

2017-10-30

Gaussian Process Regression for Virtual Metrology-enabled Run-to-Run Control in Semiconductor Manufacturing

Wan, Jian

<http://hdl.handle.net/10026.1/10123>

10.1109/TSM.2017.2768241

IEEE Transactions on Semiconductor Manufacturing

Institute of Electrical and Electronics Engineers (IEEE)

All content in PEARL is protected by copyright law. Author manuscripts are made available in accordance with publisher policies. Please cite only the published version using the details provided on the item record or document. In the absence of an open licence (e.g. Creative Commons), permissions for further reuse of content should be sought from the publisher or author.

Gaussian Process Regression for Virtual Metrology-enabled Run-to-Run Control in Semiconductor Manufacturing

Jian Wan, *Member, IEEE*, and Seán McLoone, *Senior Member, IEEE*

Abstract—Incorporating virtual metrology (VM) into run-to-run (R2R) control enables the benefits of R2R control to be maintained while avoiding the negative cost and cycle time impacts of actual metrology. Due to the potential for prediction errors from VM models, the prediction as well as the corresponding confidence information on the predictions should be properly considered in VM-enabled R2R control schemes in order to guarantee control performance. This paper proposes the use of Gaussian process regression (GPR) models in VM-enabled R2R control due to their ability to provide this information in an integrated fashion. The mean value of the GPR prediction is treated as the VM value and the variance of the GPR prediction is used as a measure of confidence to adjust the coefficient of an exponentially-weighted-moving-average (EWMA) R2R controller. The effectiveness of the proposed GPR VM-enabled R2R control approach is demonstrated using a chemical mechanical polishing process case study. Results show that better control performance is achieved with the proposed methodology than with implementations that do not take prediction reliability into account.

Index Terms—Virtual Metrology (VM), Run-to-Run (R2R) Control, Gaussian Process Regression (GPR), Exponentially-Weighted-Moving-Average (EWMA).

I. INTRODUCTION

Semiconductor manufacturing processes are highly repetitive and cyclic in nature and it is quite appealing to make use of the results from previous runs to improve the operation of subsequent ones. Those controllers that are capable of learning from past experience to improve future control performance are referred to as run-to-run (R2R) controllers in semiconductor manufacturing [1], [2], [3], [4], [5]. R2R control adjusts the process inputs or recipes run by run based on the information obtained before (pre-metrology), during and/or after (post-metrology) the process so as to compensate for the effects of process drifts, large shifts in incoming product and other disturbances. The adjustment is usually based on a process model and the process model is continuously updated using new measurements from ongoing runs. The exponentially-weighted-moving-average (EWMA) method is the most widely used R2R controller where a linear (affine) model is used to approximate the process and only the offset term in the model is updated [6]. The EWMA coefficient does not need to be fixed and it can be optimized for each

run [7]. There are several other R2R control schemes such as the double EWMA (dEWMA) method based on two coupled EWMA equations [8], [9], the least-squares estimation (LSE) method using a second-order model to approximate the process [10] and the set-valued R2R controller using sets such as an outer-bounding ellipsoid to approximate the likely model parameter set [11].

For all R2R control schemes, it is necessary to employ adequate metrology systems so as to update process models and modify recipes effectively based on metrology data. The deployment of metrology systems needs extra investment for metrology tools and metrology activities also impact negatively on production cycle time. In order to reduce the cost and the time needed for actual metrology, current metrology systems usually measure a few samples for the R2R controller. For example, the metrology systems in semiconductor manufacturing often measure one wafer from a lot of 25 wafers with the assumption that this adequately represents the quality of the whole lot. Such an approach restricts the application of R2R control to lot level. However, in practice significant within-lot wafer variation can occur, typically as a function of the wafer position within a lot (e.g. first wafer effects due to chamber seasoning in plasma etch processes). In addition, abnormal production equipment conditions can occur abruptly, which will not in general be immediately detected by one randomly sampled wafer [12]. While wafer level R2R control can address these issues with 100% pre-metrology measurement, the associated cost of metrology (capital and cycle time) generally prohibits its application.

To resolve the limitations of current sampling-based metrology systems, virtual metrology (VM) was proposed to predict the metrology values using sensor data from production equipment and actual metrology values of samples. The benefits of VM have been well reported in the literature [12], [13], [14], [15]. For instance, VM enables predictive maintenance and earlier detection of process drifts based on real-time forecast of metrology values [12]. VM also meets the need of R2R controllers by providing timely virtual metrology data both in terms of post-process feedback metrology and pre-process feedforward metrology [14], [16]. Thus the study and the implementation of VM-enabled R2R control have increasingly attracted attention from both academia and industry. For example, the interfaces between VM, other manufacturing execution components, and R2R modules were defined in [17] to develop a business model to measure the profitability of VM-related manufacturing practices, while in [18] VM and R2R solutions

J. Wan is with the School of Engineering, The University of Plymouth, Plymouth, PL4 8AA, UK. E-mail: jian.wan@plymouth.ac.uk.

S. McLoone is with the School of Electronics, Electrical Engineering and Computer Science, Queen's University Belfast, Belfast, BT9 5AH, Northern Ireland. E-mail: s.mcloone@qub.ac.uk.

are developed within a big data context. VM values were also employed to adjust equipment settings of a chemical vapor deposition process for maintenance compensation at a Taiwanese semiconductor manufacturing company [19].

There are some practical issues with applying VM-enabled R2R controllers in the semiconductor manufacturing industry. The essential one is that a virtual metrology value is not an actual metrology value and thus the error between the estimate and the actual value negatively impacts the performance of R2R controllers. The data quality of the VM values inside the VM-enabled R2R control scheme was considered in [13], [14] where the EWMA coefficient was adjusted according to the reliability of VM data. The concepts of reliance index (RI) and global similarity index (GSI) were proposed in [20] to gauge the reliability of VM data and they were further utilized within the feedback loop of R2R control in [21], [22]. By penalizing statistical measurements based on their informative distance from real metrology data, an information-theory and virtual metrology-based approach to R2R control was further proposed in [23]. The confidence information used in these VM-enabled R2R control schemes is obtained externally, i.e., the computation of the confidence information is separate from the computation of the VM estimates.

For some VM models, the confidence information of the predictions can also be obtained internally. For example, the VM model using relevance vector machine returns not only point estimates but also probabilistic intervals indicating the confidence of the predictions [24]. The relevance vector machine is actually a Gaussian process model with a specific covariance function [25]. It was observed in [26] that Gaussian process regression (GPR) performed better than multiple linear regression (MLR), least absolute shrinkage and selection operator (LASSO) and neural networks (NN) for the VM task of a benchmark semiconductor manufacturing process. Motivated by such experience this paper proposes using GPR to provide a probabilistic prediction of metrology for a confidence information enhanced R2R control scheme. Specifically, by using the mean estimate from the GPR model as the predicted VM value and the variance estimate as a measure of confidence to adjust the coefficient of an EWMA R2R controller, a novel self-contained, confidence information enhanced, VM-enabled post-metrology R2R control solution is obtained.

The rest of the paper is organized as follows. GPR is briefly introduced in Section II. The proposed VM-enabled R2R control scheme using GPR is detailed in Section III. Then, as an illustrative example, Section IV investigates the application of the control scheme to a chemical mechanical polishing process commonly used in semiconductor manufacturing. Finally, conclusions are presented in Section V.

II. GAUSSIAN PROCESS REGRESSION

Gaussian processes extend multivariate Gaussian distributions to infinite dimensionality, i.e., the joint distribution over any finite set of fixed test points is a multivariate Gaussian. According to the tutorial in [27], GPR is briefly introduced as follows.

