

Two-dimensional dopant profiling of submicron metal–oxide–semiconductor field-effect transistor using nonlinear least squares inverse modeling

N. Khalil, J. Faricelli, and C.-L. Huang

Ultra Large Scale Integration Operation Group, Digital Semiconductor, Hudson, Massachusetts 01749

S. Selberherr

Institute for Microelectronics, Technische Universität Wien, Gusshausstrasse 27-29, A-1040 Vienna, Austria

(Received 3 February 1995; accepted 20 September 1995)

We present an inverse modeling technique to determine the two-dimensional (2D) dopant profile of a metal–oxide–semiconductor field-effect transistor from electrical measurements. In our method, the profile is formulated using two tensor product splines (TPSs). This analytical representation is general, compact, and flexible. It simplifies the profile determination problem to the extraction of the TPS coefficients from experimental data. We show the results of applying the new technique on data collected from a sub-0.5 μm complementary metal–oxide–semiconductor technology with various source/drain implants. We also compare the measured and simulated $I-V$ and $C-V$ characteristics. The results illustrate the importance of accurate 2D dopant profiles for short-channel device simulation and modeling. © 1996 American Vacuum Society.

I. INTRODUCTION

With the scaling of metal–oxide–semiconductor field-effect-transistor (MOSFET) dimensions into the submicrometer regime, the influence of the dopant distribution on short-channel device characteristics increases dramatically. The complex multidimensional fields created by the doping become one of the most important factors in determining the electrical behavior of MOSFET. One-dimensional (1D) profiling tools such as spreading resistance (SRP) and secondary ion mass spectrometry (SIMS) are available. However, due to the shallow vertical and lateral junctions, the proximity effects, and the interaction between dopants species, 1D profiles are less indicative of actual 2D profiles. Attempts to extend the 1D profiling tools to higher dimensions (e.g., 2D SRP, 2D SIMS) have met with limited success when applied to state-of-the-art complementary metal–oxide–semiconductor (CMOS) technology.¹⁻³ New techniques aimed at addressing these shortcomings are under development.

Scientists in other fields, such as in geophysics, facing a similar lack of direct experimental measurements resort to inverse modeling.^{4,5} Inverse techniques deal with the determination of parameter values of a physical system from experimental measurements, along with the physical laws and theories that relate the inputs of the system to its outputs. For semiconductors, the theoretical relationship takes the form of the basic semiconductor equations, namely Poisson's equation and the current continuity equations.⁶ In the case of thermal equilibrium and negligible current flow, the solution of the continuity equations can be ignored. Hence, the space charge density within a device can be calculated by solving Poisson's equation

$$\Delta^2 \psi = -\frac{q}{\epsilon} (n - p + N_D^+ - N_A^-), \quad (1)$$

where q is the elementary charge; ϵ the semiconductor per-

mitivity, ψ the electrostatic potential; n and p the electron and hole concentrations; and N_D^+ and N_A^- the donor and acceptor concentrations.

The charges associated with the device terminals are calculated by integrating the space charge density over a device region \mathcal{A} :

$$Q_i = \int_{\mathcal{A}} (n - p + N_D^+ - N_A^-) d\mathcal{A} \quad (2)$$

or by applying Gauss's law to calculate the gate charge by evaluating the line integral of the gate electric field on a closed loop surrounding the gate:

$$Q_i = \oint (\mathbf{E} \cdot \hat{n}) dl. \quad (3)$$

The device capacitances are then approximated by numerically differentiating the terminal charges:

$$C_{ij} = \frac{\delta Q_i}{\delta V_j} \approx \frac{\Delta Q_i}{\Delta V_j}. \quad (4)$$

As an inverse problem, the profile extraction consists of finding the profile that minimizes the weighted least squares fit criterion (SSQ) between experimental and simulated capacitance values:

$$\text{SSQ} = \sum w_i (C_i^{\text{experimental}} - C_i^{\text{simulated}})^2, \quad (5)$$

where the w_i are the weights associated with each data point. This is a continuous minimization problem. The target is to determine the complete functional variation of the profile. We convert it to a discrete problem by representing the profile variation using B-splines in 1D and tensor product spline (TPS) in 2D. With a fixed sequence of breakpoints, i.e., knots, the inverse problem simplifies the profile extraction to the determination of the B-spline or TPS coefficients from the capacitance data.

