		SO2	0.2414	10	10	5.20
		SO3	0.3807	23	23	11.56
		MO	<b>0.8850</b>	<b>44</b>	<b>44</b>	<b>22.46</b>
Proposed	<b>-14.1326</b>	<b>1.4401</b>		<b>25</b>	<b>28</b>	<b>2</b>

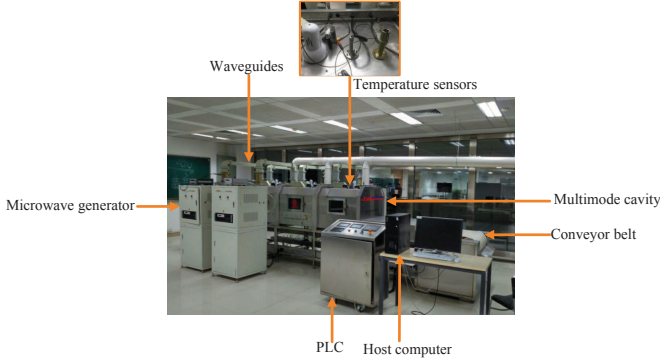


Fig. 8. An industrial microwave heating system.

control of the microwave heating process (MHP), material temperature is an important feedback information. However, establishing a temperature model from the first-principle is challenging, since the MHP involves multi-physical fields coupling. Hence, assumptions have to be made and, consequently, unmodeled dynamics always exist in a model derived from the first-principle [43], [44]. This motivates us to investigate data-driven predictive model [45]–[47].

A distributed microwave heating system [48], [49] is shown in Fig. 8. It consists of five microwave generators and waveguides. Microwave generated by each microwave generator is transmitted through the corresponding waveguide, fed into the cavity and absorbed by the heated material. The material is continuously transported through cavity by the conveyor belt, whose speed can be adjusted by a motor driver. Three fiber optical sensors (FOSs), denoted as FOS1 to FOS3, are placed at three different locations to online record multiple-points of temperature. During the realtime operation of this MHP, the control center receives the measured temperature values from the FOSs, and sends control commands, including the five microwave powers  $u_{p_i}(t)$ ,  $1 \leq i \leq 5$ , for the five microwave generators as well as the conveyor speed  $v(t)$  to the cavity. Thus, the control inputs to this MHP are given by

$$\mathbf{u}(t) = [u_{p_1}(t) \ u_{p_2}(t) \ u_{p_3}(t) \ u_{p_4}(t) \ u_{p_5}(t) \ v(t)]^T. \quad (43)$$

Each FOS measures the temperature at its location, which

is the MHP's output  $y_i(t)$ ,  $1 \leq i \leq 3$ . Because of near instantaneous response of MHP, the process's temperatures  $\mathbf{y}(t) = [y_1(t) \ y_2(t) \ y_3(t)]^T$  can be adequately modeled as

$$\mathbf{y}(t) = \mathbf{f}_{nn}(\mathbf{x}(t); t), \quad (44)$$

where  $\mathbf{f}_{nn}(\cdot; t)$  represents the unknown nonlinear time-varying system mapping with the input vector given by

$$\mathbf{x}(t) = [\mathbf{y}^T(t-1) \ \mathbf{u}^T(t-1)]^T \in \mathbb{R}^9. \quad (45)$$

3,000 process data have been collected from this distributed microwave heating system. We first normalize the microwave power inputs and the temperature measurements according to

$$\bar{u}_{p_i}(t) = \frac{u_{p_i}(t)}{1000}, \quad 1 \leq i \leq 5, \quad (46)$$

$$\bar{y}_i(t) = \frac{y_i(t) - y_{i_{\min}}}{y_{i_{\max}} - y_{i_{\min}}}, \quad 1 \leq i \leq 3, \quad (47)$$

where  $y_{i_{\min}}$  and  $y_{i_{\max}}$  are the minimum and maximum temperatures recorded by the  $i$ th FOS, respectively. We use the first 5,00 samples for training, and the last 2,500 samples for online prediction and adaptive modeling. For our proposed method, the three algorithmic parameters are empirically chosen to be  $W_G = 110$ ,  $q = 25$  and  $\xi = 0.95$ , respectively. For each S-SER model, the three algorithmic parameters are chosen to be  $W_G = 20$ ,  $q = 25$  and  $\xi = 0.5$ .