Given a relationship of the form $y = f(\mathbf{x})$ between the input \mathbf{x} and the single output y , and set of n test points, then

$$[f(\mathbf{x}_1) f(\mathbf{x}_2) \cdots f(\mathbf{x}_n)]^T \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad (1)$$

where $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ signifies a multivariate normal distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. $\Sigma_{ij} \in \boldsymbol{\Sigma}$ defines the covariance between $f(\mathbf{x}_i)$ and $f(\mathbf{x}_j)$ and can often be defined as a function of the input \mathbf{x}_i and \mathbf{x}_j , i.e., $\Sigma_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$. The mean value $\boldsymbol{\mu}$ of the normal distribution \mathcal{N} is usually assumed to be zero, i.e., $\boldsymbol{\mu} = \mathbf{0}$. The Gaussian process is then fully specified by the mean function $\boldsymbol{\mu}$ and the covariance function $\boldsymbol{\Sigma}$. It can be seen that GPR is a probabilistic non-parametric modeling technique as it does not impose a specific model structure on the function itself.

The covariance function $k(\mathbf{x}_i, \mathbf{x}_j)$ can be a function of any form as long as it generates a positive definite covariance matrix $\boldsymbol{\Sigma}$. An often used covariance function is the squared exponential covariance function [28], which has the following form:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp\left(-\frac{|\mathbf{x}_i - \mathbf{x}_j|^2}{2l^2}\right), \quad (2)$$

where σ_f and l are the hyperparameters (denoted by vector $\boldsymbol{\theta}$) for the covariance function which can be optimized to best suit the training data. If $\mathbf{x}_i \approx \mathbf{x}_j$, then $k(\mathbf{x}_i, \mathbf{x}_j)$ approaches its maximum value, indicating that $f(\mathbf{x}_i)$ is almost perfectly correlated with $f(\mathbf{x}_j)$; otherwise, distant observations have negligible effect on each other as $k(\mathbf{x}_i, \mathbf{x}_j)$ approaches zero.

Considering measurement errors and other noise sources, the observation y can be further expressed as a combination of the underlying function $f(\mathbf{x})$ and a Gaussian noise model with a mean value of 0 and a variance of σ_n^2 :

$$y = f(\mathbf{x}) + \mathcal{N}(0, \sigma_n^2). \quad (3)$$

Therefore

$$\mathbf{y} = [y_1 \ y_2 \ \cdots \ y_n]^T \sim \mathcal{N}(\mathbf{0}, \mathbf{K}), \quad (4)$$

where $\mathbf{K} = \boldsymbol{\Sigma} + \sigma_n^2 \mathbf{I}$. Here the output variable is assumed to be mean-centered.

According to the assumption made for Gaussian processes, the joint distribution over the observed targets \mathbf{y} and a test target y_* is also Gaussian:

$$\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \mathbf{K} & \mathbf{K}_*^T \\ \mathbf{K}_* & K_{**} \end{bmatrix}\right), \quad (5)$$

where $\mathbf{K}_* = [k(\mathbf{x}_*, \mathbf{x}_1) \ k(\mathbf{x}_*, \mathbf{x}_2) \ \cdots \ k(\mathbf{x}_*, \mathbf{x}_n)]$ and $K_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$. Then the conditional probability of $p(y_*|\mathbf{y})$ also follows a Gaussian distribution:

$$y_*|\mathbf{y} \sim \mathcal{N}(K_*\mathbf{K}^{-1}\mathbf{y}, K_{**} - K_*\mathbf{K}^{-1}K_*^T). \quad (6)$$

Hence, the best estimate for y_* is the mean of this distribution and the uncertainty of the estimate is captured in its variance, that is:

$$\bar{y}_* = K_*\mathbf{K}^{-1}\mathbf{y}, \quad (7)$$

and

$$\text{var}(y_*) = K_{**} - K_*\mathbf{K}^{-1}K_*^T. \quad (8)$$

The operation of GPR for a single-input single-output system is shown in Figure 1, where the regression model is

trained using just four measurement points at $x = -2, -1, 1, 2$. The values of y at other points $-2 \leq x \leq 2$ are estimated by the trained GPR model from these four measurement points. It can be seen from Figure 1 that GPR returns both mean values and confidence intervals for all test points, where the 95% confidence interval is taken as $\bar{y}_* \pm 1.96 \cdot \sqrt{\text{var}(y_*)}$. Furthermore, the confidence interval for the prediction grows with the distance from the four measurement points used for the training, which implies less confidence in the prediction. This is physically intuitive as GPR can provide more accurate predictions for those test points that are similar to the data used to train the model.

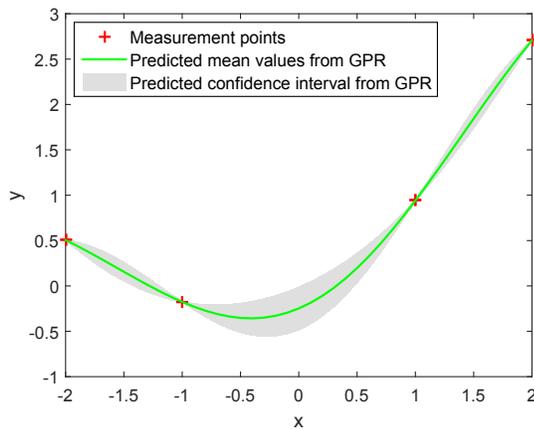


Fig. 1. Mean value and 95% confidence domain predicted by GPR

The performance of GPR depends on how well the covariance function is selected [27]. The hyperparameters θ for the selected covariance function can be optimized using the training data, i.e., by maximizing the conditional probability $p(\theta|\mathbf{x}, \theta)$, which corresponds to maximizing the log marginal likelihood $\log p(\mathbf{y}|\mathbf{x}, \theta)$:

$$\log p(\mathbf{y}|\mathbf{x}, \theta) = -\frac{1}{2} \mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}| - \frac{n}{2} \log 2\pi. \quad (9)$$

Using multivariate optimization algorithms, the optimal choice for θ can be obtained straightforwardly [25].

III. VM-ENABLED R2R CONTROL USING GPR

Using a similar R2R control diagram to that presented in [22], the proposed VM-enabled R2R control scheme using GPR is described in Figure 2, where k is the process run index; \mathbf{u}_{k+1} is the process input or recipe setting at the start of run $k+1$, which is derived from the EWMA controller; y_t is the target output; the pair $(\bar{y}_*, \text{var}(y_*))$ is the output of the GPR VM model at the end of run k ; and y_z is the sampled metrology data. Both \bar{y}_* and $\text{var}(y_*)$ from the GPR VM model are used in the EWMA controller, as discussed in the following paragraphs.

Consider a process with linear input and output relationship defined as:

$$y_k = \mathbf{b}\mathbf{u}_k + \eta_k, \quad (10)$$

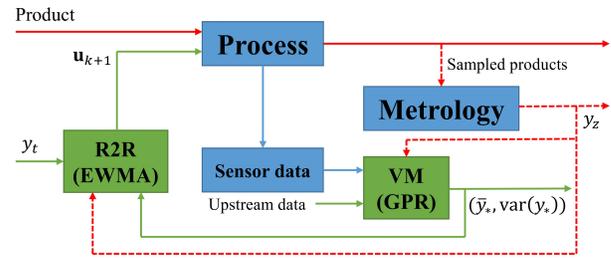


Fig. 2. VM-enabled R2R control using GPR

where y_k is the plant output, \mathbf{u}_k is the control input, \mathbf{b} is the process gain, and η_k is a disturbance. The disturbance η_k is updated recursively for the next run $k+1$ by an EWMA filter:

$$\eta_{k+1} = \alpha(y_k - \mathbf{b}\mathbf{u}_k) + (1 - \alpha)\eta_k, \quad (11)$$

where α is the exponential weighting factor or tuning parameter for the filter. When \mathbf{b} is invertible, the control input at the next run \mathbf{u}_{k+1} can be expressed as:

