# OPTIMAL CONTROL OF TRANSIENT ENHANCED DIFFUSION

**R. Gunawan, M. Y. L. Jung, E. G. Seebauer, and R. D. Braatz**

*Department of Chemical and Biomolecular Engineering*
*University of Illinois at Urbana-Champaign*
*Urbana, IL 61801*

Abstract: Transient enhanced diffusion of boron inhibits the formation of ultrashallow junctions needed in the next-generation of microelectronic devices. Reducing the junction depth using rapid thermal annealing with high heating rates comes at a cost of increasing sheet resistance. The focus of this study is to design the optimal annealing temperature program that gives the minimum junction depth while maintaining satisfactory sheet resistance. Comparison of different parameterizations of the optimal trajectories shows that linear profiles gave the best combination of minimizing junction depth and sheet resistance. Worst-case robustness analysis of the optimal control trajectory motivates improvements in feedback control implementations for these processes. *Copyright © 2003 IFAC*

Keywords: optimal control, model-based control, semiconductors

## 1. INTRODUCTION

Moore's law requires a continued shinkage of feature sizes in microelectronic devices. For example, advanced CMOS devices will require junction depths between 13 to 22 nm in the source and drain extension region by the year 2005 according to the 2001 International Technology Roadmap for Semiconductors. The current technology for the formation of such ultrashallow junctions depends on ion implantation of dopant, such as boron, into silicon. Although the junction depth can be made shallower by reducing the implant energy, the effectiveness of this approach is limited by the need to anneal out the point and/or extended defects generated by ion implantation. Silicon self-interstitial defects can mediate the diffusion of dopants during the annealing process, which leads to a significant increase of the junction depth. This phenomenon is known as "transient enhanced diffusion" (TED). For this reason, considerable efforts have been put forth in the modeling of the TED for designing appropriate post-implant annealing programs to produce the desired junction depth (see (Jain, *et al.*, 2002) and references therein).

The state-of-the-art in post-implant annealing employs a lamp-based rapid thermal annealing (RTA). Figure 1 shows a typical RTA "spike" anneal program, which consists of a stabilization step at constant temperature (~650 °C), followed by a linear heating step at a constant rate (~100 °C/s) reaching a maximum temperature (~1000 °C), and finally a radiative cooling step at a initial rate of several tens of degrees per second. In the literature, there exists conflicting experimental evidence on the efficacy of using high heating rates (up to 400 °C/s) in the spike anneal profile to reduce TED (Downey, *et al.*, 1999; Shishiguchi, *et al.*, 1997). Recent results (Gelpey, *et al.*, 2002; Mannino, *et al.*, 2001) tend to confirm the benefit of using high heating rates. The results also suggest that the reduction in the junction depth comes at the expense of an undesired increase in the sheet resistance. The tradeoff in reducing the junction depth without sacrificing the sheet resistance motivates a careful optimization of the post-implant annealing temperature program.

A recently developed comprehensive TED model consists of a set of reaction-diffusion equations combined with Poisson's equation to account for the electric field effects on charged species (Jung, *et al.*,
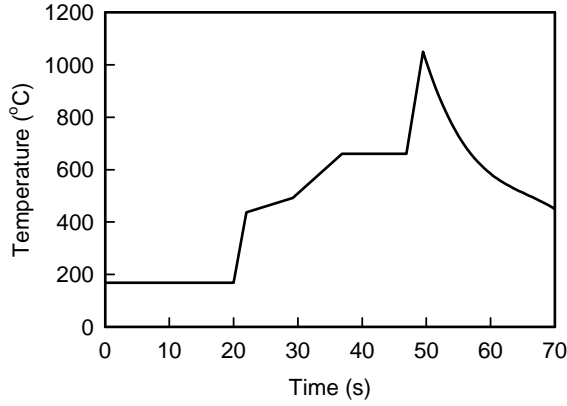
Fig. 1. A typical RTA temperature program.

1999). The activation energies arising from the reaction rate constants and diffusivities in the TED model are not exactly known, but there exist extensive experimental and computational estimates of these parameters. To resolve problems with regard to conflicting estimates in the literature, maximum likelihood (ML) estimation was applied to give the most likely values and the standard deviations for the parameters from the published parameter estimates. Furthermore, maximum *a posteriori* (MAP) estimation was applied to produce improved parameters from the ML estimates and experimental Boron profile data collected at International Sematech (Gunawan, *et al.*, 2003).

