

An Evaluation of Multiple-input Dual-output Run-to-Run Control Scheme for Semiconductor Manufacturing

Shu-Kai S. Fan[†] · Yen Lin

Department of Industrial Engineering and Management
Yuan Ze University, Taoyuan County, Taiwan 320, Republic of China
Tel: +886-3-4638800-2510, E-mail: simonfan@saturn.yzu.edu.tw

Abstract. This paper provides an evaluation of an optimization-based, multiple-input double-output (MIDO) run-to-run (R2R) control scheme for general semiconductor manufacturing processes. The controller in this research, termed adaptive dual response optimizing controller (ADROC), can serve as a process optimizer as well as a recipe regulator between consecutive runs of wafer fabrication. In evaluation, it is assumed that the equipment model could be appropriately described by a pair of second-order polynomial functions in terms of a set of controllable variables. Of practical relevance is to consider a drifting effect in the equipment model since in common semiconductor practice the process tends to drift due to machine aging and tool wearing. We select a typical application of R2R control to chemical mechanical planarization (CMP) in semiconductor manufacturing in this evaluation, and there are five different CMP process scenarios demonstrated, including mean shift, variance increase, and IMA disturbances. For the controller, ADROC, an on-line estimation technique is implemented in a self-tuning (ST) control manner for the adaptation purpose. Subsequently, an *ad hoc* global optimization algorithm based on the dual response approach, arising from the response surface methodology (RSM) literature, is used to seek the optimum recipe within the acceptability region for the execution of next run. The main components of ADROC are described and its control performance is assessed. It reveals from the evaluation that ADROC can provide excellent control actions for the MIDO R2R situations even though the process exhibits complicated, nonlinear interaction effects between control variables, and the drifting disturbances.

Keywords: Production Management, Non-linear Programming, Optimization

1. INTRODUCTION

It is well known from control engineering that the goal of an adaptive extremum controller is to seek the optimum setting of control variables that keeps the process output at the extremum value and then continuously fine-tunes the process operating at its optimum despite model mismatch and/or the influence of system dynamics and noise disturbances.

In contrast to process monitoring applied in the phase of conventional control charting for the detection and thus removal of assignable causes of variation, process adjustment is referred to as another variability-reduction tool for the process compensation or regulation, also known as engineering process control (EPC), where an adjustment to manipulable process variables is made in an attempt to keep the process output as close as possible on

some target value. One main group of EPC designs involves the notion of feedback control, a most recent application of which to microelectronics industry leads to run-to-run (R2R) control schemes (del Castillo and Hurwitz 1997) in semiconductor manufacturing.

Successful implementation of the most popular R2R controllers in semiconductor community lies in a fundamental assumption that the functional relationship associating the compensating (or recipe) variables and process output of interest is in the form of linearity (Sachs, Hu and Ingolfsson 1995 and Patel and Jenkins 2000). It implies that if this assumption holds in general, then several complex semiconductor processes can be well represented by the transfer function models (Box, Jenkins and Reinsel 1994) fitted from experimental data collected in a pre-control stage. However, nonlinear effects among semiconductor process variables often seem critical; for

[†] : Corresponding Author

example, the chemical mechanical planarization (CMP) process is a typical case of that kind (del Castillo and Yeh 1998 and Fan 2000a). If nonlinearities were noticeable, the existing R2R regulators (such as EWMA-based controllers) would not be adequate for the control purposes. The traditional method of physical principle models is ideal for the specific tasks of electronics-model building; yet sufficient knowledge of the chemistry, physics and the internal dynamics is necessary to build a model (Nanz and Camilletti 1995). In semiconductor practice, these conditions are rarely completely met or the physical models may not even exist. Therefore, if nonlinear complexity can be taken into consideration while designing a controller, better control performance can be anticipated as compared to the linear ones. The previous statement is one of the primary concerns addressed in this research.

Nearly all the R2R controllers are model-referenced, denoting that a model is first envisioned as a base from which to devise a controller. In the vast majority of realistic environments, multiple outputs are usually entailed, whereas one-attribute-at-a-time control will be a sensible choice that is much more likely to be carried out from an empirical point of view. In light of these features, a multiple-input dual-output (MIDO) R2R optimizing controller for semiconductor manufacture is presented in this paper, serving as a supervisory recipe regulator between batches (or runs) of silicon wafer production. The proposed controller is dubbed adaptive dual response optimizing controller (ADROC) where a variant of the second-order (quadratic) response surface model (Myers and Montgomery 2002) is adopted and an optimization algorithm is then developed to anchor an optimum "setpoint" of control variables for a subsequent run. In essence, ADROC follows the concept of "self-optimizing control" (or called "adaptive extremum control"), as will be seen shortly. For its details, see Golden and Ydstie (1989).

The rest of this paper is laid out as follows. Section 2 describes the ingredients of ADROC, separating into subsections of model building, optimization algorithm, control model description, and on-line estimation. Several important implementation issues and the system architecture of ADROC along with a block diagram are summarized in Section 3. Section 4 is devoted to a series of simulation studies (based upon an advanced polishing process drawn from the semiconductor industry), which illustrate the ADROC's behavior in typical experimental situations. There, the computational results of related performance measures are reported as well. Conclusions are drawn in Section 5 with some remarks of areas for further research.

2. COMPONENTS OF ADROC

For the extremum seeking functionality, an *ad hoc* constrained quadratic programming algorithm, suitable

for solving the dual response systems (DRS) arising from response surface methodology (RSM, see, *e.g.*, Khuri and Cornell 1996 and Myers and Montgomery 2002), is utilized to calculate the recipe from run to run. For the adaptation mechanism, a generalized multivariate recursive least squares (RLS) algorithm is utilized for the on-line estimation. To be discussed next is the process model considered in this study.

2.1 Process Model Description

Myers and Carter (1973) develop a useful methodology termed "dual response approach" in the case of two responses of interest. It is also assumed that the process engineer can clearly categorize these two responses by their importance as the "primary" and "secondary" response variables. The goal is to seek the optimum condition $\mu^* = [\mu_1^*, \mu_2^*, \dots, \mu_k^*]^T$ on the k control variables that optimizes the primary response while keeping the secondary response on the target value, as defined by

$$\begin{aligned} \text{Min} \quad & \hat{y}_p(\mu) = \hat{\beta}_p + \hat{\beta}'_p \mu + \mu' \hat{\mathbf{B}}_p \mu \\ \text{s.t.} \quad & \hat{y}_s(\mu) = \hat{\beta}_s + \hat{\beta}'_s \mu + \mu' \hat{\mathbf{B}}_s \mu = T \\ & \mu' \mu \leq \rho^2 \end{aligned} \quad (1)$$

where \hat{y}_p is the fitted quadratic primary response function, \hat{y}_s is the fitted quadratic secondary response function, and T is the target value for \hat{y}_s . Henceforth, the subscripts/superscripts p and s denote primary and secondary, respectively. In the regression functions \hat{y}_p and \hat{y}_s , $\hat{\beta}_p$ and $\hat{\beta}_s$ contain the parameter estimates of the constant terms (or intercepts) β_p and β_s ; $\hat{\beta}'_p$ and $\hat{\beta}'_s$ contain the parameter estimates of the first-order (or linear) terms β_p and β_s , as can be written by

$$\beta_p = [\beta_{p,1}, \beta_{p,2}, \dots, \beta_{p,k}]', \beta_s = [\beta_{s,1}, \beta_{s,2}, \dots, \beta_{s,k}]'$$

$\hat{\mathbf{B}}_p$ and $\hat{\mathbf{B}}_s$ contain the parameter estimate of the second-order (or quadratic) terms \mathbf{B}_p and \mathbf{B}_s as can be written by

$$\mathbf{B}_p = \begin{bmatrix} \beta_{p,11} & \beta_{p,12}/2 & \dots & \beta_{p,1k}/2 \\ \beta_{p,21}/2 & \beta_{p,22} & \dots & \beta_{p,2k}/2 \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{p,k1}/2 & \beta_{p,k2}/2 & \dots & \beta_{p,kk} \end{bmatrix}_{k \times k},$$

$$\mathbf{B}_s = \begin{bmatrix} \beta_{s,11} & \beta_{s,12}/2 & \dots & \beta_{s,1k}/2 \\ \beta_{s,21}/2 & \beta_{s,22} & \dots & \beta_{s,2k}/2 \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{s,k1}/2 & \beta_{s,k2}/2 & \dots & \beta_{s,kk} \end{bmatrix}_{k \times k}.$$

Note that the matrices \mathbf{B}_p and \mathbf{B}_s are symmetric. For the inequality constraint in equation (1), ρ is the radius used to restrict the search on the variables μ inside the spherical experimental region where the responses were fitted.

Here, we call model (1) as the “dual response systems (DRS).” Partly motivated by the ridge analysis procedure developed by Draper (1963), the dual response approach introduced by Myers and Carter (1973) is a contour-based approach, which, essentially, is constructed based on the Lagrange method in the context of constrained optimization from nonlinear programming (NLP). Like ridge analysis, the dual response approach produces a locus of coordinates of control variables where various values of the predicted secondary response function \hat{y}_s are considered.