$$\mathbf{u}_{k+1} = (y_t - \eta_{k+1})/\mathbf{b}. \quad (12)$$

In a similar manner to the approach proposed in [13], [22], when y_k is obtained from an actual metrology tool measurement, then $y_k = y_z$ and a conventional EWMA controller can be employed,

$$\eta_{k+1} = \alpha_0(y_z - \mathbf{b}\mathbf{u}_k) + (1 - \alpha_0)\eta_k, \quad (13)$$

with the EWMA coefficient α_0 determined using standard design methods such as those discussed in [7] [29]. When y_k is obtained from the prediction provided by the GPR VM model, then $y_k = \bar{y}_*$ and the EWMA coefficient α_0 is scaled by a confidence factor, referred to as the Gaussian Reliance Index (GRI), to give:

$$\eta_{k+1} = \alpha_1(\bar{y}_* - \mathbf{b}\mathbf{u}_k) + (1 - \alpha_1)\eta_k, \quad (14)$$

where

$$\alpha_1 = \text{GRI} \cdot \alpha_0. \quad (15)$$

GRI is derived from the coefficient of variation of the prediction distribution ($c_v(y_*)$) such that it is monotonically maps the unbounded confidence information provided by GPR to a confidence factor spanning the interval 0 and 1, with 1 indicating full confidence in the accuracy of the prediction and 0 no confidence. The coefficient of variation (CV) is defined as:

$$c_v(y_*) = \frac{\text{std}(y_*)}{\bar{y}_*}, \quad (16)$$

where $\text{std}(y_*) = \sqrt{\text{var}(y_*)}$.

Various possibilities exist for the GRI mapping function, Here two options are considered, a truncated linear decay function of the form:

$$\text{GRI}(c_v(y_*)) = [1 - \beta \frac{c_v(y_*)}{c_v^{\max}}]_+, \quad (17)$$

which yields a GRI, and hence α_1 value, that decays linearly to 0 with increasing prediction uncertainty and which is

exactly zero when $c_v(y_*) = c_v^{\max}/\beta$, and an exponential decay function:

$$\text{GRI}(c_v(y_*)) = \exp\left(-\beta \frac{|c_v(y_*)|}{c_v^{\max}}\right), \quad (18)$$

which provides a more gradual decay towards zero with increasing $c_v(y_*)$. Constant c_v^{\max} is a normalization factor employed to scale the value of $c_v(y_*)$ with respect to the expected normal range of variation, as defined by the distribution of $c_v(y_*)$ values for the reference training dataset used to generate the GPR VM model. Various options exist for defining c_v^{\max} with respect to the reference distribution, e.g., the maximum value or the 95th percentile. Here the maximum value over the training dataset is selected as the value of c_v^{\max} .

Larger values of the CV map to smaller values of the GRI and consequently smaller values of the EWMA coefficient α_1 . If $c_v(y_*) = 0$ or $y_k = y_z$, then $\text{GRI} = 1$ and $\alpha_1 = \alpha_0$, which is the case when actual metrology is used. The sensitivity of the GRI to changes in the CV can then be controlled by the parameter $\beta > 0$. As will be demonstrated in Section IV, in practice the exponential mapping provides consistent performance over a much wider range of β than the truncated linear mapping. As such the exponential function is the preferred mapping with $\beta = 1$.

Note that, while the GRI modulated α value (i.e. α_1) can be expected to yield better control performance on average than using a fixed value of α_0 or control without VM feedback (i.e. with $\alpha_1 = 0$) it is not guaranteed to be better for every run. This is because the choice of α_1 is based on a statistical quantity, GRI, and the heuristically selected mapping from α_0 to α_1 is not guaranteed to be optimal. Determination of the optimal α_0 to α_1 mapping is a topic for future research.

It can be seen that the proposed GPR VM-enabled R2R control scheme makes full use of the predictions made by GPR and the confidence information is internally available from GPR rather than from external computations of RI and GSI as in [13], [22]. Similarly to the approach in [22], it is also possible to set a baseline threshold level GRI_T for the GRI, below which α_1 is set to zero, i.e., if $\text{GRI} \leq \text{GRI}_T$ then $\alpha_1 = 0$. This can be defined with reference to the distribution of the GRI values for the GPR predictions of the training data. Here GRI_T is set as the 5th percentile of the distribution. The motivation for this modification is a desire to have no contribution from the current prediction if it is deemed very unreliable [22].

It is relatively straightforward to update GPR models with new data, hence the GPR VM model can be continuously improved as new measurements become available and accordingly $c_v(y_*) \rightarrow 0$, $\text{GRI} \rightarrow 1$, and $\alpha_1 \rightarrow \alpha_0$. The fundamental assumption of VM models, and hence also the GPR prediction model, is that variations in the metrology value being predicted must be reflected in the variations in the inputs being fed to the model in order for the model to be able to predict the metrology value correctly. If the variation in the inputs is outside the previous experience of the GP model then this will be reflected in higher variance and hence lower confidence levels in predictions. If however, a drift occurs that is not reflected in the model inputs the model will be blind to it, and hence not react, with negative consequences.

Only actual metrology measurements can detect such drifts, hence periodic metrology is advisable even when using VM models. In practice, it is not possible to capture the impact of all process drifts or maintenance related shifts in VM models and, therefore, model maintenance strategies such as those discussed in [30] are vital for successful long term deployment.

IV. ILLUSTRATIVE EXAMPLE

In order to demonstrate the proposed GPR VM-enabled R2R control scheme and compare it with existing methods in the literature, the illustrative example studied in [22], a chemical mechanical polishing (CMP) tool with a scheduled maintenance cycle of 600 wafers, is adopted. The corresponding process can be described by the following linear relationship between the input and the output:

$$y_k = \text{Pre}Y_k - r_k u_k, \quad (19)$$

where y_k is the post CMP thickness of run k and $\text{Pre}Y_k$ is the initial CMP thickness of run k ; r_k is the actual removal rate of run k and u_k is the polishing time. The target CMP thickness is denoted by y_t and $y_t = 2800$ Angstrom in the following simulation.

The actual removal rate r_k of run k is hard to measure in practice and here its value is simulated by the following formula [22]:

$$r_k = A_k \times \left(\frac{\text{Stress}_1 + \text{Stress}_2}{1000}\right) \times \left(\frac{\text{Rotspd}_1 + \text{Rotspd}_2}{100}\right) \times \left(\frac{\text{Sfuspd}_1 + \text{Sfuspd}_2}{100}\right) + (PM_1 + PM_2) + \text{Error}, \quad (20)$$

where the meanings of Error , PM_1 , PM_2 , Stress_1 , Stress_2 , Rotspd_1 , Rotspd_2 , Sfuspd_1 , Sfuspd_2 and $\text{Pre}Y_k$ are listed in Table I as in [22]; $A_k = (4 \times 10^{-6}) \cdot (PU - 1)^3 + (3.4 \times 10^{-3}) \cdot (PU - 1)^2 + (6.9 \times 10^{-3}) \cdot (PU - 1) + 1.202 \times 10^3$ is the nominal removal rate and it is assumed to be known for the controller; PU is the parts usage count between periodic maintenances (PMs).