This paper focuses on the design of the spike anneal program that optimizes the junction depth subject to a constraint on the sheet resistance. The TED model is implemented using the process simulator FLOOPS (Law and Tasch, 2000). Different parameterizations of the optimal trajectory are used to elucidate the true optimal annealing program. Worst-case analysis of the resulting optimal trajectory quantifies the performance degradation with respect to control implementation inaccuracies and model uncertainties.

## 2. TRANSIENT ENHANCED DIFFUSION MODEL

Transient enhanced diffusion arises from reaction-diffusion processes consisting of Fickian diffusion, electrical drift motion, and reaction networks including boron activation and interstitial clustering. The model comprises of coupled continuity equations (i.e., mass balances) for each species and Poisson's equation to include the electrical field effect on the charged species. The general continuity equation is

$$\frac{\partial N_i}{\partial t} = -\frac{\partial J_i}{\partial x} + G_i \qquad (1)$$

where $N_i$ denotes the concentration, $J_i$ is the flux, and $G_i$ is the net generation rate of species $i$. The flux $J_i$ includes terms from the Fickian diffusion and the electric field drift motion:

$$J_i = -D_i \frac{\partial^2 N_i}{\partial x^2} + \gamma_i \mu_i N_i E(x), \qquad (2)$$

where $D_i$ denotes the diffusivity and $E(x)$ is the electric field. The mobility $\mu_i$ follows the Einstein relation

$$\mu_i = \frac{qD_i}{kT}, \qquad (3)$$

where $q$ is the electron charge, $k$ is the Boltzmann constant, and $T$ is the temperature. The term $\gamma_i$ describes the average charge of species according to

$$\gamma_i = \sum_j z_j \gamma_{z_j}, \qquad (4)$$

where $z_j$ are the possible charge states (*i.e.,* +2, 0, $-1$, etc.) and $\gamma_{z_j}$ is the fraction of species $i$ having charge $z_j$ according to the Fermi-Dirac statistics.

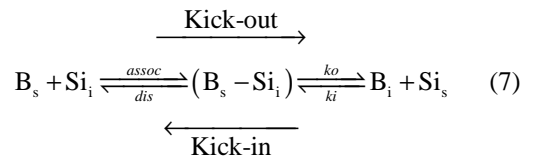Poisson's equation describes the electric field induced by the spatial charge imbalance:

$$\frac{\partial^2 \psi}{\partial x^2} = \frac{Q(x)}{\varepsilon} \qquad (5)$$

where $\varepsilon$ denotes the dielectric constant and the charge density $Q(x)$ is given by

$$Q(x) = p - n + \sum_i \gamma_i N_i \qquad (6)$$

with $p$ and $n$ denoting the hole and electron concentrations, respectively. The concentrations $p$ and $n$ are assumed to be in thermal equilibrium.
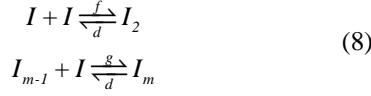
The generation term $G_i$ includes the formation and annihilation rates due to the boron activation reaction and/or clusters formation and dissociation. The boron activation reaction provides a pathway between mobile interstitial boron $B_i$ to and from immobile activated boron (*i.e.,* substitutional boron, $B_s$):

$$B_s + Si_i \underset{dis}{\overset{assoc}{\rightleftharpoons}} (B_s - Si_i) \underset{ki}{\overset{ko}{\rightleftharpoons}} B_i + Si_s \qquad (7)$$

In addition, the intermediate $(B_s - Si_i)$ acts as nucleation centers for mixed boron-silicon clusters. The rates of reactions follow reactant-limited rate expressions with the reaction rate constants adhering to the Arrhenius law.
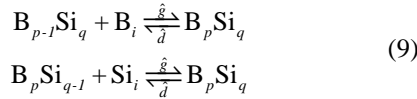
Clusters of interstitial atoms have been shown to form during TED (Collart, *et al.*, 2000; Stolk, *et al.*,