2.2 Constrained Nonlinear Optimization

Assuming a suitable constraint qualification (*i.e.*, the gradient of the binding constraints are linearly independent (see, *e.g.*, Luenberger 1984)), a necessary condition for local optimality of a feasible point μ is the existence of Lagrange multipliers μ, θ such that

$$(\mathbf{B}_p - \mu\mathbf{B}_s + \theta\mathbf{I})\mu = \frac{1}{2}(\mu\hat{\beta}_s - \hat{\beta}_p) \quad (2)$$

where μ, θ additionally satisfy $\theta \geq 0$ whenever $\mu'\mu < \rho^2$. If the matrix $(\mathbf{B}_p - \mu\mathbf{B}_s + \theta\mathbf{I})$ is positive definite, then it can be shown that μ is a global optimum for (1) (Fan 2000b). The latter result suggests taking values for μ, θ that ensure the matrix $(\mathbf{B}_p - \mu\mathbf{B}_s + \theta\mathbf{I})$ positive definite. This technique was recommended as part contour-based method in Myers and Carter (1973), and later developed into a formal algorithm, DRSALG (Semple 1997).

2.2.1 DRSALG

DRSALG was designed to search the (convex) region Γ define by

$$\Gamma = \{(\mu, \theta): (\mathbf{B}_p - \mu\mathbf{B}_s + \theta\mathbf{I}) \text{ is positive definite, } \theta \geq 0\} \quad (3)$$

for values of μ, θ that make the μ obtained from (2) feasible in (1). If μ is fixed, then a related and easier problem is to determine values for θ and μ that solve the parametric trust region subproblem

$$\begin{aligned} \text{Min } \hat{y}_p(\mu) &= \hat{\beta}_p' + \hat{\beta}_p'\mu + \mu'\hat{\mathbf{B}}_p\mu \\ &\quad - \mu(\hat{\beta}_s' + \hat{\beta}_s'\mu + \mu'\hat{\mathbf{B}}_s\mu - T) \quad (4) \\ \text{s.t. } \mu'\mu &\leq \rho^2 \end{aligned}$$

(4) has a stationary equation identical to (2), but now the stationary solution μ is only required to satisfy the radial inequality. The sole Lagrange multiplier, $\theta \geq 0$, corresponds to this inequality and must additionally satisfy $\theta \cdot (\mu'\mu - \rho^2) = 0$ and make $(\mathbf{B}_p - \mu\mathbf{B}_s + \theta\mathbf{I})$ positive semidefinite at a global optimum (see, *e.g.*, Moré and Sorensen 1983). If $(\mathbf{B}_p - \mu\mathbf{B}_s + \theta\mathbf{I})$ is positive definite, then the solution to (4), denoted by $\mu^*(\mu)$, is the unique global optimum to (4). Observe that $\mu^*(\mu)$ will not, in general, be feasible for (1), *i.e.*, $\hat{y}_s(\mu^*(\mu)) - T \neq 0$. However, should μ be determined such that $\hat{y}_s(\mu^*(\mu)) - T = 0$, then $\mu^*(\mu)$ solves (1) as well. Calculating μ so that $\hat{y}_s(\mu^*(\mu)) - T = 0$ can be accomplished whenever the function $\hat{y}_s(\mu^*(\mu))$ is continuous on some interval $[a, b]$ with $\hat{y}_s(\mu^*(a)) < T < \hat{y}_s(\mu^*(b))$. Continuity, in turn, hinges on the eigenstructure of the matrix $(\mathbf{B}_p - \mu\mathbf{B}_s)$. If the eigenvalues of this matrix are ordered $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ (where the dependence on μ has been suppressed for notational simplicity), then it can be shown that $\hat{y}_s(\mu^*(\mu))$ is continuous on $[a, b]$ provided $(\mu\hat{\beta}_s - \hat{\beta}_p)$ is never perpendicular to the first eigenspace define by $E_1 = \{q: (\mathbf{B}_p - \mu\mathbf{B}_s)q = \lambda_1 q\}$ on this interval (see, *e.g.*, Semple 1997), where the eigenvector q associated with the smallest eigenvalue λ_1 of $(\mathbf{B}_p - \mu\mathbf{B}_s)$.

DRSALG algorithm guarantees global optimal solutions for nondegenerate dual response problems. However, even the majority of dual response problems are nondegenerate, the degenerate problems can occur and cannot be solved by DRSALG. del Castillo, Fan and Semple (1999) devised a generalized algorithm, called DR2, that computes global optimal solutions for nondegenerate problems and approximate global optimal solutions for degenerate problems.

2.2.2 DR2

In order to solve the degenerate case, a single axis-henceforth termed the *grid axis*-is selected. The variable associated with this axis will be termed the grid variable. Values of the grid variable are restricted to the interval $[-\rho, \rho]$. Since ρ^2 is bounded by n in dual response optimization problems, this interval is relatively narrow. The interval $[-\rho, \rho]$ along the grid axis can then be divided into regular subintervals. Each interstitial point-henceforth termed the *grid point*-represents a value where the grid variable will be temporarily fixed. When the grid variable is fixed at a grid point, a subproblem of

the form (1) having $n-1$ process variables is created. New matrices and vectors, each deflated by one dimension and adjusted to account for the fixed grid variable, replace the old values for $\hat{\beta}_p$, $\hat{\beta}_s$, $\hat{\mathbf{B}}_p$, and $\hat{\mathbf{B}}_s$ in (1). The values of ρ^2 , $\hat{\beta}_p$, and $\hat{\beta}_s$ require similar adjustments. Each $n-1$ dimensional subproblem is solved, and the best solution found among the subproblems is recorded. If desired, the grid can be refined locally around the best solution found to improve accuracy. Since degeneracy indicates $(\mu_d \beta_s - \beta_p) \perp \mathbf{E}_1^{n-1}$, thus the original axes, say $\mu^* = [\mu_1^*, \mu_2^*, \mu_3^*, \mu_4^*, \dots, \mu_k^*]^T$, are rotated so that the eigen-vector \mathbf{q}_1 associated with the smallest eigenvalue λ_1 of $(\mathbf{B}_p - \mu_d \mathbf{B}_s)$ becomes the first axis \mathbf{z}_1 in the new coordinate system based on the orthonormal basis $\mathbf{Q} = [\mathbf{q}_1 : \mathbf{q}_2 : \dots : \mathbf{q}_k]$. Herein, \mathbf{q}_1 can be quickly computed via a few steps of Inverse Iteration. Then, \mathbf{Q} can be formed by utilizing the Gram-Schmidt orthogonalization procedure (see, e.g., Golub and van Loan 1984). By the rotational transformation via $\mu = \mathbf{Q}\mathbf{z}$ illustrated in figure 1, (1) becomes

$$\begin{aligned} \text{Min } & \hat{y}_p(\mathbf{z}) = \hat{\beta}_p + \hat{\beta}'_p \mathbf{z} + \mathbf{z}' \hat{\mathbf{B}}_p \mathbf{z} \\ \text{s.t. } & \hat{y}_s(\mathbf{z}) = \hat{\beta}_s + \hat{\beta}'_s \mathbf{z} + \mathbf{z}' \hat{\mathbf{B}}_s \mathbf{z} = T \\ & \mathbf{z}' \mathbf{z} \leq \rho^2 \end{aligned} \quad (5)$$

where $\tilde{\beta}'_p = \beta'_p \mathbf{Q}$, $\tilde{\mathbf{B}}_p = \mathbf{Q}' \mathbf{B}_p \mathbf{Q}$, $\tilde{\beta}'_s = \beta'_s \mathbf{Q}$, and $\tilde{\mathbf{B}}_s = \mathbf{Q}' \mathbf{B}_s \mathbf{Q}$. The radial constraint does not change in the new \mathbf{z} -coordinates because \mathbf{Q} is orthonormal ($\mathbf{Q}' \mathbf{Q} = \mathbf{I}$).

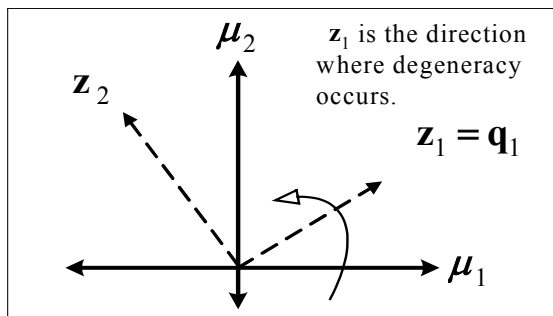


Figure 1. The Rotation Procedure

Degeneracy results from orthogonality on the direction of \mathbf{q}_1 , so removing \mathbf{z}_1 's effect from (5) can effectively rectify this computational difficulty around μ_d . Consequently, \mathbf{z}_1 is selected as the grid variable; namely, \mathbf{z}_1 satisfies $-\rho \leq \mathbf{z}_1 \leq \rho$ and this interval can be equally divided into many subintervals by many grid points (illustrated in figure 2). If \mathbf{z}_1 is, in turn, fixed at each grid point, then (5) is decomposed into a series of

DR subproblems having $(k-1)$ control factors. New matrices, vectors, intercepts, and radii need adjustments to account for the fixed grid variable. Each subproblem with lower dimensions is created and swiftly solved by DRSALG, and the best solution found among the subproblems is earmarked.