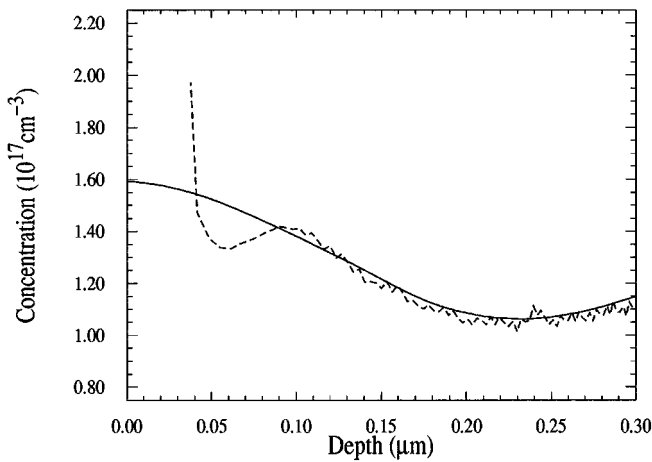


FIG. 1. MOS capacitor channel doping extracted from deep depletion $C-V$ data: inverse modeling (solid line) and analytical results (dashed line).

The rest of the article is divided as follows: The extraction procedures and related topics are presented in the next section. In Sec. III, the extracted profiles of a subhalf micrometer CMOS process that has various source/drain implants are shown. Section III also shows a comparison of experimental $I-V$ and $C-V$ characteristics with simulations using the extracted profiles. We discuss the limitations and accuracy of our method in Sec. IV and present a preliminary assessment of its resolution. Finally, in the conclusion, we offer a list of open questions for future work.

II. EXTRACTION PROCEDURE

In Refs. 7 and 8 a technique for the determination of 2D MOSFET dopant profiles from gate and source/drain (S/D) capacitance measurements is presented. A brief overview of the method for a P -channel MOSFET is given as follows:

(1) Input parameters determination: Several important parameter values are obtained by independent experimental means. The oxide thickness (t_{ox}) is determined from capacitance measurement in the accumulation region as suggested in Ref. 9. The polysilicon gate length (L_p) and electrically active polysilicon gate concentration (N_p) are extracted by matching experimental gate-to-channel capacitances (C_{gc}) and simulated results that take into account quantum mechanical and polysilicon depletion effects.⁹⁻¹¹ The S/D diode acceptor profile is determined using SIMS measurements.

(2) Starting 2D profile generation: The SIMS S/D acceptor profile, the S/D donor profile extracted from diode capacitances, and the channel profile extracted from deep depletion capacitance data are combined into an initial 2D profile. We use a subdiffusion factor for rotating the S/D acceptor profile to obtain the measured effective electrical channel length for the device.

(3) 2D extraction: Using the gate overlap and S/D diode capacitance measurements taken on a fingered polysilicon structure over active region with varying bias voltages, the TPS coefficients are adjusted to achieve a good fit between

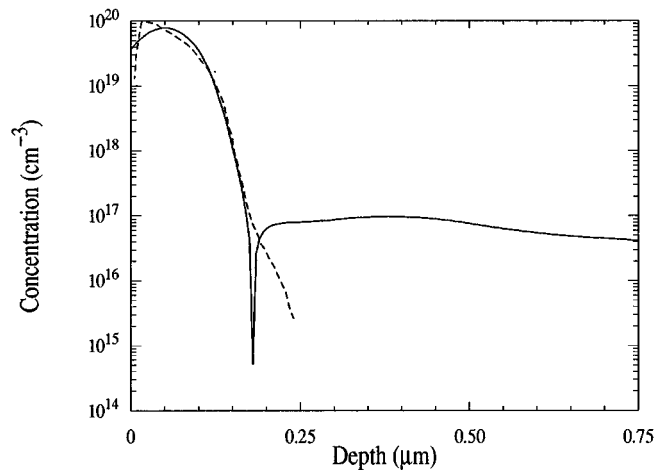


FIG. 2. Source/drain doping profiles: SIMS acceptors profile (dashed line) and net doping extracted by inverse modeling from area diode capacitance.

simulated and experimental values. In the rest of this section, some important aspects of our implementation are presented.