1997). There is evidence supporting the formation of clusters consisting of pure boron (Collart, *et al.*, 2000), pure silicon (Stolk, *et al.*, 1997), and mixed boron-silicon (Haynes, *et al.*, 1996). During thermal annealing, the clusters can act as reservoir (during the stabilization step) and source (during ramp up) for mobile interstitial boron and silicon. The formation and dissolution of pure interstitial clusters follow the reactions

$$I + I \underset{d}{\overset{f}{\rightleftharpoons}} I_2$$
$$I_{m-1} + I \underset{d}{\overset{g}{\rightleftharpoons}} I_m$$

(8)

where $I$ denotes the interstitials (boron and silicon) and the indices $m$ denote the sizes of the clusters. The cluster formation rate assumes a reactant diffusion-limited reaction in agreement with much of the literature (see for example (Laidler, 1987)). On the other hand, the dissolution rate follows a first-order kinetic expression with rate constant according to the Arrhenius law.

The formation and dissolution of mixed boron-silicon clusters is described by:

$$B_{p-1}Si_q + B_i \underset{\hat{d}}{\overset{\hat{g}}{\rightleftharpoons}} B_pSi_q$$
$$B_pSi_{q-1} + Si_i \underset{\hat{d}}{\overset{\hat{g}}{\rightleftharpoons}} B_pSi_q$$

(9)

where $p$, $q$ are integers larger than or equal to 1. The formation and dissolution rates of mixed clusters again follow diffusion-limited and first-order kinetics, respectively, as in the pure cluster dynamics.

The TED model requires a set of activation energies associated with the diffusivities and kinetic rate constants for the boron activation reaction and cluster dissociation dynamics. These activation energies are difficult to directly measure experimentally and determine computationally. Experimental and *ab initio* density functional theorem (DFT) estimates of the activation energies are scattered throughout the literature. For many of the activation energies, the published values show significant variation. To resolve problems in regard to conflicting estimates in the literature, maximum likelihood (ML) estimation was applied to determine the most likely values and the standard deviations from the published parameter estimates (Gunawan, *et al.*, 2003).

Maximum *a posteriori* estimation takes a Bayesian approach which combines experimental data with the *a priori* information, in this case, from maximum likelihood estimation of published experimental and/or DFT values (Gunawan, *et al.*, 2003). Figure 2 presents the after-anneal experimental data used in the MAP estimation along with simulation profiles using the MAP parameters employing various RTA programs. Figure 3 shows the agreement between the TED model using the MAP estimates and the
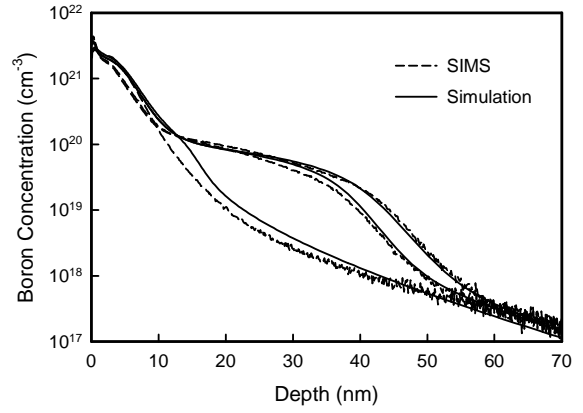


Fig. 2. Comparison of experimental and simulation boron profiles using the TED model with the MAP parameters. The junction depth is defined as the spatial penetration of boron at a total concentration of $10^{18}$ atoms/cm$^3$

experimental observations compiled from the literature (Agarwal, *et al.*, 1999).

## 4. OPTIMAL CONTROL FORMULATION

In the literature, control of transient enhanced diffusion through manipulation of RTA programs adopted an *ad hoc* approach through trial and error (Jain, *et al.*, 2002), due to incomplete understanding of TED mechanisms and correspondingly undependable models for describing dopant diffusion and activation. In contrast, this work employs a model-based control approach for designing an optimal temperature program that minimizes the junction depth while maintaining a suitable sheet resistance. The optimization variable is the RTA temperature profile, in particular, the heating and cooling profiles and the annealing temperature. A conventional RTA only employs a radiative cooling step, but there exists evidence (Agarwal, 2000;
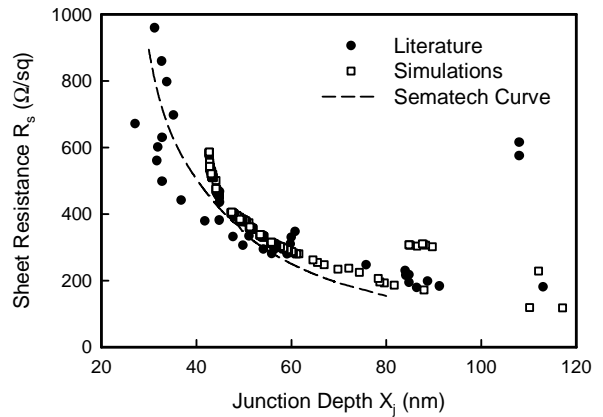