In order to increase accuracy, a new working (or bracketing) interval for \mathbf{z}_1 is constructed by two grid points adjacent to the best solution obtained in the first pass of grid search. Then, the best solution found in the second pass (local refinement) can be transformed back to the original variables through $\mu = \mathbf{Q}\mathbf{z}$.

DR2 will find global optimal solutions of the nondegenerate DR problems by calling DRSALG as a subroutine, and will use the procedures aforementioned to return an approximate global optimal solution (due to the computational accuracy dictated by the mesh used in the grid search) in degenerate cases. While solving too many subproblems would appear to compromise the speed of the procedure, it has been found that this is not the case in practice. First, many of the subproblems are infeasible, and these are screened quickly and discarded in the initial phase of DRSALG. Second, if the grid axis is constructed carefully, the majority of subproblems will be nondegenerate, and convergence will be quite swift.

In contrast to DRSALG for nondegenerate problem, it is impossible to claim that grid point procedure (del Castillo, Fan and Semple 1999) will produce an exact global optimal solution for degenerate problems. However, since the grid method surveys as much of the feasible region as possible, it is reasonable to expect a more accurate approximation to the global optimum than would be obtained using a local search procedure. Since each nondegenerate problem is actually solved in the initial phase as screen for degeneracy, the implementation is already a comprehensive solver for general quadratic dual response systems. For each nondegenerate problem, the DR2 solution satisfies the sufficient conditions for global optimality. For degeneracy case, the fact that DR2 lacks of speed on this type of problems is not surprising giving the sheer number of subproblems that are solved. Fortunately, nondegenerate problems occur quite frequently in common practice.

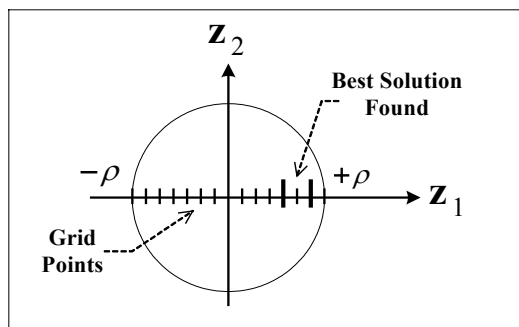


Figure 2. Grid Point Search in the Direction \mathbf{z}_1

2.3 Control Model Description

Most of the results available on extremum control deem an engineering process (or scientific system) as a static map, and the static input-output response surface at each discrete time period is nonlinear in nature. In modelbased optimization control schemes, the uncertainties existing in the input-output mapping make it necessary to use some sort of “adaptation” to determine the unknown parameters of the referenced model, and then on-line identification techniques such as the recursive least squares (RLS) apply. Invoking the certainty equivalence principle in the control literature, a specific optimization algorithm is provided to find the best operating point of a process.

2.3.1 Multiple-Input Dual-Output Control Model

To lend the static DRS itself to an R2R formulation, now we consider the MIDO system at discrete time t , as defined by

$$\begin{aligned} y_p(t) &= \beta_p(t) + \beta'_p(t)\mu(t-1) \\ &\quad + \mu'(t-1)\mathbf{B}_p(t)\mu(t-1) + \eta_p(t), \\ y_s(t) &= \beta_s(t) + \beta'_s(t)\mu(t-1) \\ &\quad + \mu'(t-1)\mathbf{B}_s(t)\mu(t-1) + \eta_s(t), \\ t &= 0, 1, 2, \dots, \end{aligned} \quad (6)$$

where $y_p(t)$ and $y_s(t)$ denote the “actual” primary outputs and secondary outputs at time t ; $\eta_p(t)$ and $\eta_s(t)$ denote the process disturbances and/or model misspecifications at time t unaccounted for by using the quadratic primary- and secondary-process functions. As can clearly be seen from (6), the process outputs produced at time t are subject to the recipe of previous run (*i.e.*, $\mu(t-1)$), showing time lag 1 in a feedback control sense. This is a typical EPC convention applied in R2R applications—a “dead-beat” control policy. Thus, EPC practitioner may view R2R control as a supervisory controller that purely adjusts the set-point of the machine tool controller in trying to bring the process output back to target from run to run. Note that, in semiconductor manufacturing, a run (or batch) is referred to a specific process, which deals with a single wafer, a single “lot” of wafers (typically 25 wafers in cassette), or several lots at once. If the DRS as in (1) is assumed adequate to describe the process, the computation of recipe update between runs is equivalent to solving the following constrained quadratic programming problem:

$$\begin{aligned} \text{Min.} \quad & \hat{y}_{p,t+|t} = \hat{\beta}_{p,t} + \hat{\beta}'_{p,t}\mu_t + \mu'_t \hat{\mathbf{B}}_{p,t} \mu_t \\ \text{s.t.} \quad & \hat{y}_{s,t+|t} = \hat{\beta}_{s,t} + \hat{\beta}'_{s,t}\mu_t + \mu'_t \hat{\mathbf{B}}_{s,t} \mu_t = T \\ & \mu'_t \mu_t \leq \rho^2, \quad t = 0, 1, 2, \dots, \end{aligned} \quad (7)$$

where $\hat{y}_{p,t+|t}$ and $\hat{y}_{s,t+|t}$ indicate the predicted process

models, for the next discrete time $t+1$, built upon the control outputs up to time t ; $\mu(t)$ is the computed set-point at time t and will serve as the recipe of running the next batch at time $t+1$. The radius ρ in the inequality constraint can be considered as the allowable adjustment range pre-specified by the control engineer, which is sometimes subject to the practical limitation of equipment setting.

The concept of adaptive extremum control is essentially related to optimization techniques, many of which have been transferred from numerical optimization. In this study, we use the Dual Response II (DR2) algorithm presented in Fan (2000b) to solve the optimization problem in (7). DR2 solves the DRS by two different ways. If the DRS is a nondegenerate case, then the algorithm DRSALG (Semple 1997) directly solves the problem and guarantees a unique “global” solution. If the DRS is detected to be numerically degenerate, then the procedure AXIS rotates a degenerate problem and then decomposes it into a finite sequence of non-degenerate sub-problems of lower dimension. The non-degenerate subproblems are solved by DRSALG. The simulation analysis reveals that, for degenerate cases, DR2 obtains the global solution over 98% of the time. Computational results based on large simulations also show that DR2 is more effective at locating global (or near-global) solutions for the DRS than several optimization algorithms (such as the Generalized Reduced Gradient (GRG) algorithm and the Sequential Quadratic Programming (SQP) algorithm) that have been frequently used in RSM applications. The algorithms DRSALG and DR2 involve a great deal of algebraic mathematics and optimization methodology. For detailed discussion, interested readers can refer to Semple (1997) for DRSALG, and to del Castillo, Fan and Semple (1999) and Fan (2000b) for DR2.

2.4 On-Line Estimation Technique

The self-tuning (ST) control system, based upon an idea of separating the parameter estimation from the controller design tasks (known as the “separation theorem” in the classical control theory), has been widely applied (see, *e.g.*, Åström and Wittenmark 1973, and Seborg *et al.* 1986). The distinctive characteristic of combining a recursive estimation algorithm with a controller synthesis is due to its suitability for processes that vary with time. That is, the process under investigation can be assumed to have time-varying process parameters or has constant but initially unknown process parameters, meaning that ST mechanism might be advantageous for discrete-part manufacturing (DPM). Then, the control function uses these estimates as if they were “true” in the referenced model, called the certainty equivalence prin-

ciple. There have various controller design strategies involved in the controller synthesis stage, such as minimum-variance (MV), moving-average (MA), pole-placement, LQG, *etc.* However, extremum-seeking methods are still of great practical interest since even small improvement in the performance function can possibly result in large savings in manufacturing cost. In addition, recent development in computers has led to a rejuvenated attention in extremum-seeking method combined with adaptive control (Åström and Wittenmark 1995).

At the core of ST algorithms is the recursive estimation method. The assumed model is linear in parameter; therefore, the ordinary recursive least squares (RLS) algorithm will be adequate for this estimation situation. As a result, the input data in the TR-algorithm should be replaced with $\hat{\beta}_{p,t}$, $\hat{\beta}_{s,t}$, $\hat{\beta}'_{p,t}$, $\hat{\beta}'_{s,t}$, $\hat{\mathbf{B}}_{p,t}$, and $\hat{\mathbf{B}}_{s,t}$, denoting their LS estimates at iteration t . Herein, a clear-cut scheme is opted, which estimates the parameters in the model and then these estimates are indirectly used for function optimization, therefore leading to the so-called indirect (or explicit) self-tuning control (Åström and Wittenmark 1995). To prevent an abrupt deterioration of estimation performance due to the “parameter windup” phenomenon, an RLS with constant trace (Shah and Cluett 1991) is employed and its formal procedure with respect to (7) is expressed shortly.