A. Profile TPS representation

As stated previously, the determination of the dopant profile is converted into a discrete parameter extraction problem using the TPS representation. A TPS is a generalization of the one-dimensional polynomial piecewise B-spline functions to a multidimensional space.¹² It is defined by the two knot sequences t_x and t_y and the values of the coefficients c_{ij} :

$$f(x, y) = \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} c_{ij} \times B_{i,k_x,t_x}(x) \times B_{j,k_y,t_y}(y), \quad (6)$$

where $B_{i,k,t}$ is the i th B-spline of order k for the knot sequence t , and n_x and n_y are the number of knots in the X and Y direction, respectively. Using a fixed number of knots and locations, a 2D function can then be written as

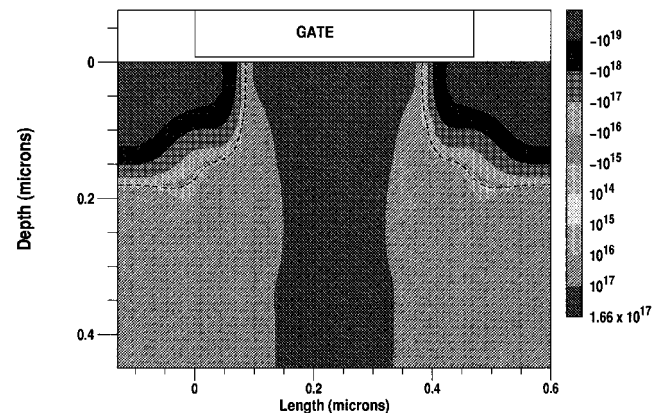


FIG. 3. P -channel extracted 2D net doping for a device with $L_p=0.45 \mu\text{m}$.

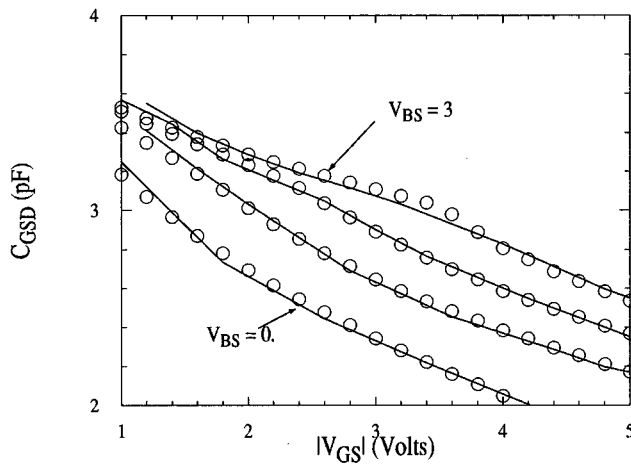


FIG. 4. Comparison of experimental (symbols) and simulated (solid lines) gate-to-source-drain capacitance ($W/L_p=3696/0.42 \mu\text{m}$).

$$f(x,y) = \sum_{i=1}^{nx \times ny} c_i \times B_{i,xy}. \quad (7)$$

In view of the wide range of the doping concentrations, the logarithmic dependency of the acceptor and donor concentrations are represented by two TPSs. Each TPS uses a different sequence of knots to accommodate the varying concentration fields. The net doping can then be written as

$$\rho(x,y) = \mathcal{F}(\boldsymbol{\alpha}, \boldsymbol{\delta}), \quad (8)$$

where $\boldsymbol{\alpha}$ and $\boldsymbol{\delta}$ are the vectors of TPS coefficients for the acceptors and donors, respectively. This analytical formulation of the profile is more compact than a straightforward mesh representation. As a result, the amount of computation required is decreased, and a better condition for the least squares problem is achieved. Moreover, this representation does not assume *a priori* knowledge of the profile functional variation. To ensure the smoothness and continuity of the profile, quadratic or cubic splines are used to represent the profile functional variation in the horizontal and vertical directions. Choosing the appropriate number of knots in each direction is important. Although a large number of parameters will result in a better fit to the data, this does not ensure accurate determination of the profile as the variance errors and the computation time are increased. We are presently investigating ideas similar to Ref. 13 for the determination of the number and location of knots. We typically use four or five knots in each direction and determine their location based on process information such as junction depth, spacer width, and gate length. Another feature of the TPS representation is the ease of including 1D information. For example, in generating the 2D starting profile, the coefficients of the 1D long channel profile B-spline become the TPS coefficients in the middle of the short-channel device.