Fig. 3. Comparison of junction depth and sheet resistance data from experimental works employing various annealing programs as summarized in the Sematech curve, and from the TED simulations.

Agarwal, *et al.*, 1999) supporting the importance of the ramp down trajectory, especially in high heating rate applications (>150 °C/s).

The optimal control objective chosen here is to minimize the junction depth while maintaining a satisfactory sheet resistance, which is equivalent to the minimization problem:

$$\min_{\substack{T(t) \\ R_s \le R_{s,\max} \\ \beta_{\min} \le \frac{dT}{dt} \le \beta_{\max} \\ T \le T_{\max}}} X_j \qquad (10)$$

where $X_j$ denotes the junction depth, $R_s$ denotes the sheet resistance, and $T(t)$ is the RTA temperature trajectory. The sheet resistance is given by:

$$R = \frac{1}{q \int \mu(C) C(x) dx} \left[ \frac{\Omega}{\text{sq}} \right] \qquad (11)$$

where $q$ denotes the carrier charge, $\mu(C)$ denotes the mobility (concentration dependent), and $C(x)$ is the spatial concentration of charge carrier (*i.e.,* activated dopant). The following empirical formula gives the mobility $\mu(C)$ for boron (Zeghbroeck, 2002):

$$\mu(C) = 44.9 + \frac{425.6}{1 + \left( \frac{C}{2.23 \times 10^{17}} \right)^{0.719}} \left[ \frac{\text{cm}^2}{\text{Vs}} \right] \qquad (12)$$

In this work, it is desired to produce junctions with the sheet resistance below $R_{s,\max}$ of 350 $\Omega$/sq. The constraints on the temperature gradient, *i.e.,* $\beta_{\min}$ and $\beta_{\max}$, describe the limits for the cooling and heating rates, respectively. The state-of-the-art lamp-based RTA can produce heating rates up to 400 °C/s (Shishiguchi, *et al.*, 1997), while recent advances in RTA technology can achieve cooling rates up to 200 °C/s (Vortek Industries Ltd., 2002). The maximum temperature of thermal anneal $T_{\max}$ is set to the melting point temperature of silicon at 1410 °C.

## 5. WORST CASE ANALYSIS

Worst case analysis (Ma and Braatz, 2001) provides tools for quantifying the robustness of the optimal control performance to uncertainties in model parameters and control implementation. The information can be used to assess whether a more accurate model and thus more experiments are needed, or to give the desired performance and accuracy of the lower level control loops and control equipment, respectively. The parametric and control uncertainties are described as norm bounded perturbations $\delta u$ and $\delta \theta$, that is,

$$E_\theta = \left\{ \theta : \theta = \hat{\theta} + \delta \theta, \ \| W_\theta \delta \theta \| \le 1 \right\} \qquad (13)$$

$$E_u = \left\{ u : u = \hat{u} + \delta u, \ \| W_u \delta u \| \le 1 \right\} \qquad (14)$$

where $W_\theta$ and $W_u$ are positive-definite weighting matrices. This formulation includes uncertain parameters lying within a hyperellipsoid as well as independent bounds on each element.

For brevity, only the simplest techniques for the worst case analysis of batch processes is summarized here. A first-order expansion of the performance objective with respect to the model parameters gives

$$\delta \Phi = L \left( \theta - \hat{\theta} \right) = L \delta \theta \qquad (15)$$

where $L$ denotes the sensitivity coefficients given by

$$L_i = \left. \frac{\partial \Phi}{\partial \theta_i} \right|_{\theta = \hat{\theta}} \qquad (16)$$

Based on this expansion, the worst-case deviation of the performance is defined by (Ma and Braatz, 2001)

$$\delta \Phi_{wc} = \max_{\| W_\theta \delta \theta \| \le 1} \left| L \delta \theta \right| \qquad (17)$$