TABLE I
DEFINITION OF SIMULATION PARAMETERS AND SETTING VALUES [22]

Item	Definition
Error	Random error
PM_1	Error due to tool-parts' variation from periodic maintenance
PM_2	Random disturbance of tool-parts' variation
Stress_1	Tool stress error due to re-assembly during PM
Stress_2	Random disturbance of tool stress
Rotspd_1	Tool rotation-speed error due to reassembly during PM
Rotspd_2	Random disturbance of tool rotation speed
Sfuspd_1	Slurry fluid-speed error due to reassembly during PM
Sfuspd_2	Random disturbance of slurry fluid speed
$\text{Pre}Y_k$	Pre-process (etching depth) value that affects the k th run

According to (20), the actual removal rate r_k of run k can be treated as a function of six process variables:

$$r_k = f(\text{Stress}, \text{Rotspd}, \text{Sfuspd}, PU, PU^2, PU^3), \quad (21)$$

where $\text{Stress} = (\text{Stress}_1 + \text{Stress}_2)$, $\text{Rotspd} = (\text{Rotspd}_1 + \text{Rotspd}_2)$ and $\text{Sfuspd} = (\text{Sfuspd}_1 + \text{Sfuspd}_2)$. Hence, using these six variables as inputs VM models can be identified from simulated runs of the CMP process to estimate the actual remove rate.

The estimation for r_k is denoted as \bar{r}_k . Two VM models identified using MLR and GPR are used for the computation of the RI, where the MLR model is the reference model and the GPR model is the conjecture model. Adopting the GPR model as the common VM model for all VM-enabled R2R control schemes ensures that any performance differences between controllers can be attributed exclusively to if and how confidence information is used in the control schemes.

The GPR model returns both the estimated removal rate \bar{r}_k and the corresponding variance $\text{var}(r_k)$. Then $\bar{y}_k = \text{Pre}Y_k - \bar{r}_k u_k$ is the estimated CMP thickness of run k with the standard deviation $\text{std}(y_k) = u_k \cdot \text{std}(r_k)$. Based on the proposed VM-enabled R2R control scheme shown in Figure 2, the control input u_{k+1} for the run $k+1$ should be:

$$u_{k+1} = (\text{Pre}Y_{k+1} - y_t + \eta_{k+1}) / \bar{r}_{k+1}. \quad (22)$$

If $y_k = y_z$, i.e., y_k is obtained from the actual metrology, then

$$\eta_{k+1} = \alpha_0 (y_z - \text{Pre}Y_k + \bar{r}_k u_k) + (1 - \alpha_0) \eta_k. \quad (23)$$

If $y_k \approx \bar{y}_k$, i.e., y_k is estimated from the GPR VM model, then

$$\eta_{k+1} = \alpha_1 (\bar{y}_k - \text{Pre}Y_k + \bar{r}_k u_k) + (1 - \alpha_1) \eta_k, \quad (24)$$

where $\alpha_1 = \text{GRI} \cdot \alpha_0$ and GRI is computed according to equation (17) or (18).

TABLE II
SIMULATION PARAMETER VALUES AS USED IN [22]

Item	Normal Distribution	
	Mean	Variance
<i>Error</i>	0	300
<i>PM</i> ₁	0	100
<i>PM</i> ₂	0	6
<i>Stress</i> ₁	1000	2000
<i>Stress</i> ₂	0	20
<i>Rotspd</i> ₁	100	25
<i>Rotspd</i> ₂	0	1.2
<i>Sfusp</i> ₁	100	25
<i>Sfusp</i> ₂	0	1.2
<i>PreY</i> _k	3800	2500

Since the quality of the confidence information plays an essential role in the proposed GPR VM-enabled R2R controller, the first simulation investigates the quality of the confidence information provided by GPR via the GRI for the exponential mapping with $\beta = 1$ and c_v^{\max} equal to the maximum value of CV over the training data. For comparison purposes the RI confidence metric is also evaluated. Ideally, the confidence levels defined by GRI and RI should be proportional to the percentage prediction errors. To investigate this MLR and GPR VM models are identified from process data of sample runs $PU = 1 : 4 : 401$ for the setup where all process variables are normally distributed with parameter values as listed in Table II. Only a fraction of the data from a full run is used to train the models so as to limit their ability to generalize, increasing the likelihood of larger prediction errors on test runs. This enables the GRI and RI to be tested over a wider range of errors than would be facilitated by a more comprehensively trained model. The two VM models are then tested on a full cycle of runs $PU = 1 : 1 : 600$ to obtain \bar{r}_k , GRI and RI, respectively.

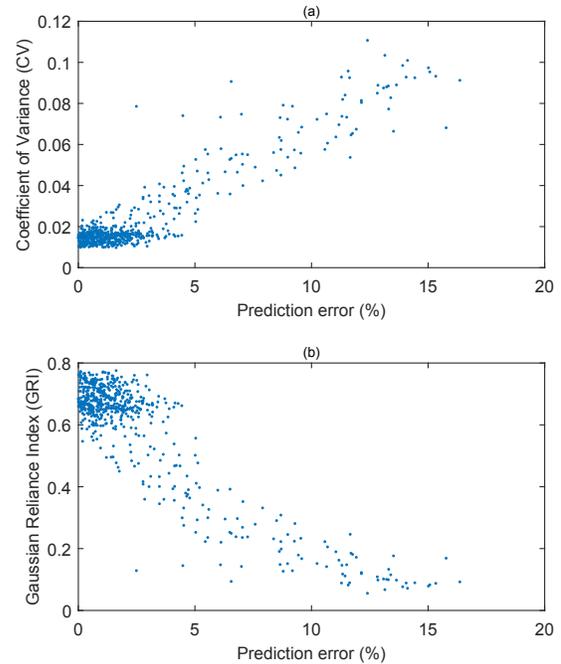


Fig. 3. A scatter plot of the CV and GRI values versus the GPR VM model percentage prediction errors

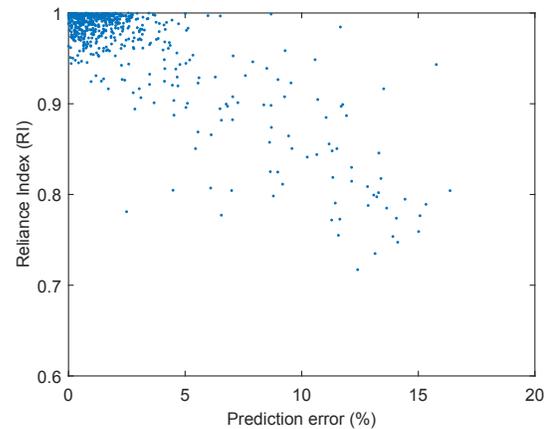


Fig. 4. A scatter plot of the RI values versus the GPR VM model percentage prediction errors

The coefficients of variation (CV) of the GPR model predictions, $c_v(r_k)$, and the corresponding GRI values are plotted as a function of the percentage prediction errors in removal rate in Figure 3(a) and 3(b), respectively. It can be seen that in general the CV values increase while the GRI values decrease as the prediction errors increase, and that the GRI is widely spread over the interval [0,1]. Hence, the GRI provides an effective measurement of confidence in the GPR VM predictions.

Similarly, the RI values (derived from the area of intersection between the statistical distribution of the GPR-based VM estimate from the conjecture model and the statistical distribution of the MLR-based reference prediction value as

described in [20]) are plotted in Figure 4 as a function of the percentage prediction errors obtained with the GPR conjecture model. The RI values obtained also reflect the confidence of the predictions made by GPR with bigger percentage prediction errors corresponding to smaller RI values, although the spread of values is less than observed with the GRI. It can also be seen from Figure 3(b) and Figure 4 that the degree of correlation between the computed confidence levels and the percentage prediction errors is stronger for the GRI than for the RI. Specifically, the magnitude of Pearson's correlation coefficient with respect to the percentage prediction errors is 0.898 for the GRI compared to 0.826 for RI. Hence, the GRI measure is more sensitive to prediction errors than RI in this case and, therefore, a stronger indicator of confidence in VM model predictions.