In a sense of feedback control, the predicted function for (7) can be expressed by a simplified general linear model (GLM) of $\hat{y}_{p,t} = \hat{\mathbf{a}}'_t \boldsymbol{\mu}_{t-1}$ and $\hat{y}_{s,t} = \hat{\mathbf{b}}'_t \boldsymbol{\mu}_{t-1}$, where $\hat{\mathbf{a}}_t$ is a $[(n^2 + 3n + 2)/2] \times 1$ vector, which contains the parameter estimates $\hat{\beta}_{p,t}$, $\hat{\beta}'_{p,t}$ and $\hat{\mathbf{B}}_{p,t}$; $\hat{\mathbf{b}}'_t$ is also a $[(n^2 + 3n + 2)/2] \times 1$ vector, which contains the parameter estimates $\hat{\beta}_{s,t}$, $\hat{\beta}'_{s,t}$ and $\hat{\mathbf{B}}_{s,t}$ at discrete time t , if the full quadratic model is considered. The RLS algorithm with constant trace can thus be formulated as follows:

Initialization: Let $t=1$ and $(n^2 + 3n + 2)/2 = p$. Let $\boldsymbol{\mu}_0$ be the initial recipe at time 0, \mathbf{P}_0 an initial matrix and λ the constant of discounting factor.

- 1) $\mathbf{K}_t = \frac{\mathbf{P}_{t-1} \boldsymbol{\mu}_{t-1}}{\lambda + \boldsymbol{\mu}'_{t-1} \mathbf{P}_{t-1} \boldsymbol{\mu}_{t-1}}$;
- 2) $\hat{\mathbf{a}}_t = \hat{\mathbf{a}}_{t-1} + \mathbf{K}_t (y_{p,t} - \boldsymbol{\mu}'_{t-1} \hat{\mathbf{a}}_{t-1})$;
 $\hat{\mathbf{b}}_t = \hat{\mathbf{b}}_{t-1} + \mathbf{K}_t (y_{s,t} - \boldsymbol{\mu}'_{t-1} \hat{\mathbf{b}}_{t-1})$;
- 3) $\mathbf{P}_t = \frac{(\mathbf{I}_p - \mathbf{K}_t \boldsymbol{\mu}'_{t-1}) \mathbf{P}_{t-1}}{\lambda} + \frac{\mathbf{I}_p (\mathbf{K}'_t \mathbf{P}_{t-1} \boldsymbol{\mu}_{t-1})}{p}$;
- 4) $t := t + 1$ and go to (1); end.

In the foregoing algorithm, the discounting factor λ

with a value less than 1.0 is very useful for tracking the parameters in time-varying systems and during the initial transient phase of self-tuning. The usual range of λ is set between 0.95 and 1.0. \mathbf{K}_t is a $p \times 1$ vector of weights and \mathbf{P}_t is a $p \times p$ matrix proportional to the variance-covariance matrix of the parameter estimates. Measurements of a performance function in extremum control are typically noise-corrupted. It is then necessary to compensate for the influence of the noise. Thus, the intrinsic noise-resistant features of ST are of great value in the concept of an ST extremum controller addressed in this study.

3. ADROC Algorithm

The preceding subsection has examined all the components of adaptive dual response optimizing controller (ADROC). The aim of the ADROC that follows is to use a straightforward idea from dual response applications and then match parameterized input-output data to the approximate representation given by (7). The matching process proceeds in a manner to construct a well-behaved, “adaptive” extremum objective function (via recursive estimation) that accurately represents the nonlinearities of the process locally about the current operating conditions. Subsequently, an optimization step based on the TR-based (or dual-response) approach is carried out so as to try to achieve the eventual goal of extremum control. For illustration, the block diagram of the ADROC in an R2R control simulation scheme is demonstrated in figure 3.

For the clarity of presentation, the flow of the ADROC algorithm is summarized as follows:

- (i) Provide the quadratic model in (7) with a set of initial parameter estimates from an off-line design of experiments; that is, $\hat{\beta}_{p,0}$, $\hat{\beta}'_{p,0}$, $\hat{\beta}_{s,0}$, $\hat{\beta}'_{s,0}$, $\hat{\mathbf{B}}_{p,0}$, and $\hat{\mathbf{B}}_{s,0}$, and compute the initial bounds on the Lagrange multiplier μ via computation of three TR subproblems. Start the time index with $t = 0$.
- (ii) If the DRS is found degenerate, then use AXIS algorithm to decompose and solve every DRS subproblems with lower dimension. Provide the recipe $\boldsymbol{\mu}_t$ to the R2R system and go to step (iv). Otherwise, for the nondegenerate case go directly to step (iii).
- (iii) Compute the recipe $\boldsymbol{\mu}_t$ via the DRSALG algorithm.
- (iv) Set $t := t + 1$. Evaluate the “actual” performance (or objective) function values $y_p(\boldsymbol{\mu}_t(\mu))$ and $y_s(\boldsymbol{\mu}_t(\mu))$ from the true process or plant.
- (v) Calculate the updated parameter estimates ($\hat{\beta}_{p,t}$, $\hat{\beta}'_{p,t}$, $\hat{\beta}_{s,t}$, $\hat{\beta}'_{s,t}$, $\hat{\mathbf{B}}_{p,t}$, and $\hat{\mathbf{B}}_{s,t}$).

- $\hat{\beta}_{p,t}$, $\hat{\beta}_{s,t}$, $\hat{\mathbf{B}}_{p,t}$, and $\hat{\mathbf{B}}_{s,t}$) through the estimation algorithm.
- (vi) Form an updated, predicted response function $\hat{y}_p(\mu_{t+1}(\mu))$ and $\hat{y}_s(\mu_{t+1}(\mu))$ using the results obtained from step (v).
- (vii) Return to step (ii) and then compute the incumbent recipe μ_t via the optimization algorithm for the next control step at $t+1$.

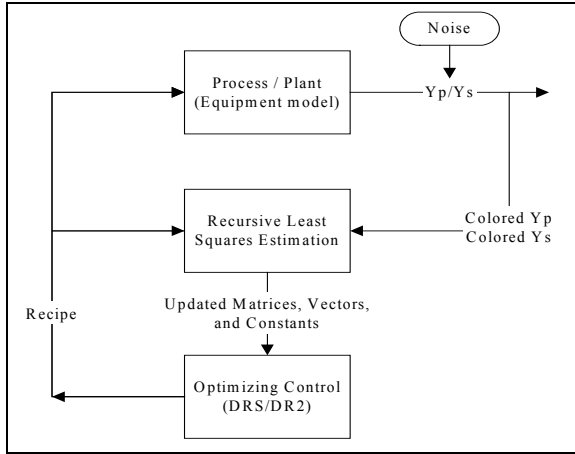


Figure 3. Block Diagram of the ADROC

4. EXPERIMENTAL STUDY OF ADROC BASED ON REAL EQUIPMENT MODELS

In order to investigate the performance of the ADROC, a CMP process was adopted. The chemical mechanical planarization (CMP) process is a typical R2R control situation in semiconductor manufacturing. Most CMP applications consider two critical, possibly “conflicting” quality characteristics; one of them is removal rate (RR) of silicon oxide and the other is within-wafer non-uniformity (WIWNU) on the wafer surface after polishing. CMP process in practice has been expected to reach around a specific removal rate and also to accomplish as minimum WIWNU as possible, leading to the multiple-input dual-output (MIDO) optimization case. Due to this particular model configuration, the new optimizing controller, termed ADROC, integrating the dual response systems with on-line estimations technique can be a perfect fit to these MIDO R2R control scenarios.

4.1 CMP Process Model

The simulated CMP process is based on the equipment models given by del Castillo and Yeh (1998). Control variables (scaled in the $(-1, 1)$ coding convention) consist of platen speed (μ_1), back pressure (μ_2), polishing downforce (μ_3), and the profile of the con-

ditioning system (μ_4). Each factor is constrained to inside the $(-1, 1)$ range. The two responses of interest are within-wafer nonuniformity (WIWNU) and removal rate (RR). The primary performance function of interest is WIWNU, denoted by y_p , and the secondary performance function of interest is RR, denote by y_s . In the simulation of the 4×2 CMP process, a target value of 1730 was set for y_s and y_p was to be minimized. The time variable t indicates the number of silicon wafers that have already been polished by the polish pad currently mounted. In this zoom CMP practice, the process constraints were considered adequate to be $y_s > 1700$, $y_p < 200$, and $-1 \leq \mu_i \leq 1$ for $i = 1, 2, \dots, 4$. The simulated equipment models (treated as the true production system) are given by

$$y_p = 254 + 32.6\mu_1 + 113.2\mu_2 + 32.6\mu_3 + 37.1\mu_4 - 36.8\mu_1\mu_2 + 57.3\mu_4t' - 2.42t' + \varepsilon_{1,t}, \quad (8)$$

$$y_s = 1563.5 + 159.3\mu_1 - 38.2\mu_2 + 178.9\mu_3 + 24.9\mu_4 - 67.2\mu_1\mu_2 - 46.2\mu_1^2 - 19.2\mu_2^2 - 28.9\mu_3^2 - 12\mu_4t' + 116\mu_4t' - 50.4t' + 20.4t'^2 + \varepsilon_{2,t}, \quad (9)$$

where $t' = (t - 53)/52 \in [-1, +1]$, implying that the age of the polishing pad range from 1 to 105, and afterwards a new pad is switched on. It needs to reset the scaled time variable t' once the first 105 wafers polishing would have been done. The random error term $\varepsilon_{1,t}$ for the primary response is assumed a white series to obey $N(0, 30^2)$, and the random error term $\varepsilon_{2,t}$ for the secondary response follows $N(0, 60^2)$, both estimated from the mean squared error (MSE) of the raw data. The deterministic drift is $-2.42t'$ for y_p and $-50.4t'$ for y_s . From (8-9), we note that the models are considerably more complex than quadratic approximations since the t' variable appears in quadratic and 2-factor interaction terms. Equations (8-9) are simulated under 5 different scenarios and the overall simulation results are tabulated in Tables 1-5, respectively.