Similar worst case analysis with respect to the control implementation inaccuracies requires a second-order series expansion:

$$\delta \Phi = M \delta u + \delta u^T H \delta u \qquad (18)$$

where

$$M_j = \left. \frac{\partial \Phi}{\partial u_j} \right|_{u = \hat{u}} \qquad (19)$$

$$H_{ij} = \left. \frac{\partial^2 \Phi}{\partial u_i \partial u_j} \right|_{u = \hat{u}} \qquad (20)$$

Then the worst-case performance deviation due to control errors is:

$$\delta \Phi_{wc} = \max_{\delta u_{\min} \le \delta u \le \delta u_{\max}} \left| M \delta u + \delta u^T H \delta u \right| \qquad (21)$$

This optimization problem is equivalent to

$$\max_{\mu_\Delta(N) \ge k} k \qquad (22)$$

where $k$ is any real number, the perturbation $\Delta = \text{diag} \{ \Delta_r, \Delta_r, \delta_c \}$ consists of independent real scalar blocks $\Delta_r$ and a complex scalar $\delta_c$, and

$$N = \begin{bmatrix} 0 & 0 & kw \\ kH & 0 & kHz \\ z^T H + M & w^T & z^T Hz + Mz \end{bmatrix} \qquad (23)$$

where

$$w = \tfrac{1}{2}\left(\delta u_{\max} - \delta u_{\min}\right) \qquad (24)$$

$$z = \tfrac{1}{2}\left(\delta u_{\max} + \delta u_{\min}\right) \qquad (25)$$

and $\delta u_{\max}$ and $\delta u_{\min}$ are the upper and lower bounds for the control implementation inaccuracies. Tight upper and lower bounds for $k$ can be computed in polynomial-time using iterative $\mu$-calculations or skewed-$\mu$ analysis.

## 6. RESULTS AND DISCUSSION

The wafers were implanted with $2\times10^{15}$ ions/cm² of boron at 0.60 keV with 0° tilt, which gave a junction depth of 40 nm. The total boron was assumed to be initially 20% substitutional boron and 80% interstitial boron (Kobayashi, *et al.*, 2001). The initial conditions for Si interstitials agreed with the "+1" model, where Si interstitial concentration tracked the total boron concentration. The clusters and the $B_s$-$Si_i$ complex were assumed not present initially. Boundary conditions at the surface for all species assumed no flux (*i.e.*, $J_i|_{\text{surface}} = 0$) with no surface Fermi level pinning (Jung, *et al.*, 2001). The optimization is solved by extending the golden search method (Press, *et al.*, 1992) to multidimensional problems.

Figure 4 presents the optimal RTA programs using linear and quadratic parameterizations of the temperature trajectory, which give junction depths of 51.3 and 48.2 nm, respectively (see Fig. 5), and the same sheet resistance of 350 $\Omega$/sq. The optimal linear heating and cooling rates were 400 °C/s and 200 °C/s, respectively, indicating that the optimal RTA program was to effectively heat and cool as quickly as possible to the annealing temperature of 1111 °C, in agreement with experimental studies (Agarwal, 2000; Mannino, *et al.*, 2001). The use of a high annealing temperature with fast heating and cooling can be explained by the lower effective activation energy of TED compared to boron activation.
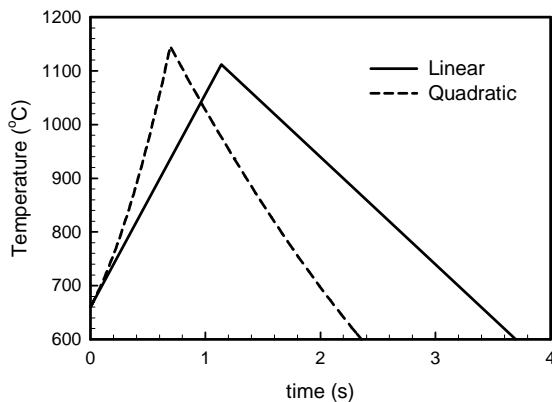


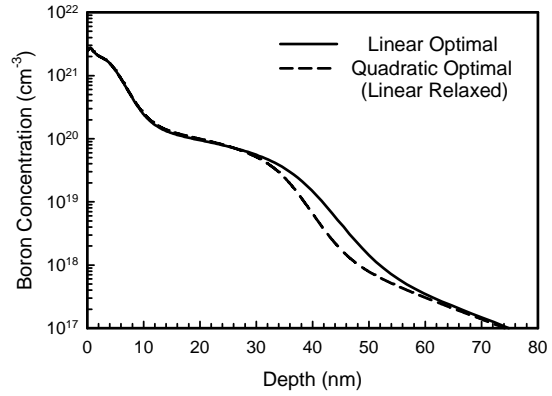Fig. 4. Optimal RTA programs employing linear and quadratic parameterizations.