Using the identified VM models, four different R2R control schemes are compared:

- Case 1: R2R with in-situ metrology (actual metrology) and a constant weighting parameter $\alpha = \alpha_0$;
- Case 2: R2R with a GPR VM model, and with a constant $\alpha = \alpha_0$;
- Case 3: R2R with a GPR VM model and with the GPR derived GRI used to adjust the α parameter, such that: $\alpha = \alpha_1 = g(\text{GRI}, \text{GRI}_T) \cdot \alpha_0$ where $g(\text{GRI}, \text{GRI}_T) = 0$ if $\text{GRI} < \text{GRI}_T$ and $g(\text{GRI}, \text{GRI}_T) = \text{GRI}$ otherwise;
- Case 4: R2R with GPR used as the VM model and the RI used to adjust α as given in [22], that is: $\alpha = \alpha_1 = f(\text{RI}, \text{RI}_T, \text{GSI}_T) \cdot \alpha_0$ with $f(\text{RI}, \text{RI}_T, \text{GSI}_T) = 0$ if $\text{RI} < \text{RI}_T$ or $\text{GSI} > \text{GSI}_T$ and $f(\text{RI}, \text{RI}_T, \text{GSI}_T) = \text{RI}$ otherwise, where RI_T and GSI_T are the thresholds for the RI and global similarity index (GSI), respectively.

The process mean-absolute-percentage error (MAPE) with respect to the target value:

$$\text{MAPE} = \frac{\sum_{k=1}^N |(y_k - y_t)/y_t|}{N} \times 100\% \quad (25)$$

is used to evaluate the performance of these four R2R control schemes.

Following [22], the first comparison for these four R2R control schemes is performed on the CMP process under the assumption that all process variables are normally distributed with distribution parameters as listed in Table II. These nominal conditions provide the reference process data for training the GPR VM model. Controller performance is then evaluated for test data where an offset of 20 is added to *Rotspd* for runs $PU = 100 : 1 : 400$ of the maintenance cycle to simulate an undetected process shift during the control process. Figure 5 shows a plot of the evolution of the removal rate under these conditions together with a plot of the ideal model removal rate and the GPR VM model prediction.

Using various values of α_0 , and setting $\text{RI}_T = 0.7$ and $\text{GSI}_T = 9$ in accordance with [22], the MAPEs of these four approaches for a typical maintenance cycle (600 wafers) are shown in Figure 6. It can be seen that the R2R controller using actual metrology performs better than the other three VM-enabled R2R controllers for small values of α_0 and that the R2R controllers using confidence information also perform better than the R2R controller which does not take confidence

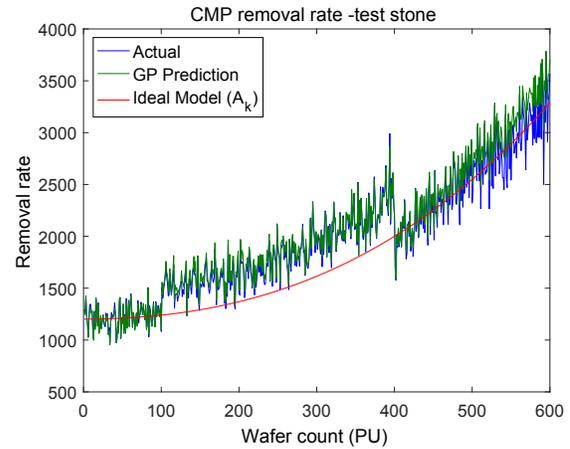


Fig. 5. Typical removal rate evolution over a maintenance cycle for the Gaussian distribution CMP model with a process drift/offset introduced during the processing of wafer 100-400. The plot also shows the VM GPR prediction and ideal model removal rates

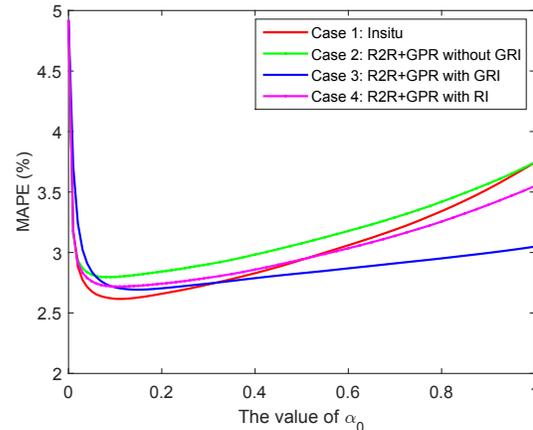


Fig. 6. MAPE performance of the four R2R control schemes for normally distributed process variables over a full maintenance cycle (600 wafers)

information into account. For larger values of α_0 , the proposed GRI enhanced controller of Case 3 performs much better than all other cases, demonstrating a degree of robustness in relation to the choice of α_0 . As a lower value of α_0 implies better R2R control performance for this specific process, the adjustment from α_0 to α_1 for Case 3 contributes to its superior performance to Case 1 for larger values of α_0 .

A plot of the value of the thresholded GRI (i.e. α_1/α_0) for Case 3 is given in Figure 7. As can be seen the GRI drops to zero above run 475, which is the point where the VM model begins to diverge from the actual removal rate, and has a higher frequency of thresholded values (i.e. $\text{GRI} < \text{GRI}_T$) in the period with the offset (c.f. Figure 5).

In order to further confirm these observations in Figure 6, a second comparison of these four R2R control schemes is performed on the CMP process under the new assumption that all process variables follow Weibull distributions. The Weibull

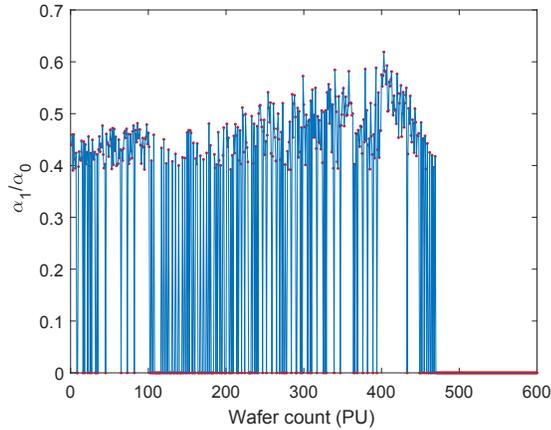


Fig. 7. Evolution of the thresholded GRI (i.e. α_1/α_0) for Gaussian distributed process variables over a full maintenance cycle (600 wafers)

TABLE III
WEIBULL-DISTRIBUTION PARAMETER VALUES AS USED IN [22]

Item	Weibull Distribution	
	α_ω	β
<i>Error</i>	108.5	7.2
<i>PM</i>	2800.1	850.1
<i>Stress</i>	0.951	150
<i>Rotspd</i>	0.986	100
<i>Sfusp</i>	0.917	250
<i>PreY_k</i>	3826.3	77.83

distribution is described by the equation [22]:

$$f(x) = \frac{\beta}{\alpha_\omega} \times \left(\frac{x}{\alpha_\omega}\right)^{\beta-1} \times e^{-\left(\frac{x}{\alpha_\omega}\right)^\beta}, \quad (26)$$

where α_ω and β are referred to as the scale parameter and shape parameter, respectively. The nominal parameters of the *Error*, *PM*, *Stress*, *Rotspd*, *Sfusp* and *PreY_k* Weibull distributions are listed in Table III, where $PM=PM_1+PM_2$ and biases of 100 and 2800 are deducted from every *Error* and *PM* value created by the Weibull distributions to generate potentially negative values [22].

The VM models are also identified using a fraction of the maintenance cycle data sampled at $PU = 1 : 4 : 401$ under nominal conditions. Since the Weibull CMP model has an offset with respect to the ideal model over the full maintenance cycle, an additional disturbance in the form of increased variation in removal rate is introduced for wafers 100-400 for the test runs. This is achieved by increasing the variance of *Rotspd* by 0.1 at $PU = 100 : 1 : 400$. Figure 8 shows a plot of the resulting removal rate, the ideal model, and the GPR VM model approximation for this scenario. The MAPEs of the four R2R controllers for a typical simulation run are shown in Figure 9. The results show a similar trend to Figure 6 with the proposed GPR-enabled R2R controller outperforming the other controllers for most values of α_0 . The optimum value occurs at $\alpha_0 = 0.3$. For completeness a plot of the evolution of the thresholded GRI is presented in Figure 10. This also follows a similar pattern to the results obtained with the Gaussian distributed process variables (Figure 7).