4.2 Simulated 4×2 CMP Process

To begin with experimental studies, 50 independent simulations of 200 wafers (batch) each were run for collecting the computational results. The performance measures utilized herein contain the averages and standard deviations of the open-loop and closed-loop performance function values (\bar{y}_{open} , \bar{y}_{closed} , S_{open} , S_{closed}), and the average and standard deviation of the i^{th} controllable variable

Table 1. Simulation Results of 4×2 CMP Process: Scenario 1

Scenario 1	\bar{y}_p	S_{y_p}	\bar{y}_s	S_{y_s}	$\bar{\mu}_1$	$\bar{\mu}_2$	$\bar{\mu}_3$	$\bar{\mu}_4$
ADROC	190.8161 (5.4463)	36.9546 (3.7859)	1714.0281 (3.8149)	62.6662 (3.4609)	0.2988 (0.0267)	-0.8101 (0.0276)	0.4112 (0.0208)	-0.1546 (0.0491)
Fix-recipe-loop	178.3741 (1.8245)	31.3461 (1.4366)	1701.1761 (4.2085)	74.3084 (3.8766)	0.2600 (0)	-0.8650 (0)	0.3803 (0)	-0.1983 (0)
Open-loop	254.4501 (1.8245)	30.338270 (1.3884)	1570.0298 (4.2085)	66.413206 (3.7283)	—	—	—	—

($\bar{\mu}_i, S_{\mu_i}$). Besides, the standard errors of all the evaluation statistics were also assessed. In addition, the RLS method with the constant trace algorithm is applied with a fixed discounting factor $\lambda = 1.0$ and $P_0 = \mathbf{I}$ in the estimation block (see figure 3). Note that the CMP models are simulated under five different scenarios: (1) approximate initial models given and quadratic models selected, (2) Over-estimation in parameters of initial models, (3) Under-estimation in parameters of initial models, (4) under the presence of variance increase, (5) under in the presence of IMA(1, 1) disturbance.

Scenario 1: Approximate initial models given

To provide initial models, y_s was modeled by a full quadratic polynomial plus drift and y_p by a model with linear terms and 2-factor interactions plus drift. The initial models (see, e.g., del Castillo and Yeh 1998) were specified as follows:

$$y_p = 250 + 30\mu_1 + 100\mu_2 + 20\mu_3 + 35\mu_4 - 30\mu_1\mu_2 + 0.05t, \quad (10)$$

$$y_s = 1600 + 150\mu_1 - 40\mu_2 + 180\mu_3 + 25\mu_4 - 60\mu_1\mu_2 - 30\mu_1^2 - 20\mu_2^2 - 25\mu_3^2 - 0.9t, \quad (11)$$

which illustrate the case where reasonable initial models are provided to start the ADROC. The equality constraints on the outputs together with strict target values are useful in this example as removal rate is usually a “closer-the-better” response and nonuniformity a “smaller-the-better” response.

For illustration, a simulation of 200 wafers is pictorially demonstrated in figure 4, including both the closed-loop and open-loop control outputs. In lower part of figure 4, it was found that the closed-loop control outputs by the ADROC satisfy both process output constraints ($y_s > 1700, y_p < 200$ as before) but the open-loop control outputs cannot meet these requirements (see also table 1). The upper part of figure 4 presented all four controllable variables of the closed-loop in this scenario. Obviously, the ADROC is able to keep all control variables inside the constraint $-1 \leq \mu_i \leq 1$.

In table 1, there show two groups of control outputs:

open-loop and closed-loop (ADROC). For every group, the averages of $\bar{y}_p, \bar{y}_s, S_{y_p}$ and S_{y_s} were computed based on the 50 independent runs. Note that the standard errors of each evaluation statistic are indicated in parenthesis. From table 1 and figure 4, simulation results of the closed-loop on both responses are better than those of the open-loop. For the ADROC, though the target (1730) of removal rate is not attainable, however, the recipes generated by ADROC are still able to produce much smaller control outputs on nonuniformity (than the open-loop strategy) and keep removal rate as close as possible to the target (see figure 4).

The control outputs of the fix-recipe-loop were obtained by only using single fixed-recipe computed from the initial model. Notice that the control outputs of the fixed-recipe-loop are still within acceptable ranges even though the average of standard deviation of the secondary response (74.3084) is higher than that of the closed-loop (62.6662). Nonetheless, the fix-recipe-loop control strategy depends strongly on the accuracy of the initial model built upon an earlier off-line experiment. We will discuss this issue in next two scenarios.

Scenario 2: Over-estimation parameters of initial models

For the moment, we allow 0% ~ 20% over estimation error in each parameter of the initial model; namely, the magnitude of every parameter in this scenario will be increased by 0% ~ 20% units (which was sampled from a uniform random number).

It was particularly pointed out from Table 2 that ADROC can still produce satisfactory control outputs even if the bias error of initial models is present, but the control outputs of the fixed-recipe-loop, especially on removal rate (where the average of \bar{y}_s is 1485.8319 and the stand error of \bar{y}_s is 76.3905), were deteriorated resulting from a poor initial model. Note that even though the fixed-recipe-loop can achieve a lower level of the average of \bar{y}_p (122.9123), however, the standard errors of \bar{y}_p and \bar{y}_s are the way much greater than these by ADROC (5.9482 and 7.0148, respectively). It clearly indicates the robustness and control stability delivered by ADROC for this scenario. With respect to the open-loop outcomes, both average of \bar{y}_p and \bar{y}_s are not even acceptable to the process requirements.