Fig. 5. Simulations of after-anneal Boron profiles employing the optimal RTA programs.

Although experimental studies (Agarwal, 2000; Mannino, *et al.*, 2001) had alluded to using fast heating and cooling rates, the determination of the maximum annealing temperature was made through extensive trial and error. In contrast, the optimal control formulation using the TED model can directly determine the annealing temperature, and therefore reduce the number of costly experiments.

The quadratic parameterization was applied to the heating ramp, while the cooling rate was kept at the optimal linear case. The slope of the quadratic profile was limited to 1000 °C/s, which gave the optimal trajectory with an annealing temperature of 1144 °C. The quadratic heating profile only gave minimal improvement of the junction depth over a linear heating profile. If the optimization problem for linear heating profile was solved using a relaxed constraint on the heating rate $\beta_{\max}$ of 1000 °C/s, the optimal annealing temperature increased to 1146 °C giving a junction depth of 48.2 nm (see Fig. 5). In other words, if the same maximum heating rate is used, the quadratic and linear parameterizations give the same minimum junction depth. Since heating rate is the true constraint in practice, this indicates that there is no benefit to using quadratic over linear heating and cooling profiles.

Worst case analysis was applied to the linear optimal trajectory. The model parameter uncertainties were quantified by the MAP covariance estimate (Gunawan, *et al.*, 2003). The analysis on control implementation inaccuracies used control trajectory perturbations of 5 °C, 10 °C, and 15 °C at five temperatures along the heating and cooling ramps (660 °C, 800 °C, 950 °C, 1050 °C, 1100 °C, 1112 °C). Table 1 presents the worst-case junction depth increases due to uncertainties in the model parameters and control implementation.

Table 1. Worst-case junction depth increases (in nm) from parameter uncertainty and control errors.

| $\delta\theta$ | $|\delta u| \leq 5$ °C | $|\delta u| \leq 10$ °C | $|\delta u| \leq 15$ °C |
|---|---|---|---|
| 0.13 | 1.89 | 5.15 | 9.78 |

The analysis results indicate that the deviations from the optimal junction depth were minimal for model parameter uncertainties and moderate to significant for control inaccuries. These results indicate that the MAP estimation gave parameter estimates with sufficient accuracy for use in optimal control studies. The typical RTA controllers make ~20 control moves every second (Bratschun, 1999), which translates to every 20 ºC in the heating step. Further accounting of nonuniformity of temperature across the wafer, the control inaccuracies could exceed 15 ºC at any given time. The analysis indicates that existing feedback controllers for implementing RTA programs need improvement as future junction depth requirements necessitate further reduction of the junction depth.

## 7. CONCLUSIONS

This paper has shown that the optimal RTA program for minimizing TED while achieving the desired sheet resistance consisted of fast linear heating and cooling profiles, as suggested in many experimental studies. Worst case analysis on the optimal junction depth deviations suggested the need of improvements in existing RTA controllers and advances in RTA technology to ensure temperature uniformity across the wafer.

## REFERENCES

Agarwal, A. (2000). *Ultra-shallow junction formation using conventional ion implantation and rapid thermal annealing.* Paper presented at the *Int. Conf. on Ion Impl. Tech.*, Austria.

Agarwal, A., Gossmann, H.-J., & Fiory, A. T. (1999). Effect of ramp rates during rapid thermal annealing of ion implanted boron for formation of ultra-shallow junctions. *J. Elec. Mat.,* **28**(12), 1333-1339.

Bratschun, A. (1999). The application of rapid thermal processing tehcnology to the manufacture of integrated circuits - An overview. *J. Elec. Mat.,* **28**, 1328-1332.