B. Capacitance calculation

The charge integration method is used to calculate the device capacitances by numerical differentiation. This proce-

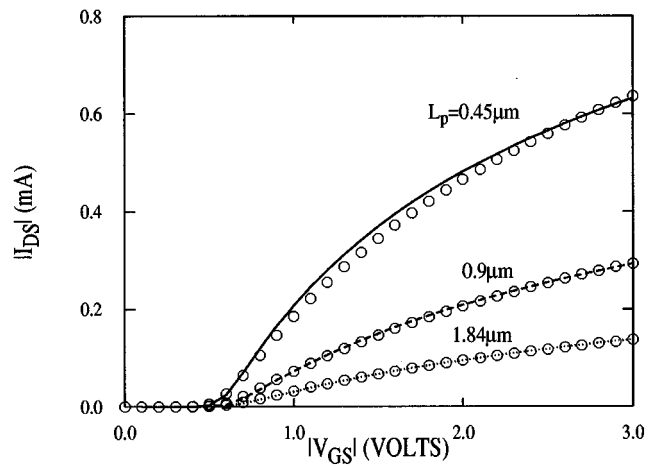


FIG. 5. Comparison of measured (symbols) and simulated (lines) $I-V$ characteristics in the linear region ($V_{DS}=-50 \text{ mV}$) for three gate lengths ($L_p=0.45, 0.9, \text{ and } 1.84 \mu\text{m}$) with $W=64 \mu\text{m}$, $N_p=2.7 \times 10^{19} \text{ cm}^{-3}$, and $t_{ox}=73 \text{ \AA}$.

cedure is inherently prone to numerical roundoff and integration errors. We resort to the following strategies to minimize their effects:

(1) We use central differences for numerical differentiation:

$$C_i \approx \frac{Q_2 - Q_1}{V_2 - V_1}, \quad (9)$$

where Q_2 and Q_1 are the terminal charges at V_2 and V_1 , respectively, $V_2 = V + \delta V$, and $V_1 = V - \delta V$. The size of the voltage step δV should be small enough to achieve a good linearized approximation of the derivative. However, a small

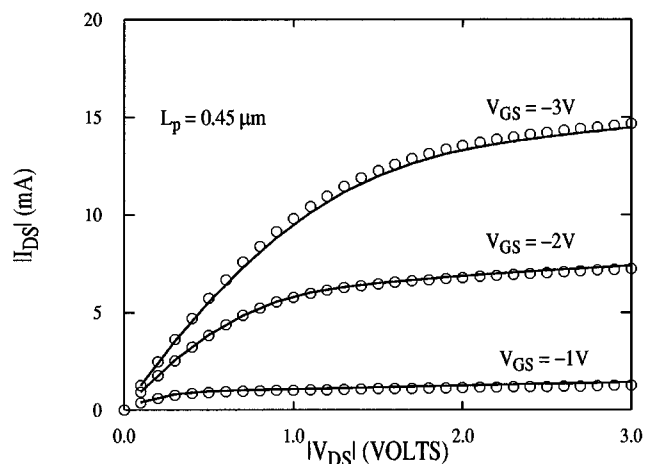


FIG. 6. Comparison of measured (symbols) and simulated (lines) $I-V$ results in the linear and saturation regions for device with $L_p=0.45 \mu\text{m}$, $W=64 \mu\text{m}$, $N_p=2.7 \times 10^{19} \text{ cm}^{-3}$, and $t_{ox}=73 \text{ \AA}$.

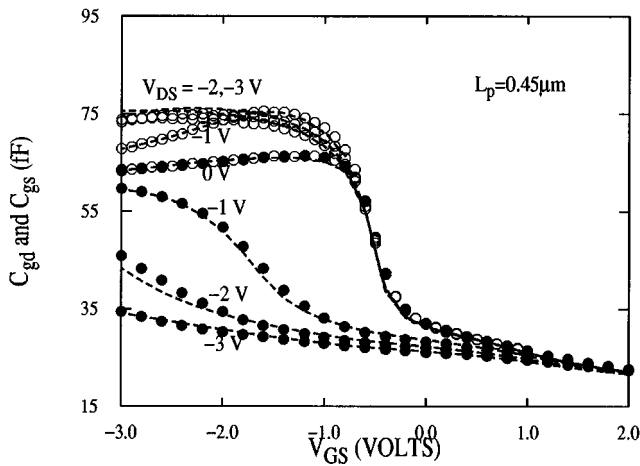


FIG. 7. Comparison of measured (symbols) C_{gd} (lower three curves) and C_{gs} (upper three curves) capacitances with simulated results (lines) as a function of V_{GS} and V_{DS} for device with $L_p=0.45 \mu\text{m}$, $W=64 \mu\text{m}$, $N_p=2.7 \times 10^{19} \text{cm}^{-3}$, and $t_{ox}=73 \text{ \AA}$.