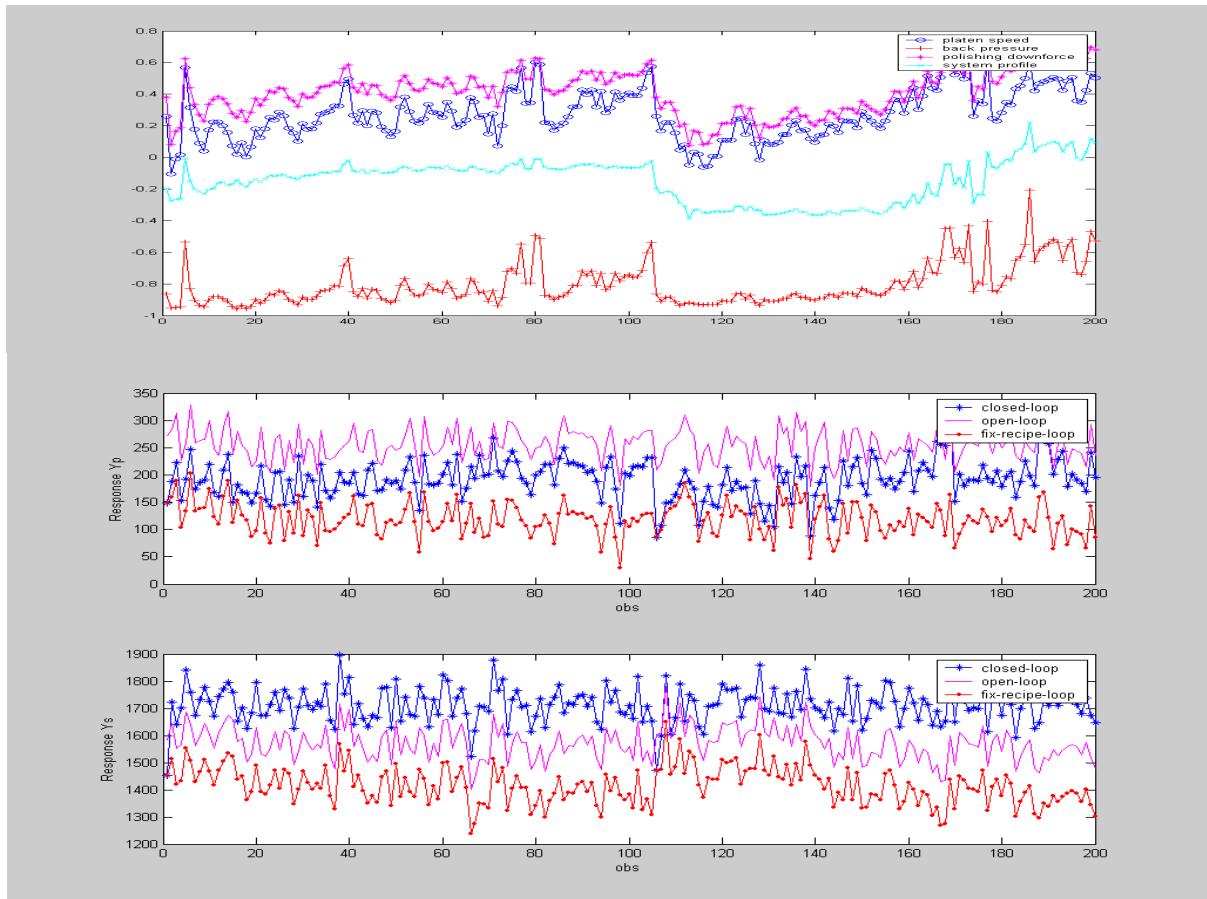


Figure 4. Single 200-Wafer Realization of 4×2 CMP Process under Scenario 1

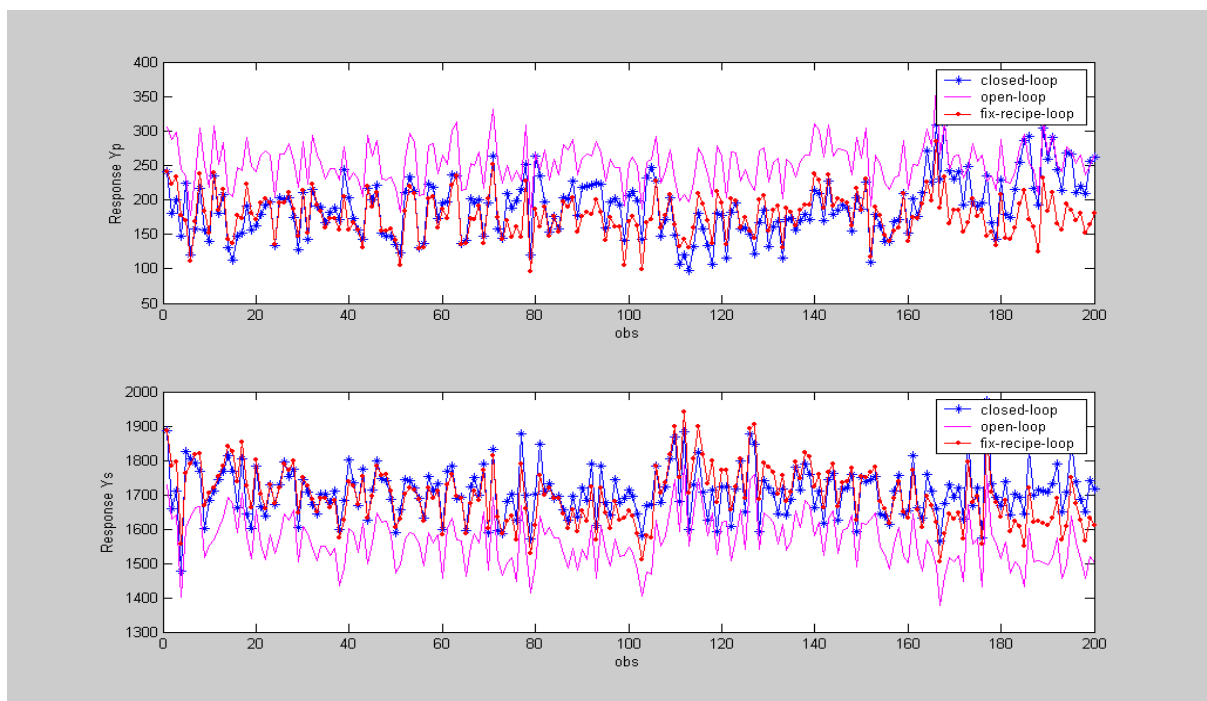


Figure 5. Single 200-Wafer Realization of 4×2 CMP Process under Scenario 2

Table 2. Simulation Results of 4×2 CMP Process: Scenario 2

Scenario 2	\bar{y}_p	S_{y_p}	\bar{y}_s	S_{y_s}	$\bar{\mu}_1$	$\bar{\mu}_2$	$\bar{\mu}_3$	$\bar{\mu}_4$
ADROC	187.7606 (5.9482)	41.2711 (5.2849)	1697.6627 (7.0148)	64.4344 (3.8337)	0.2631 (0.1906)	-0.7993 (0.1440)	0.3436 (0.1365)	-0.1532 (0.1546)
Fix-recipe-loop	122.9123 (12.6249)	31.4092 (1.0880)	1485.8319 (76.3905)	73.8516 (3.8255)	-0.1841 (0.1852)	-0.9356 (0.1523)	-0.0154 (0.1127)	-0.3006 (0.1733)
Open-loop	254.4537 (1.5328)	29.8592 (1.1466)	1572.5270 (3.7557)	65.9551 (3.9166)	—	—	—	—

Table 3. Simulation Results of 4×2 CMP Process: Scenario 3

Scenario 3	\bar{y}_p	S_{y_p}	\bar{y}_s	S_{y_s}	$\bar{\mu}_1$	$\bar{\mu}_2$	$\bar{\mu}_3$	$\bar{\mu}_4$
ADROC	192.7480 (5.9555)	36.4703 (3.2342)	1714.4090 (5.5249)	67.1745 (4.6031)	0.2883 (0.1474)	-0.7864 (0.1217)	0.4301 (0.1194)	-0.1409 (0.1293)
Fix-recipe-loop	263.0668 (24.8988)	30.3851 (1.4596)	1784.0845 (17.9392)	66.4410 (4.5661)	0.6417 (0.1098)	-0.3376 (0.1452)	0.6787 (0.1523)	0.1158 (0.1359)
Open-loop	254.4968 (1.8433)	30.3228 (1.5086)	1571.7618 (3.9597)	66.9126 (3.8906)	—	—	—	—

Table 4. Simulation Results of 4×2 CMP Process: Scenario 4

Scenario 4	\bar{y}_p	S_{y_p}	\bar{y}_s	S_{y_s}	$\bar{\mu}_1$	$\bar{\mu}_2$	$\bar{\mu}_3$	$\bar{\mu}_4$
ADROC	190.4475 (9.9386)	44.5329 (7.1375)	1710.1592 (5.0959)	95.1914 (6.3974)	0.2828 (0.1500)	-0.7956 (0.1361)	0.3976 (0.1273)	-0.1615 (0.1265)
Fix-recipe-loop	178.7326 (2.5281)	37.9778 (1.8635)	1704.5772 (6.6973)	101.2771 (6.6228)	0.2600 (0)	-0.8650 (0)	0.3803 (0)	-0.1983 (0)
Open-loop	254.3438 (2.5281)	37.288406 (1.8534)	1572.3625 (6.6973)	96.081984 (6.4204)	—	—	—	—

Scenario 3: Under-estimation parameters of initial models

In a similar way, we allow 0% ~ 20% under estimation error in each parameter of the initial model. The magnitude of every parameter in this scenario will be decreased by 0% ~ 20% units (which was also sampled from a uniform random number as before). It can be seen from Table 3, ADROC can still generate the primary response of the closed-loop under 200 level, and which satisfy process constraint $y_p < 200$. On the contrary, the average of WIWNU (y_p) of the fixed-recipe-loop is too high (263.0668) to be acceptable. Overall, among these three control strategies, ADROC still performed best from every statistic.

According to the discussions presented in scenarios 2 and 3, ADROC can generally provide an optimal, at least practically feasible, recipe from run to run and appropriate control outputs adaptable in changing conditions, even though the initial models suffered bias errors that the fix-recipe-loop control strategy can not cope with.

Scenario 4: Performance in the presence of variance increase

In this scenario, the same initial models (10-11) are provided to ADROC. A moderate variance increase of magnitude 40 is added to y_s at time 20 and a second variance increase of magnitude 10 is added to y_p at time 30. Figure 7 shows the control outputs of ADROC and closed-loop for a single 200-wafer realization. As indicated from lower part of Figure 7 (see also S_{y_s} in Table 4), there exhibits larger deviation (or fluctuation) of RR from target since time 20. Even so, ADROC is still able to control both process outputs comparably well.

Table 4 presents the comparison results of ADROC and open-loop computed over 50 independent simulations. In sum, the open-loop performance is absolutely impractical ($\bar{y}_p \approx 254$, $\bar{y}_s \approx 1572$), but, by contrast, the closed-loop strategy by ADROC yields acceptable control outputs on both responses as the process variance inflates. In sum, ADROC brings WIWNU further down to 190 while controlling the RR response very well ($\bar{y}_s \approx 1710$).