Two factors contribute to the performance of the GRI based R2R controller (Case 3) with respect to the real metrology based R2R controller (Case 1). The applied α with real metrology is always fixed at α_0 whereas the applied α with GRI will in general be less than α_0 . Hence to get the same effective α the optimum α_0 for GRI will always be greater than the value with real metrology. The superior performance of GRI, especially for higher values of α_0 can be partly attributed to this. The other factor is that larger variations in removal rate generally produce lower confidence predictions, and hence a lower value of applied α . Hence, in effect the GRI based controller is a nonlinear controller with a gain that reduces as the variance in the removal rate increases. This enhances the performance of the controller making it possible for it to outperform the in-situ metrology based fixed gain linear R2R controller.

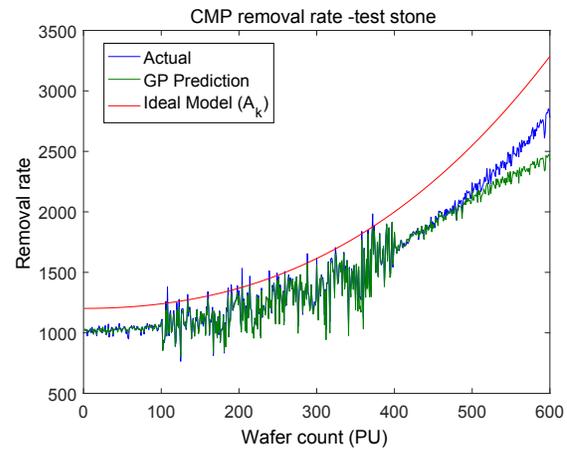


Fig. 8. Typical removal rate evolution over a maintenance cycle for the Weibull distribution CMP model with an increased process variance introduced during the processing of wafer 100-400. The plot also shows the VM GPR prediction and ideal model removal rates

In order to account for stochastic variation and provide a statistically robust comparison of the controllers, 100 simu-

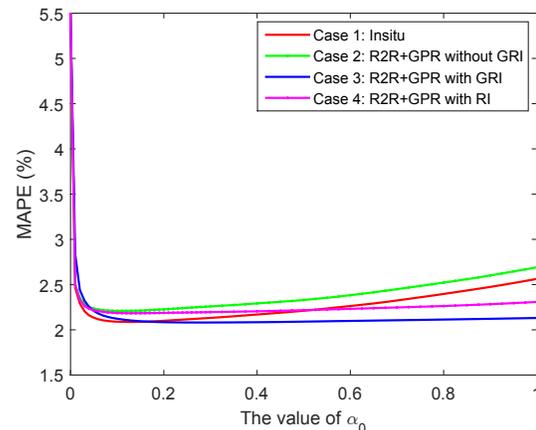


Fig. 9. MAPE performance of the four R2R control schemes for Weibull distributed process variables over a full maintenance cycle (600 wafers)

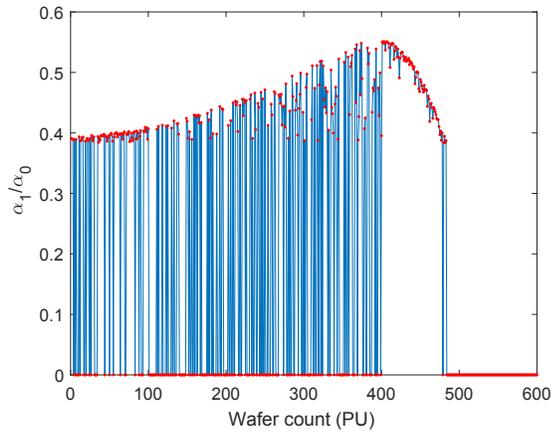


Fig. 10. Evolution of the thresholded GRI (i.e. α_1/α_0) for Weibull distributed process variables over a full maintenance cycle (600 wafers)

lation runs (i.e full maintenance cycles of 600 wafers) were performed for the normal distribution and Weibull distribution setups. The average MAPEs over the 100 repetitions for each setup are shown in Figure 11 and 12, respectively. The 95% confidence interval for the mean estimates for Case 3 and Case 4 are also plotted in the figures (dashed lines). The confidence intervals for Case 1 and Case 2, which are of similar scale to those of Case 3 and Case 4, respectively, have been omitted in the interest of clarity. It can be seen that the Monte Carlo simulation results confirm the observations in Figure 6 and Figure 9. Furthermore, the 95% confidence interval for the proposed approach of Case 3 are narrower than for Case 4, reflecting the fact that the performance of the proposed GRI enhanced R2R controller is more consistent over the 100 repetitions.

Plots of the evolution of the value of the thresholded GRI (Case 3) averaged over the 100 Monte Carlo simulations are given in Figures 13 and 14, respectively, for the Gaussian distribution and Weibull distribution CMP simulation models. The operation of the GRI is clearly evident in these plots, with both the period of process perturbation ($PU = 100 - 400$) and the period where the model is extrapolating outside its training range ($PU > 450$) yielding substantially lower average GRI values than during normal operation.

As a final experiment, to evaluate the performance of the truncated linear decay and exponential GRI mapping functions, and to assess the sensitivity of the resulting R2R controllers to β , a 100 run Monte Carlo study was conducted using each mapping function for β values in the range 0 to 5. The mean and standard deviation of the MAPE of the optimum controller (Case 3) for each value of β over the 100 simulation runs are plotted in Figures 15 and 16, for the Gaussian and Weibull CMP models, respectively. The results show that for values of $\beta < 1$ the linear and exponential GRI mapping functions yield similar performance. Beyond the value of 1 the linear mapping deteriorates rapidly while in contrast the exponential function provides consistent performance up to $\beta = 3$ with a gradual degradation in performance thereafter.

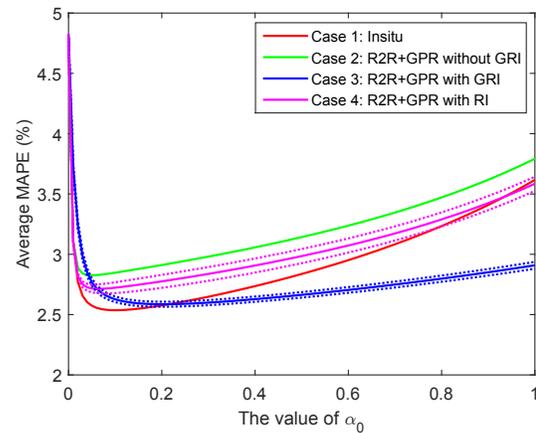


Fig. 11. Monte Carlo simulation results (100 runs) for R2R control of the CMP process assuming process variables follow normal-distributions

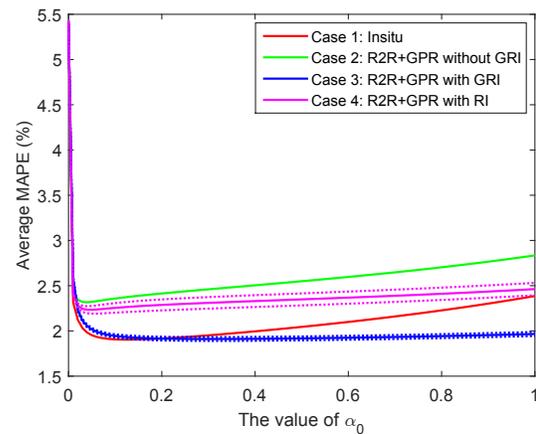


Fig. 12. Monte Carlo simulation results (100 runs) for R2R control of the CMP process assuming process variables follow Weibull-distributions