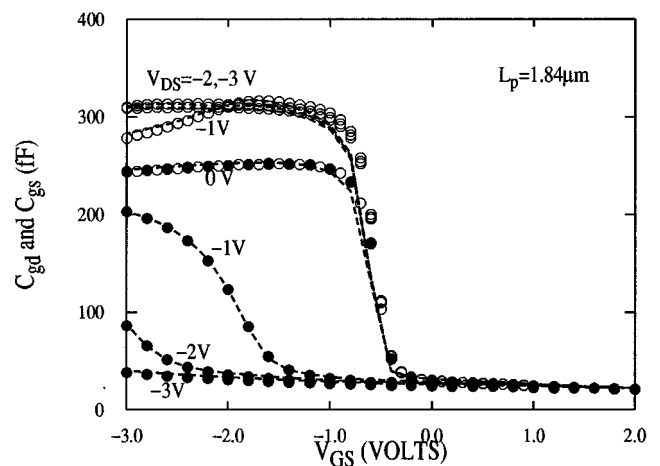


FIG. 8. Comparison of measured (symbols) C_{gd} (lower three curves) and C_{gs} (upper three curves) capacitances with simulated results (lines) as a function of V_{GS} and V_{DS} for device with $L_p=1.84 \mu\text{m}$, $W=64 \mu\text{m}$, $N_p=2.7 \times 10^{19} \text{cm}^{-3}$, and $t_{ox}=73 \text{ \AA}$.

change in the voltage yields a small change in the charge. This could result in numerical roundoff errors in calculating the quotient due to subtractive cancellation. We found that for room temperature measurements, a 20 mV step is a reasonable compromise between the two conflicting requirements.

(2) Discretization is a major source of errors in the calculations. In principle one could increase the size of the grid to effectively eliminate numerical approximation errors. While this is convenient and feasible for 1D problems, the computational demands in the 2D case prohibit such an approach. For that reason, we carefully refine the Poisson solver grid where the space charge density varies rapidly.

(3) We use the same grid at V_1 and V_2 in our solution. This reduces the integration errors since some of the errors in the calculation of the charges are canceled out by taking the difference.

(4) When applicable, we only integrate the charges in the device that contribute to the variation. This minimizes the loss of significant digits in the capacitance calculation. For example, in calculating the bulk charge, the ionized dopant atoms in the S/D diode region are not included in the summation because they remain constant when a small voltage perturbation is applied.

C. Least squares optimization

We solve the nonlinear least squares problem using the well-established Levenberg–Marquardt algorithm with linear constraints.^{14,15} This is a Gauss–Newton algorithm with trust region modification that is versatile and robust. We use finite difference derivatives to approximate the Jacobian. In the following, we present some of the specific characteristics of our solution:

(1) Starting guess: A good initial guess is important to avoid the trapping of the solution in local minimas and to

limit the number of iterations required before convergence. This is especially critical in 2D where the computational demands are great.

(2) Parameter redundancy: Poisson’s equation contains a net doping term only. In the presence of p - n junctions, this introduces a direct correlation between the acceptor and donor coefficients, especially for knots located near the junction. This could result in a rank deficient Hessian. We bypass this problem by a two-step iteration scheme in which we fix one set of coefficients while allowing the coefficients of the other type to change.

(3) Restriction of the number of parameters: We only extract the coefficients at knots when they are of the same type as the net doping at that location. Moreover, at each step, we analyze the eigenvalues of the approximate Hessian. Coefficients associated with small eigenvalues are not allowed to change.

(4) Parameters bounds: To avoid deviations outside the expected range during the initial stages of the solution, we enforce linear bounds on the parameters.

(5) Computation time: The amount of computation needed to extract the 2D profile is approximately 36 h of CPU time on a desktop Alpha AXP workstation model 3000-400. Software strategies to distribute the capacitance calculations on multiple workstations can significantly reduce this.