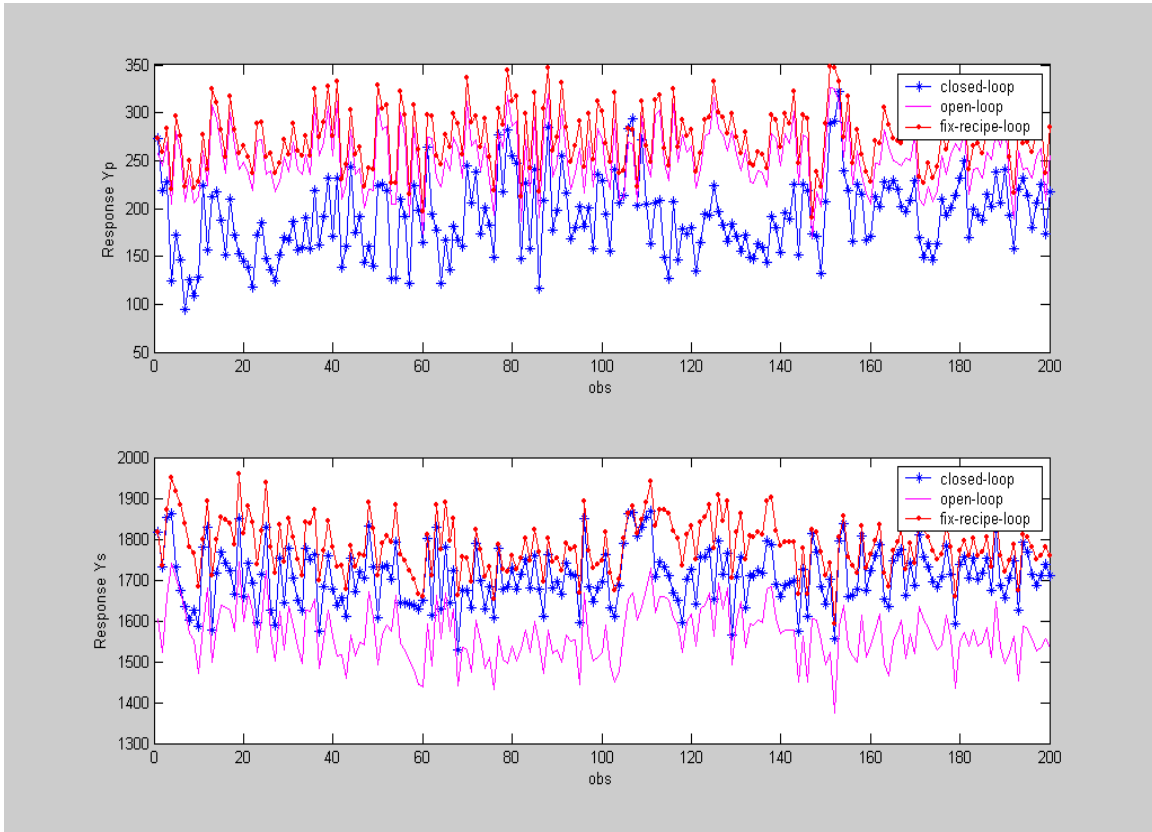


Figure 6. Single 200-Wafer Realization of 4×2 CMP Process under Scenario 3

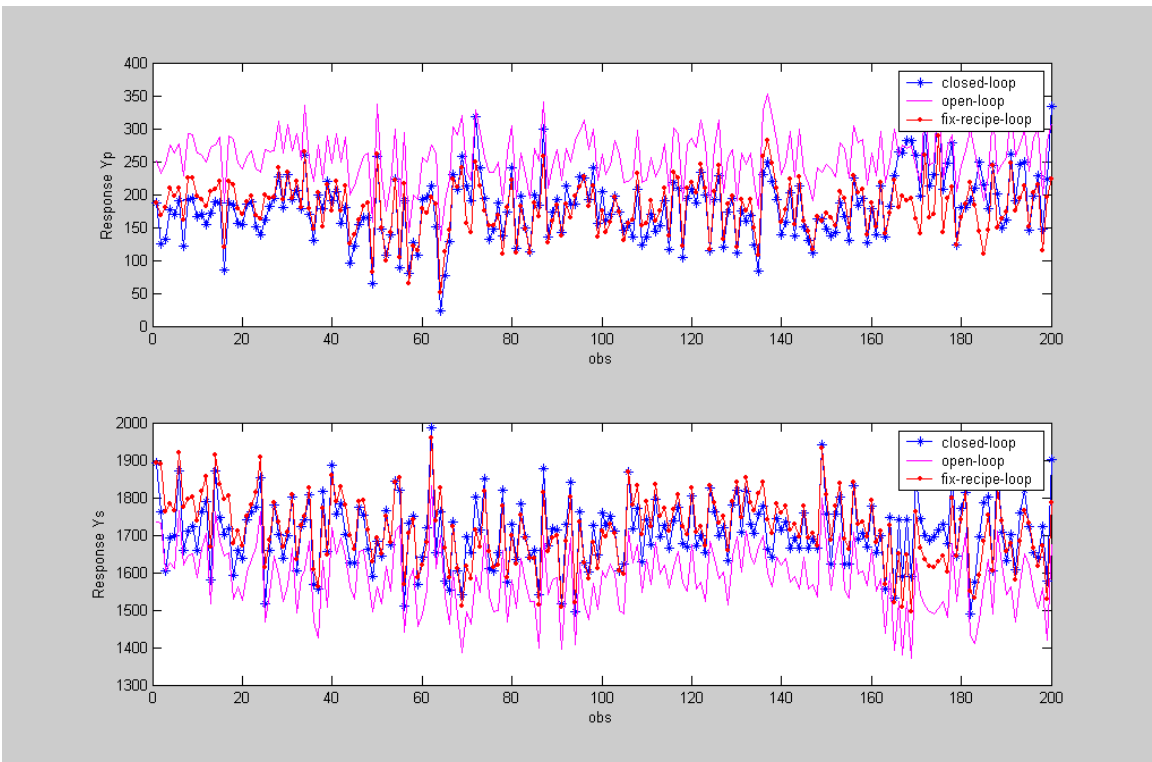


Figure 7. Single 200-Wafer Realization of 4×2 CMP Process under Scenario 4

Scenario 5: Performance in the presence of IMA disturbance

Again, initial models (10-11) are provided to ADROC, and an IMA(1, 1) disturbance is considered.

The IMA(1, 1) disturbance, denoted by $D(t)$, for two responses obeys the form of

$$D_p(t) = \eta_{p,t}(1-\theta B) + \eta_{p,t-1} \tag{12}$$

$$D_s(t) = \eta_{s,t}(1-\theta B) + \eta_{s,t-1} \tag{13}$$

where $\theta = 0.70$, $\eta_{p,t} \sim N(0, 1^2)$ and $\eta_{s,t} \sim N(0, 2^2)$ are assumed. In common R2R practice, an IMA series is usually used to model machine aging and tool wearing. With this CMP process configuration, Figure 8 shows the control outputs of open-loop and ADROC for a single 200-wafer simulation. Evidently, the closed-loop performance is the way better than open-loop on secondary process response. As can be seen from Figure 8, ADROC

generates extremely steady control outputs on RR response, meaning that the drifting process of IMA model can be compensated for effectively on secondary response. The ADROC results exhibit additional merits of having high stability on the recipe profile of $\bar{\mu}_1 \sim \bar{\mu}_4$, even if it returns the average mean response and standard deviation of WIWNU a little inferior to those of the fix-recipe-loop ($\bar{y}_p \approx 194$, $S_{y_p} \approx 27.7$).

5. Conclusions and Further Research

This paper provided an evaluation of a new multiple-input double-output controller for CMP process in semiconductor manufacturing. This controller, termed ADROC, is founded mainly on the theory of adaptive extremum control. In the ADROC, a global optimization algorithm