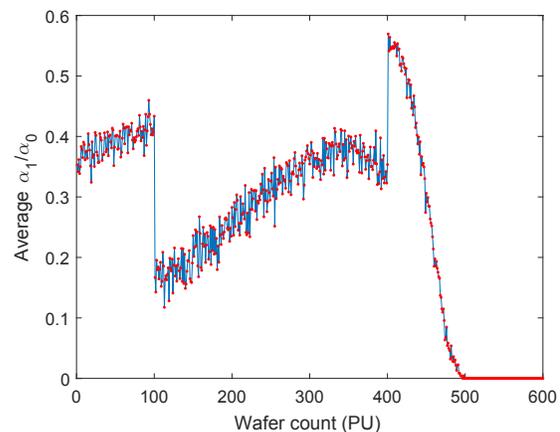


Fig. 13. Evolution of the thresholded GRI (i.e. α_1/α_0) for Gaussian distributed process variables over a full maintenance cycle (600 wafers) averaged over 100 Monte Carlo simulations

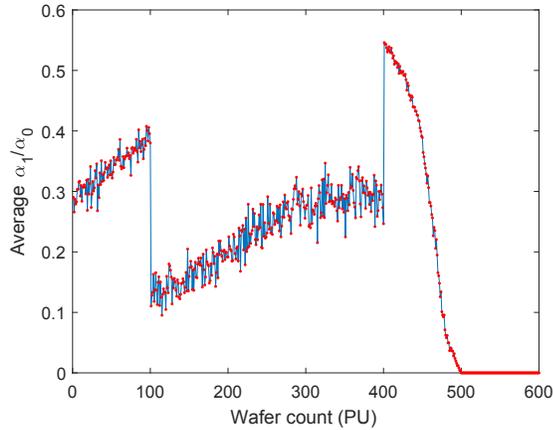


Fig. 14. Evolution of the thresholded GRI (i.e. α_1/α_0) for Weibull distributed process variables over a full maintenance cycle (600 wafers) averaged over 100 Monte Carlo simulations

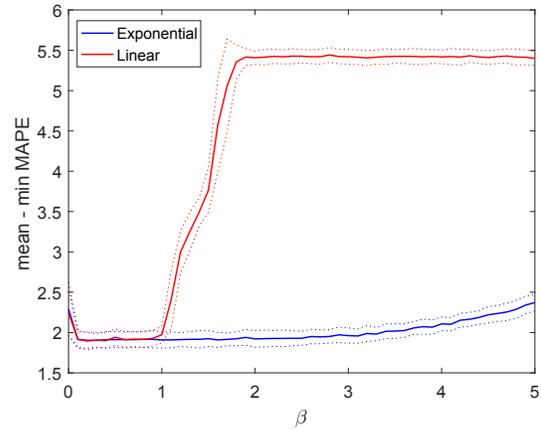


Fig. 16. Performance of the Case 3 R2R controller on the Weibull CMP process when employing truncated linear decay and exponential decay based GRI as a function of β . Results are computed over 100 Monte Carlo simulations (dashed lines indicate one standard deviation)

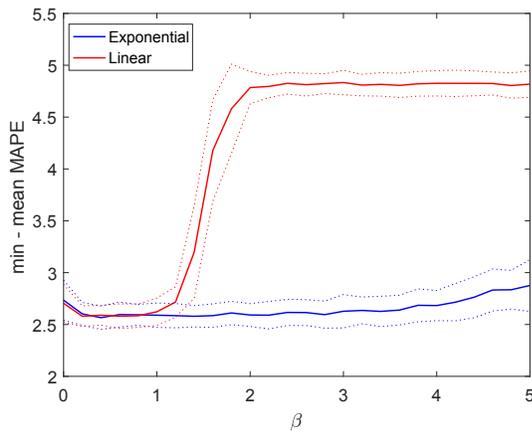


Fig. 15. Performance of the Case 3 R2R controller on the Gaussian CMP process when employing truncated linear decay and exponential decay based GRI as a function of β . Results are computed over 100 Monte Carlo simulations (dashed lines indicate one standard deviation)

V. CONCLUSIONS

This paper has proposed the use of GPR to implement a VM-enabled EWMA R2R control scheme that adapts its response based on the confidence in the prediction provided by the VM model. A similar approach has been proposed by [22] using the externally computed RI and GSI to provide a measure of the confidence in model predictions. The advantage of employing GPR is that, as a probabilistic modeling paradigm, it directly generates a mean prediction and variance estimate. Hence, the GPR based solution is self-contained, providing both the VM prediction and the variance estimate used to adjust the coefficient of the EWMA controller.

Simulation results for a CMP case study confirm that, as expected, the variance estimate provided by GPR is indicative of the accuracy of the VM prediction, and comparisons with the RIs suggest that GRIs can be a stronger indicator of confidence in predictions than RIs, especially when the controller operating space is not adequately represented by the

training space used to estimate the VM model, and hence the accuracy of the VM model is limited. Such scenarios are likely to be common place in practice due to the challenges with collecting suitable datasets for training VM models *a priori* in a production environment. The proposed GRI confidence information enhanced VM-enabled R2R controllers can achieve better control performance over a wide range of α_0 values, providing a degree of robustness to sub-optimal selection of this control parameter.

It should be noted that in this treatise only feedback (post-metrology) R2R is considered. Future work will explore the use of GRI information with feedforward VM (pre-metrology) R2R control. Other areas to investigate include incorporating on-line VM model updating into the VM-enhanced R2R control framework using confidence information informed dynamic metrology sampling, and developing strategies for VM model maintenance in the presence of process drifts and maintenance shifts that are not reflected in the model inputs.

ACKNOWLEDGMENT

This work was supported by the Irish Centre for Manufacturing Research (ICMR) and Enterprise Ireland (Grant: CC/2011/2001).

REFERENCES

- [1] J. Moyne, E. del Castillo, and A. M. Hurwitz, *Run-to-Run Control in Semiconductor Manufacturing*. CRC Press, 2001.
- [2] Y. Wang, F. Gao, and F. J. D. III, "Survey on iterative learning control, repetitive control, and run-to-run control," *Journal of Process Control*, vol. 19, no. 10, pp. 1589 – 1600, 2009.
- [3] J. Ringwood, S. Lynn, G. Bacelli, B. Ma, E. Ragnoli, and S. McLoone, "Estimation and control in semiconductor etch: Practice and possibilities," *Semiconductor Manufacturing, IEEE Transactions on*, vol. 23, pp. 87–98, Feb 2010.
- [4] J. Moyne, "Run-to-run control in semiconductor manufacturing," in *Encyclopedia of Systems and Control* (J. Baillieul and T. Samad, eds.), pp. 1248–1254, Springer London, 2015.
- [5] F. Tan, T. Pan, Z. Li, and S. Chen, "Survey on run-to-run control algorithms in high-mix semiconductor manufacturing processes," *IEEE Transactions on Industrial Informatics*, vol. 11, pp. 1435–1444, Dec 2015.