Table V compares the performance of different predictive models. The predictor outputs versus the actual process outputs by various models are depicted in Fig. 9, where the EOS-ELM has total of 2500 hidden nodes and the MOLS has 10 hidden nodes. Furthermore, the comparison of error covariance learning curves is given in Fig. 10. In this case, the OS-ELM with 500 nodes and the EOS-ELM with total of 2500 nodes outperform both the S-SER and MOLS, but they impose extremely higher online computational complexity. In particular, the EOS-ELM with total of 2500 nodes requires the highest ACTpS of 57.57 ms. By contrast, our proposed multi-output SER achieves the best prediction performance with a reasonable online computation time of ACTpS = 1.44 ms. The comparison of online local model set learning curves of the

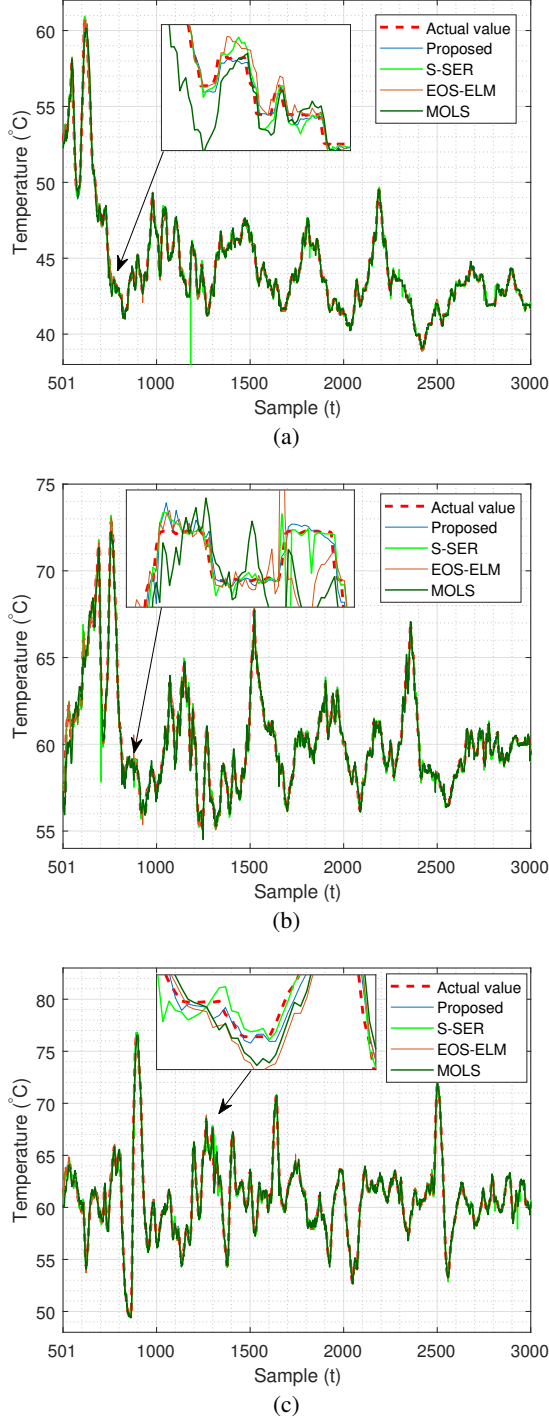


Fig. 9. Online identification of microwave heating process by various models: (a) FOS1, (b) FOS2, and (c) FOS3.

two SER methods is given in Fig. 11, which also confirms that our proposed method needs a smaller model set and smaller average ensemble, compared with the three S-SER models.

#### IV. CONCLUSIONS

In this paper, a novel multi-output SER evolving model has been developed for multi-output nonlinear and nonstationary process identification. Our proposed online adaptive learner has been inspired by the fundamental biological system learning principle, namely, the ability to acquire new knowledge

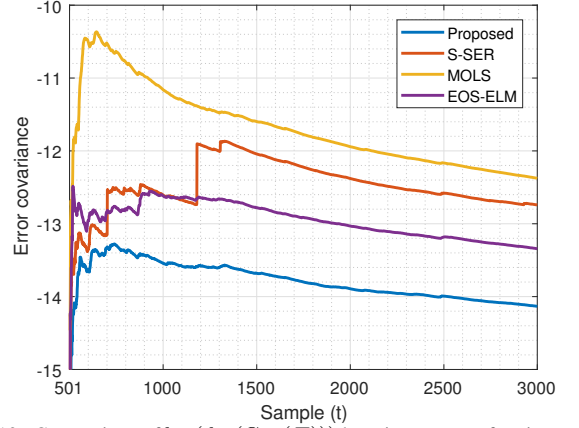


Fig. 10. Comparison of  $\log(\det(\text{Cov}(\mathbf{E})))$  learning curves of various models for MHP identification. The EOS-ELM has total of 2500 hidden nodes, while the MOLS has 10 hidden nodes.