Collart, E. J. H., Murrell, A. J., Foad, M. A., van den Berg, J. A., Zhang, S., Armour, D., Goldberg, R. D., Wang, T. S., Cullis, A. G., Clarysse, T., & Vandervorst, W. (2000). Cluster formation during annealing of ultra-low-energy boron-implanted silicon. *J. Vac. Sci. & Tech. B,* **18**(1), 435-439.

Downey, D. F., Falk, S. W., Bertuch, A., F., & Marcus, S. D. (1999). Effects of "fast" rapid thermal anneals on sub-keV boron and $BF_2$ ion implants. *J. Elec. Mat.,* **28**(12), 1340-1344.

Gelpey, G., Elliot, K., Camm, D., McCoy, S., Ross, J., Downey, D. F., & Arevalo, E. A. (2002). *Advanced annealing for sub-130 nm junction formation.* Paper presented at the *201st Meeting of the ECS*, Philadelphia, PA.

Gunawan, R., Jung, M. Y. L., Seebauer, E. G., & Braatz, R. D. (2003). Maximum *a posteriori* estimation of transient enhanced diffusion energetics. *AIChE J.*, in press.

Haynes, T. E., Eaglesham, D. J., Stolk, P. A., Gossmann, H. J., Jacobson, D. C., & Poate, J. M. (1996). Interactions of ion-implantation-induced interstitials with boron at high concentrations in silicon. *Appl. Phys. Lett.,* **69**(10), 1376-1378.

Jain, S. C., Schoenmaker, W., Lindsay, R., Stolk, P. A., Decoutere, S., Willander, M., & Maes, H. E. (2002). Transient enhanced diffusion of boron in Si. *J. Appl. Phys.,* **91**(11), 8919-8941.

Jung, M. Y. L., Gunawan, R., Braatz, R. D., & Seebauer, E. G. (1999). Detailed TED modeling of transient enhanced diffusion in implanted Si. *AIChE Annual Meeting*, Dallas, TX. Paper 189d.

Jung, M. Y. L., Gunawan, R., Braatz, R. D., & Seebauer, E. G. (2001). Surface Fermi level pinning: An electrical "valve" in transient enhanced diffusion. *Materials Research Society Spring Meeting*, San Francisco, CA. Paper J4.21.

Kobayashi, H., Nomachi, I., Kusanagi, S., & Nishiyama, F. (2001). Lattice site location of ultra-shallow implanted B in Si using ion beam analysis. In: *Si Front-end Processing - Physics & Technology of Dopant-Defect Interactions III,* pp. J5.3. MRS, Inc., Warrendale, PA.

Laidler, K. J. (1987). *Chemical Kinetics.* Harper & Row, New York, NY.

Law, M. E., & Tasch, A. (2000). Florida object oriented process simulator (FLOOPS) 2000.

Ma, D. L., & Braatz, R. D. (2001). Worst-case analysis of finite-time control policies. *IEEE Trans. on Control Sys. Tech.,* **9**(5), 766-774.

Mannino, G., Stolk, P. A., Cowern, N. E. B., Boer, W. B. d., Dirks, A. G., Roozeboom, F., Berkum, J. G. M. v., Woerlee, P. H., & Toan, N. N. (2001). Effect of heating ramp rates on transient enhanced diffusion of ion-implanted silicon. *Appl. Phys. Lett.,* **78**(7), 889-891.

Press, W. H., Flannery, B. P., Teukolsky, S. A., and Vetterling, W. T. (1992). *Numerical Recipes in C: The Art of Scientific Computing.* Cambridge University Press, New York, NY.

Shishiguchi, S., Mineji, A., Hayashi, T., & Saito, S. (1997). Boron implanted shallow junction formation by high-temperature/short-time/high-ramping-rate (400 °C/sec) RTA. *Symposium on VLSI Technology*, Japan.

Stolk, P. A., Gossmann, H. J., Eaglesham, D. J., Jacobson, D. C., Rafferty, C. S., Gilmer, G. H., Jaraiz, M., Poate, J. M., Luftman, H. S., & Haynes, T. E. (1997). Physical mechanisms of transient-enhanced dopant diffusion in ion-implanted silicon. *J. Appl. Phys.,* **81**(9), 6031-6050.

Vortek Industries Ltd. (2002). Private communication.

Zeghbroeck, B. V. (2002). http://ece-www.colorado.edu/~bart/book/.