- [6] E. Sachs, A. Hu, and A. Ingolfsson, "Run by run process control: combining spc and feedback control," *Semiconductor Manufacturing, IEEE Transactions on*, vol. 8, no. 1, pp. 26–43, 1995.
- [7] M. D. Ma and J. Y. Li, "Improved variable ewma controller for general arima processes," *IEEE Transactions on Semiconductor Manufacturing*, vol. 28, pp. 129–136, May 2015.
- [8] S. W. Butler and J. Stefani, "Supervisory run-to-run control of polysilicon gate etch using in situ ellipsometry," *Semiconductor Manufacturing, IEEE Transactions on*, vol. 7, no. 2, pp. 193–201, 1994.
- [9] E. D. Castillo and R. Rajagopal, "A multivariate double ewma process adjustment scheme for drifting processes," *IIE Transactions*, vol. 34, no. 2, pp. 1055–1068, 2002.
- [10] E. Del Castillo and J.-Y. Yeh, "An adaptive run-to-run optimizing controller for linear and nonlinear semiconductor processes," *Semiconductor Manufacturing, IEEE Transactions on*, vol. 11, no. 2, pp. 285–295, 1998.
- [11] C. Zhang, H. Deng, and J. S. Baras, "Run-to-run control methods based on the dhobe algorithm," *Automatica*, vol. 39, pp. 35–45, 2003.
- [12] P. Kang, D. Kim, H. joo Lee, S. Doh, and S. Cho, "Virtual metrology for run-to-run control in semiconductor manufacturing," *Expert Systems with Applications*, vol. 38, no. 3, pp. 2508 – 2522, 2011.
- [13] A. Khan, J. R. Moyne, and D. Tilbury, "An approach for factory-wide control utilizing virtual metrology," *Semiconductor Manufacturing, IEEE Transactions on*, vol. 20, no. 4, pp. 364–375, 2007.
- [14] A. A. Khan, J. Moyne, and D. Tilbury, "Virtual metrology and feedback control for semiconductor manufacturing processes using recursive partial least squares," *Journal of Process Control*, vol. 18, no. 10, pp. 961 – 974, 2008.
- [15] F.-T. Cheng, C.-A. Kao, C.-F. Chen, and W.-H. Tsai, "Tutorial on applying the vm technology for tft-lcd manufacturing," *Semiconductor Manufacturing, IEEE Transactions on*, vol. 28, pp. 55–69, Feb 2015.
- [16] S.-K. S. Fan and Y.-J. Chang, "An integrated advanced process control framework using run-to-run control, virtual metrology and fault detection," *Journal of Process Control*, vol. 23, no. 7, pp. 933 – 942, 2013.
- [17] F.-T. Cheng, J.-C. Chang, H.-C. Huang, C.-A. Kao, Y.-L. Chen, and J.-L. Peng, "Benefit model of virtual metrology and integrating avm into mes," *Semiconductor Manufacturing, IEEE Transactions on*, vol. 24, pp. 261–272, May 2011.
- [18] T. Tsuda, S. Inoue, A. Kayahara, S. i. Imai, T. Tanaka, N. Sato, and S. Yasuda, "Advanced semiconductor manufacturing using big data," *IEEE Transactions on Semiconductor Manufacturing*, vol. 28, pp. 229–235, Aug 2015.
- [19] K.-Y. Lin, C.-Y. Hsu, and H.-C. Yu, "A virtual metrology approach for maintenance compensation to improve yield in semiconductor manufacturing," *International Journal of Computational Intelligence Systems*, vol. 7, no. sup2, pp. 66–73, 2014.
- [20] F.-T. Cheng, Y.-T. Chen, Y.-C. Su, and D.-L. Zeng, "Evaluating reliance level of a virtual metrology system," *Semiconductor Manufacturing, IEEE Transactions on*, vol. 21, no. 1, pp. 92–103, 2008.
- [21] C.-A. Kao, F.-T. Cheng, and W.-M. Wu, "Preliminary study of run-to-run control utilizing virtual metrology with reliance index," in *Automation Science and Engineering (CASE), 2011 IEEE Conference on*, pp. 256–261, Aug 2011.
- [22] C.-A. Kao, F.-T. Cheng, W.-M. Wu, F.-W. Kong, and H.-H. Huang, "Run-to-run control utilizing virtual metrology with reliance index," *Semiconductor Manufacturing, IEEE Transactions on*, vol. 26, no. 1, pp. 69–81, 2013.
- [23] G. Susto, A. Schirru, S. Pampuri, G. De Nicolao, and A. Beghi, "An information-theory and virtual metrology-based approach to run-to-run semiconductor manufacturing control," in *Automation Science and Engineering (CASE), 2012 IEEE International Conference on*, pp. 358–363, 2012.
- [24] S. Hwang, M. Jeong, and B.-J. Yum, "Robust relevance vector machine with variational inference for improving virtual metrology accuracy," *Semiconductor Manufacturing, IEEE Transactions on*, vol. 27, pp. 83–94, Feb 2014.
- [25] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. The MIT Press, 2006.
- [26] J. Wan, S. Pampuri, P. O'Hara, A. Johnston, and S. McLoone, "On regression methods for virtual metrology in semiconductor manufacturing," in *Irish Signals Systems Conference 2014 and 2014 China-Ireland International Conference on Information and Communications Technologies (ISSC 2014/CICT 2014). 25th IET*, pp. 380–385, June 2014.
- [27] M. Ebden, "Gaussian processes for regression: a quick introduction," tech. rep., 2008.
- [28] S. Lynn, J. Ringwood, and N. Macgearailt, "Gaussian process regression for virtual metrology of plasma etch," in *Signals and Systems Conference (ISSC 2010), IET Irish*, pp. 42–47, 2010.
- [29] S.-T. Tseng, A. B. Yeh, F. Tsung, and Y.-Y. Chan, "A study of variable ewma controller," *IEEE Transactions on Semiconductor Manufacturing*, vol. 16, pp. 633–643, Nov 2003.
- [30] J. Iskandar and J. Moyne, "Maintenance of virtual metrology models," in *2016 27th Annual SEMI Advanced Semiconductor Manufacturing Conference (ASMC)*, pp. 393–398, May 2016.



Jian Wan (S'06-M'12) received the B.Eng and M.Eng in Engineering from Northwestern Polytechnical University, Xi'an, China, in 1997 and 2000, respectively. He obtained his Ph.D. degree in Control Engineering from the University of Girona, Catalonia, Spain in 2007. He was a Research Assistant from 2007 to 2009 at the University of Leeds, UK, a Research Associate from 2009 to 2011 at the University of Manchester, UK, a Research Fellow from 2012 to 2014 at Maynooth University, Ireland, and also a Research Associate in 2014 at the University of Strathclyde, UK. He joined the University of Plymouth, UK as a lecturer in control systems engineering in 2015. His research interests include statistical processing monitoring and control, set-membership methods for control, constrained optimization and control with applications to manufacturing processes, renewable energy, robotics and autonomous systems.



Seán McLoone (S'94-M'96-SM'02) received the M.Eng. degree (Hons.) in Electrical and Electronic Engineering and the Ph.D. degree in Control Engineering from Queen's University Belfast (QUB), Belfast, U.K., in 1992 and 1996, respectively. He was a Post-Doctoral Research Fellow, from 1996 to 1997, and a Lecturer, from 1998 to 2002, with QUB. He joined the Department of Electronic Engineering, NUI Maynooth, Maynooth, Ireland, in 2002, where he served as a Senior Lecturer, from 2005 to 2012, and as the Head of Department, from 2009 to 2012,

before returning to QUB in 2013 to take up his current post as a Professor and the Director of the Energy Power and Intelligent Control Research Cluster in the School of Electronics, Electrical Engineering, and Computer Science. His research interests are in the general area of intelligent systems, with a particular focus on data based modeling and analysis of dynamical systems. This encompasses techniques ranging from classical system identification, fault diagnosis and statistical process control to modern computational intelligence based adaptive learning algorithms and optimization techniques. His research has a strong application focus, with many projects undertaken in collaboration with industry in areas such as process monitoring, control and optimization, time-series prediction, and inline sensor characterization. Prof. McLoone is a Chartered Engineer and a Fellow of the Institution of Engineering and Technology. He is a Past Chairman of the U.K. and Republic of Ireland Section of the IEEE.