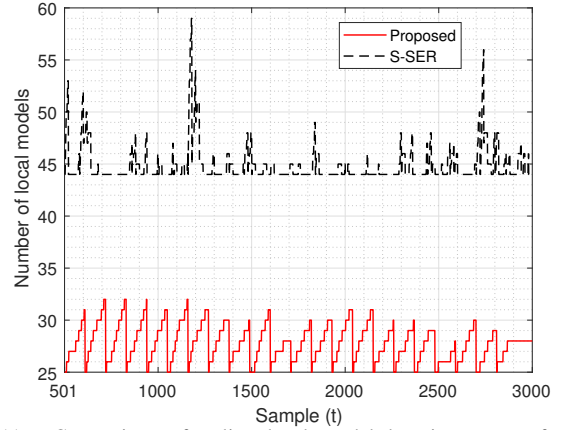


Fig. 11. Comparison of online local model learning curves for MHP identification.

to the memory and to remove out-of-date knowledge from the memory so that intelligent decision can be made based on the most new and relevant knowledge in the memory. Hence, our contribution has been threefold. First, our adaptive local learning enables automatically identifying every newly emerging process state and constructing a matching local multi-output linear model via multivariate statistic hypothesis testing. Second, the optimal online prediction of the system's multi-outputs is obtained by a selective ensemble of the most relevant subset local linear models. Third, an effective pruning strategy removes the most out-of-date local linear models that are no longer needed in modeling the process in order to free up the memory for fast acquiring the new process knowledge. This pruning strategy also significantly reduces online computational complexity without scarifying the prediction accuracy. Extensive studies have been conducted, including a simulated nonlinear system, a simulator based CSTR process and a real-world microwave heating system identification. The results obtained have demonstrated that our proposed multi-output selective ensemble identification technique attains the best online modeling accuracy, compared with a range of the state-of-art methods for online identification of nonlinear and nonstationary multi-output processes, while imposing a reasonably low online computational complexity which meets

the real-time operation constraint. This study therefore has provided a reliable and accurate online system model for designing efficient real-time control strategy for multi-output nonlinear and nonstationary processes.