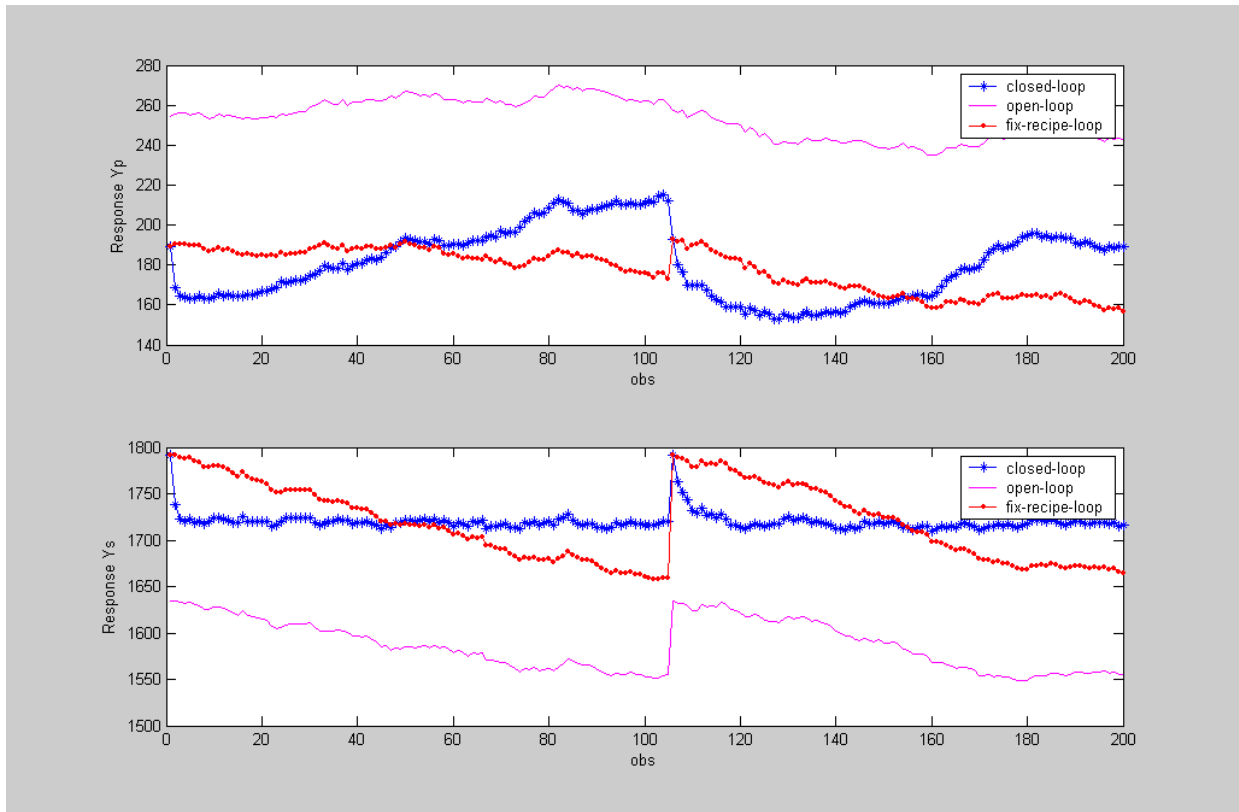


Figure 8. Single 200-Wafer Realization of 4×2 CMP Process under Scenario 5

Table 5. Simulation Results of 4×2 CMP Process: Scenario 5

Scenario 5	\bar{y}_p	S_{y_p}	\bar{y}_s	S_{y_s}	$\bar{\mu}_1$	$\bar{\mu}_2$	$\bar{\mu}_3$	$\bar{\mu}_4$
ADROC	194.7876 (9.3949)	27.7284 (9.3386)	1721.0368 (2.8759)	15.3202 (2.0253)	0.3012 (0.1437)	-0.7930 (0.1376)	0.4341 (0.1251)	-0.1098 (0.1169)
Fix-recipe-loop	177.3440 (4.8239)	10.3211 (3.5529)	1702.3081 (8.8523)	47.2430 (5.6744)	0.2600 (0)	-0.8650 (0)	0.3803 (0)	-0.1983 (0)
Open-loop	252.9552 (4.8239)	6.769904 (3.4626)	1570.0934 (8.8523)	33.646803 (5.7078)	—	—	—	—

based on theories of dual response systems was employed serving as an extremum-seeking controller and the RLS method with the constant trace algorithm was utilized acting as an estimation self-tuner in the closed-loop structure. Thus, the R2R controller can act both as a controller and as an optimizer. Treating equipment models built upon the real data as real production control systems, a CMP process was simulated to illustrate the performance of the proposed controller in typical experimental R2R circumstances. The findings of this research can be summarized as follows:

- A R2R controller based on a model of dual response systems was able to control nonlinear MIMO process through the approximation of the nonlinearity with quadratic polynomial in terms of input variables.
- ADROC could rapidly achieve the extremum of controlled outputs nearly without transients and thereafter kept consistent control results by using continuously updated quadratic models. It indicates the merit that the control knowledge addressed in this research is familiar to applied statisticians and quality engineers, as quadratic approximations are common in the area of response surface methodology and on-line quality control.
- Multivariate adaptive control has been shown to be a feasible control strategy in run-to-run problem due to its on-line estimation nature. This allows not only the control of some equipment based on given models developed previously, but also the optimization of the equipment as well.
- It was observed that the better the initial models provided to the ADROC algorithm, the better the control output and the less severe the transient effect is. But, the ADROC can allow the initial models parameters 20% deviation (estimation error) from the true process model such that ADROC still reached the extremum quickly by dual response systems and RLS algorithm.
- It was shown by example of how the ADROC could cope with system with two responses that are confounded with the correlated input variables. It reveals from the evaluation that ADROC can provide excellent control actions for the MIMO R2R situations even though the process exhibits complicated, nonlinear interaction effects between control variables, and the drifting disturbances.
- The DRS/DR2 algorithm in the ADROC provides the solution (process recipe) to system that minimizes primary response and keeps secondary response on the target. Even when this is not the case, the ADROC would try “the best it can” in the least squares sense, which is the base of the recursive estimation method utilized.

An interesting opportunity for future research would be the extension of ADROC to the multiple-input

multiple-output (MIMO) case. Further work can be devoted to studying the control performance of ADROC when there exhibit other process and/or noise dynamics. An area that also needs further investigation is that a “live” R2R project allowing us to exercise the controller on-line should take place to complement the research results shown here.

REFERENCES

- Åström, K. J. and Wittenmark, B. (1973), On Self-Tuning Regulators, *Automatica*, 9, 185-199.
- Åström, K. J. and Wittenmark, B. (1995), *Adaptive Control*, 2nd edition, Addison-Wesley: Reading, Mass.
- Box, G. E. P., Jenkins, G. M. and Reinsel, G. C. (1994), *Time Series Analysis: Forecasting and Control*, 3rd edition, Prentice Hall: Englewood Cliffs, NJ.
- Del Castillo, E. and Hurwitz, A. M. (1997), Run-to-Run Process Control: Literature Review and Extensions, *Journal Quality Technology*, 29, 184-496.
- Del Castillo, E. and Yeh, J.-Y. (1998), An Adaptive Run-to-Run Optimizing Controller for Linear and Nonlinear Semiconductor Processes, *IEEE Transactions on Semiconductor Manufacturing*, 11, 285-295.
- Del Castillo, E., Fan, S.-K. S. and Semple, J. (1999), Optimization of Dual Response Systems: A Comprehensive Procedure for Degenerate and Nondegenerate Problems, *European Journal Operation Research*, 112, 174-186.
- Drajer, N. R. (1963), Ridge Analysis of Response Surfaces, *Technometrics*, 5, 469-479.
- Fan, S.-K. S. (2000a), Quality Improvement of Chemical-mechanical Wafer Planarization Process in Semiconductor Manufacturing Using a Combined Generalized Linear Modeling-Nonlinear Programming Approach, *International Journal Production Research*, 38, 3011-3029.
- Fan, S.-K. S. (2000b), A Generalized Global Optimization Algorithm for Dual Response Systems, *Journal Quality Technology*, 32, 444-456.
- Golden, M. P. and Ydstie, B. E. (1989), Adaptive Extremum Control Using Approximate process Models, *AIChE Journal*, 35, 1157-1169.
- Golub, G. H. and van Loan, C. F. (1984), *Matrix Computations*, The John Hopkins University Press: Baltimore, MD.
- Khuri, A. I. and Cornell, J. A. (1996), *Response Surface: Designs and Analyses*, 2nd edition, Marcel Dekker: New York, NY.
- Luenberger, D. G. (1984), *Linear and Nonlinear Programming*, 2nd edition, Addison-Wesley: Reading, MA.
- Moré, J. J. and Sorensen, D. C. (1983), Computing a Trust Region Step, *SIAM Journal on Scientific and*

- Statistical Computing*, 4, 553-572.
- Myers and Montgomery (2002), *Response Surface Methodology: Process and Product Optimization Using Designed Experiments*, 2nd edition, Wiley, Inc, CA.
- Myers, R. H. and Carter, W. H. (1973), Response Surface Techniques for Dual Response System, *Technometrics*, 15, 301-317.
- Nanz, G. and Camilletti, L. E. (1995), Modeling of Chemical-mechanical Polishing: A Review, *IEEE Transactions on Semiconductor Manufacturing*, 8, 382-389.
- Patel, N. S. and Jenkins, S. T. (2000), Adaptive Optimization of Run-to-Run Controller: The EWMA Example, *IEEE Transactions on Semiconductor Manufacturing*, 13, 97-107.
- Sachs, E., Hu, A. and Ingolfsson, A. (1995), Run by Run Process Control: Combining SPC and Feedback Control, *IEEE Transactions on Semiconductor Manufacturing*, 8, 26-43.
- Seborg, D. E., Edgar, T. F. and Shah, S. L. (1986), Adaptive Control Strategies for Process Control: A Survey, *AIChE Journal*, 32, 881-913.
- Semple, J. (1997), Optimality Conditions and Solution Procedures for Nondegenerate Dual Response Systems, *IEE Transactions*, 29, 743-752.
- Shah, S. and Cluett, W. R. (1991), Recursive Least Squares Based Estimation Schemes for Self-Tuning Control, *The Canadian Journal of Chemical Engineering*, 69, 89-96.