## REFERENCES

- [1] D. E. Goldberg, *Genetic Algorithms in Search, Optimization and Machine Learning*. Reading, MA, USA: Addison-Wesley, 1989.
- [2] K. Price, R. Storn, and J. Lampinen, *Differential Evolution: A Practical Approach to Global Optimization*. Berlin, Germany: Springer-Verlag, 2005.
- [3] Z. Zhou, N. V. Chawla, Y. Jin, and G. J. Williams, "Big data opportunities and challenges: Discussions from data analytics perspectives," *IEEE Computational Intelligence Magazine*, vol. 9, no. 4, pp. 62–74, Nov. 2014.
- [4] G. Ditzler, M. Roveri, C. Alippi, and R. Polikar, "Learning in non-stationary environments: A survey," *IEEE Computational Intelligence Magazine*, vol. 10, no. 4, pp. 12–25, Oct. 2015.
- [5] L. Rutkowski, "Generalized regression neural networks in time-varying environment," *IEEE Trans. Neural Networks*, vol. 15, no. 3, pp. 576–596, May 2004.
- [6] J. Liu and D.-S. Chen, "Nonstationary fault detection and diagnosis for multimode processes," *AIChE Journal*, vol. 56, no. 1, pp. 207–219, Jan. 2010.
- [7] H. Ning, G. Qing, T. Tian, and X. Jing, "Online identification of nonlinear stochastic spatiotemporal system with multiplicative noise by robust optimal control-based kernel learning method," *IEEE Trans. Neural Networks and Learning Systems*, vol. 30, no. 2, pp. 389–404, Feb. 2019.
- [8] J. Shan, H. Zhang, W. Liu, and Q. Liu, "Online active learning ensemble framework for drifted data streams," *IEEE Trans. Neural Networks and Learning Systems*, vol. 30, no. 2, pp. 486–498, Feb. 2019.
- [9] J. L. Lobo, *et al.*, "Evolving spiking neural networks for online learning over drifting data streams," *Neural Networks*, vol. 108, pp. 1–19, Dec. 2018.
- [10] S. Chen and S. Billings, "Recursive prediction error parameter estimator for non-linear models," *Int. J. Control*, vol. 49, no. 2, pp. 569–594, 1989.
- [11] S. Chen, "Nonlinear time series modelling and prediction using Gaussian RBF networks with enhanced clustering and RLS learning," *Electronics Letters*, vol. 31, no. 2, pp. 117–118, Jan. 1995.
- [12] F. Ding, P. X. Liu, and G. Liu, "Multiinnovation least-squares identification for system modeling," *IEEE Trans. Systems, Man, Cybernetics: Part B*, vol. 40, no. 3, pp. 767–778, Jun. 2010.
- [13] G. B. Huang, Q. Y. Zhu, and C. K. Siew, "Extreme learning machine: Theory and applications," *Neurocomputing*, vol. 70, nos. 1–3, pp. 489–501, Dec. 2006.
- [14] N. Liang, G. B. Huang, P. Saratchandran, and N. Sundararajan, "A fast and accurate online sequential learning algorithm for feedforward networks," *IEEE Trans. Neural Networks*, vol. 17, no. 6, pp. 1411–1423, Nov. 2006.
- [15] G. B. Huang and L. Chen, "Enhanced random search based incremental extreme learning machine," *Neurocomputing*, vol. 71, nos. 16–18, pp. 3460–3468, Oct. 2008.
- [16] B. Krawczyk, *et al.*, "Ensemble learning for data stream analysis: A survey," *Information Fusion*, vol. 37, pp. 132–156, Sep. 2017.
- [17] L. L. Minku, A. P. White, and X. Yao, "The impact of diversity on online ensemble learning in the presence of concept drift," *IEEE Trans. Knowledge and Data Engineering*, vol. 22, no. 5, pp. 730–742, May 2010.
- [18] S. Wang, L. L. Minku, and X. Yao, "A systematic study of online class imbalance learning with concept drift," *IEEE Trans. Neural Networks and Learning Systems*, vol. 29, no. 10, pp. 4802–4821, Oct. 2018.
- [19] Y. Sun, *et al.*, "Online ensemble learning of data streams with gradually evolved classes," *IEEE Trans. Knowledge and Data Engineering*, vol. 28, no. 6, pp. 1532–1545, June 2016.
- [20] S. Wang, L. L. Minku, and X. Yao, "Resampling-based ensemble methods for online class imbalance learning," *IEEE Trans. Knowledge and Data Engineering*, vol. 27, no. 5, pp. 1356–1368, May 2015.
- [21] Y. Sun, *et al.*, "Concept drift adaptation by exploiting historical knowledge," *IEEE Trans. Neural Networks and Learning Systems*, vol. 29, no. 10, pp. 4822–4832, Oct. 2018.
- [22] D. Brzezinski and J. Stefanowski, "Reacting to different types of concept drift: The accuracy updated ensemble algorithm," *IEEE Trans. Neural Networks and Learning Systems*, vol. 25, no. 1, pp. 81–94, Jan. 2014.
- [23] G. Song, *et al.*, "Dynamic clustering forest: An ensemble framework to efficiently classify textual data stream with concept drift," *Information Sciences*, vol. 357, no. 1, pp. 125–143, 2016.
- [24] Y. Lan, Y. C. Soh, and G. B. Huang, "Ensemble of online sequential extreme learning machine," *Neurocomputing*, vol. 72, nos. 13–15, pp. 3391–3395, Aug. 2009.
- [25] T. Liu, S. Chen, S. Liang, and C. J. Harris, "Selective ensemble based multiple local model learning for nonlinear and nonstationary systems," *Neurocomputing*, vol. 378, pp. 98–111, Oct. 2020.
- [26] T. Liu, S. Chen, S. Liang, and C. J. Harris, "Growing and pruning selective ensemble regression for nonlinear and nonstationary systems," *IEEE Access*, vol. 8, pp. 73278–73292, Apr. 2020.
- [27] H. Chen, Y. Gong, and X. Hong, "A new adaptive multiple modelling approach for non-linear and non-stationary systems," *Int. J. Systems Science*, vol. 47, no. 9, pp. 2100–2110, Aug. 2016.
- [28] H. Chen, Y. Gong, X. Hong, and S. Chen, "A fast adaptive tunable RBF network for nonstationary systems," *IEEE Trans. Cybernetics*, vol. 46, no. 12, pp. 2683–2692, Dec. 2016.
- [29] H. Chen, Y. Gong, X. Hong, "Online modeling with tunable RBF network," *IEEE Trans. Cybernetics*, vol. 43, no. 3, pp. 935–947, May. 2013.
- [30] T. Xiong, Y. K. Bao, Z. Y. Hu, "Multiple-output support vector regression with a firefly algorithm for interval-valued stock price index forecasting," *Knowledge-Based Systems*, vol. 55, pp. 87–100, Oct. 2013.
- [31] D. J. Du, *et al.*, "A multi-output two-stage locally regularized model construction method using the extreme learning machine," *Neurocomputing*, vol. 128, pp. 104–112, Mar. 2014.
- [32] S. Chen, P. M. Grant, and C. F. N. Cowan, "Orthogonal least squares algorithm for training multi-output radial basis function networks," *IEE Proc. Part F*, vol. 139, no. 6, pp. 378–384, Dec. 1992.
- [33] S. Chen, "Multi-output regression using a locally regularised orthogonal least-squares algorithm," *IEE Proc. - Vision, Image and Signal Processing*, vol. 149, no. 4, pp. 185–195, Aug. 2002.
- [34] S. Chen, X. Hong, and C. J. Harris, "Sparse multioutput radial basis function network construction using combined locally regularised orthogonal least square and D-optimality experimental design," *IEE Proc. - Control Theory and Applications*, vol. 150, no. 2, pp. 139–146, Mar. 2003.
- [35] S. A. Billings and S. Chen, "The determination of multivariable non-linear models for dynamic systems," in *Control and Dynamic Systems, Volume 7 of Neural Network Systems Techniques and Applications* (ed. C. L. Leondes), San Diego: Academic Press, 1998, pp. 231–278.
- [36] W. Shao, S. Chen, and C. J. Harris, "Adaptive soft sensor development for multi-output industrial processes based on selective ensemble learning," *IEEE Access*, vol. 6, pp. 55628–55642, Oct. 2018.
- [37] C. Cheng and M. S. Chiu, "A new data-based methodology for nonlinear process modeling," *Chemical Engineering Science*, vol. 59, no. 13, pp. 2801–2810, Jul. 2004.
- [38] M. Nikraves, A. E. Farell, and T. G. Stanford, "Control of nonisothermal CSTR with time varying parameters via dynamic neural network control (DNNC)," *Chemical Engineering Journal*, vol. 76, no. 1, pp. 1–16, Jan. 2000.
- [39] C. A. Vriezanga, S. Sánchez-Pedreño, and J. Grasman, "Thermal runaway in microwave heating: A mathematical analysis," *Applied Mathematical Modelling*, vol. 26, no. 11, pp. 1029–1038, Nov. 2002.
- [40] B. P. Korin, "On the distribution of a statistic used for testing a covariance matrix," *Biometrika*, vol. 55, no. 1, pp. 171–178, Mar. 1968.
- [41] S. Jung, B. Moon, and D. Han, "Unsupervised learning for crowdsourced indoor localization in wireless networks," *IEEE Trans. Mobile Computing*, vol. 15, no. 11, pp. 2892–2906, Nov. 2016.
- [42] Y. You and M. Nikolaou, "Dynamic process modeling with recurrent neural networks," *AIChE J.*, vol. 39, pp. 1654–1667, 1993.
- [43] J. Zhong, S. Liang, Y. Yuan, and Q. Xiong, "Coupled electromagnetic and heat transfer ODE model for microwave heating with temperature-dependent permittivity," *IEEE Trans. Microwave Theory & Techniques*, vol. 64, no. 8, pp. 2467–2477, Aug. 2016.
- [44] J. Zhong, S. Liang, and Q. Xiong, "Improved receding horizon  $H_\infty$  temperature spectrum tracking control for Debye media in microwave heating process," *J. Process Control*, vol. 71, pp. 14–24, Nov. 2018.
- [45] T. Liu, S. Liang, Q. Xiong, and K. Wang, "Two-stage method for diagonal recurrent neural network identification of a high-power continuous microwave heating system," *Neural Processing Letter*, pp. 1–22, Feb. 2019.
- [46] T. Liu, S. Liang, Q. Xiong, and K. Wang, "Data-based online optimal temperature tracking control in continuous microwave heating system by adaptive dynamic programming," *Neural Processing Letter*, vol. 51, no. 1, pp. 167–191, Feb. 2020.



- [47] T. Liu, S. Liang, Q. Xiong, and K. Wang, "Integrated CS optimization and OLS for recurrent neural network in modeling microwave thermal process," *Neural Computing & Applications*, vol. 32, no. 16, pp. 12267–12280, Aug. 2020.
- [48] T. Liu, S. Liang, Q. Xiong, and K. Wang, "Adaptive critic based optimal neurocontrol of a distributed microwave heating system using diagonal recurrent network," *IEEE Access*, vol. 6, pp. 68839–68840, Dec. 2018.
- [49] T. Liu, S. Liang, and J. L. Hu, "Expert control system based hierarchical control strategy for tunnel microwave rice drying," *Proc. ECC 2019* (Naples, Italy), Jun. 25–28, 2019, pp. 3619–3624.



**Tong Liu** received the B.Sc. degree in automation from the College of Automation, Chongqing University, Chongqing, China, in 2016, where he is currently pursuing the Ph.D. degree in control theory and control engineering. From September 2018 to September 2019, he was a visiting Ph.D. student with the School of Electronics and Computer Science, University of Southampton, Southampton, UK. His current research interest include online learning, system identification, neural networks, machine learning and intelligent control system design.



**Sheng Chen** (M'90-SM'97-F'08) received his BEng degree from the East China Petroleum Institute, Dongying, China, in 1982, and his PhD degree from the City University, London, in 1986, both in control engineering. In 2005, he was awarded the higher doctoral degree, Doctor of Sciences (DSc), from the University of Southampton, Southampton, UK.

From 1986 to 1999, He held research and academic appointments at the Universities of Sheffield, Edinburgh and Portsmouth, all in UK. Since 1999, he has been with the School of Electronics and

Computer Science, the University of Southampton, UK, where he holds the post of Professor in Intelligent Systems and Signal Processing. Dr Chen's research interests include neural network and machine learning, wireless communications, and adaptive signal processing. He has published over 700 research papers. Professor Chen has 15,200+ Web of Science citations with h-index 54, and 30,700+ Google Scholar citations with h-index 75.

Dr. Chen is a Fellow of the United Kingdom Royal Academy of Engineering, a Fellow of IET, a Distinguished Adjunct Professor at King Abdulaziz University, Jeddah, Saudi Arabia, and an original ISI highly cited researcher in engineering (March 2004).



**ShanLiang** received his M.Sc. degree in control science and engineering from the College of Automation, Chongqing University, Chongqing, China, in 1995, and the Ph.D. degree from the Department of Mechanical Systems Engineering, Kumamoto University, Kumamoto, Japan, in 2004. His current research interests include nonlinear system modeling and adaptive control, with the applications to microwave heating processes and intelligent transportation systems. He is a Member of IEEE.



**Shaojun Gan** received the BEng degree in automation from Huainan normal University in 2010, and the PhD degree in control engineering from the School of Automation, Chongqing University, China, in 2016. From 2014 to 2016, He was a visiting PhD student at the School of Electronics, Electrical Engineering and Computer Science, Queen's University Belfast, UK. From 2017 to 2019, He was a research fellow at the School of Electronic and Electrical Engineering, University of Leeds, UK. Currently, he is an assistant professor with the

College of Metropolitan Transportation, Beijing University of Technology, China. His main research interests include system modelling and machine learning methods, with the applications to manufacturing energy systems and intelligent transportation systems.



**Chris J. Harris** received his BSc and MA degrees from the University of Leicester and the University of Oxford in UK, respectively, and his PhD degree from the University of Southampton, UK, in 1972. He was awarded the higher doctoral degree, the Doctor of Sciences (DSc), by the University of Southampton in 2001. He is Emeritus Research Professor at the University of Southampton, having previously held senior academic appointments at Imperial College, Oxford and Manchester Universities, as well as Deputy Chief Scientist for the UK

Government.

Professor Harris was awarded the IEE senior Achievement Medal for Data Fusion research and the IEE Faraday Medal for distinguished international research in Machine Learning. He was elected to the UK Royal Academy of Engineering in 1996. He is the co-author of over 500 scientific research papers during a 50 year research career.