III. RESULTS

We applied our method to data collected from devices fabricated using a retrograde n -well, salicided dual-gate CMOS process. All experimental I - V and C - V characteristics were obtained using an HP 4145B parameter analyzer and an HP 4275A LCR meter. The measurement frequency of the HP 4275A LCR meter was set to 100 kHz. The resolution of the system is around 0.1 fF. In order to reduce the noise level in the measured results, the experimental C - V

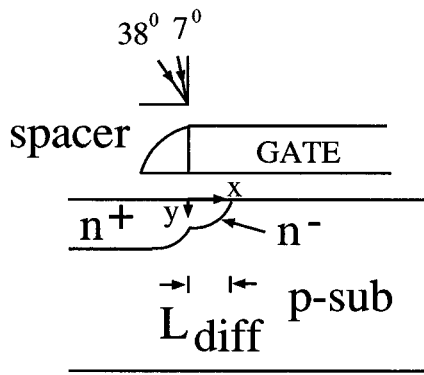


FIG. 9. A schematic diagram of half an N -channel MOSFET with an n -implant angle.

data for sub- $0.5\text{-}\mu\text{m}$ devices were averaged for several measurements. As a result, the actual resolution of experimental data is better than 0.1 fF .

Using deep depletion capacitance data taken on a large MOS capacitor we first extract the B-spline coefficients of the vertical channel doping. Figure 1 shows the extracted profile together with an analytically extracted profile. In contrast to the analytical results that fail to determine the doping near the interface, the inverse modeling profile is extracted up to the Si/SiO₂ interface. We note, however, that there is a correlation between the gate work function, oxide, and interface charge densities (Q_{ox} and Q_{it}), and the doping near the interface. One strategy we have recently begun to use is to assume a small fixed oxide and interface charge and use a polysilicon work function that includes band gap narrowing effects. We then let the channel doping adjust itself to fit the measured threshold.

For the S/D diode, the acceptor B-spline coefficients were determined by curve fitting the 1D SIMS profile. The SIMS profile was also used in the simulation of the reverse junction capacitance data. The B-spline coefficients of the S/D donors profiles were extracted by matching the experimental diode area capacitance. Figure 2 shows the resultant S/D net doping as well as the SIMS acceptor profiles.

The channel and S/D profiles were then combined into an initial 2D profile. By carefully matching the experimental capacitances, we ensure that the initial profile is in the immediate proximity of the solution. The extracted 2D dopant profile for the P -channel MOSFET is shown in Fig. 3. Figure 4 illustrates the good fit achieved between the experimental data used in the extraction and simulation with the extracted 2D profile.

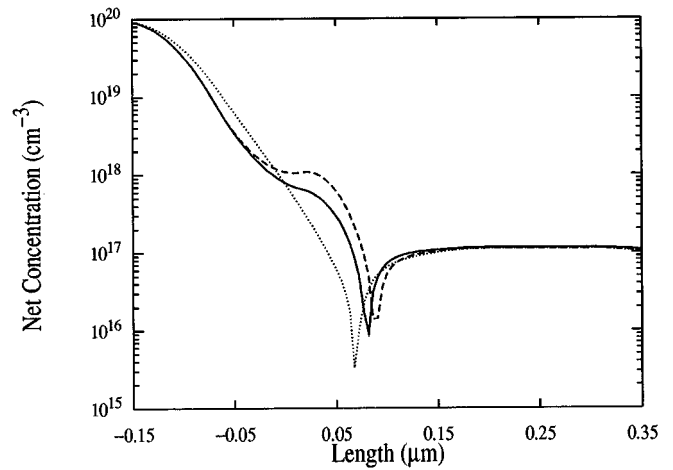


FIG. 10. Extracted one-sided net doping at the SiO₂/Si interface (device A—solid line; device B—dashed line; device C—dotted line).

For validation, we compare simulated and experimental $I-V$ and $C-V$ characteristics using the extracted profile as input. Since the 2D profile, polysilicon gate concentration (N_p), oxide thickness (t_{ox}), and polysilicon gate length (L_p) are obtained experimentally, 2D device simulation is expected to accurately reproduce experimental $I-V$ and $C-V$ characteristics over a wide range of biases and devices. To ensure a good agreement we adjust the mobility parameters using long-channel device data, as well as the source/drain external resistance R_t , and the gate workfunction of the short-channel devices. Figure 5 shows measured and simulated results for the linear region currents for three gate lengths. Reasonable good agreement is found in all regions of bias for all three lengths. Figure 6 shows a comparison of measured and simulated results for the device with $L_p=0.45\text{ }\mu\text{m}$. Note that a good agreement was also obtained for longer devices (not shown).

We also compare the gate-to-source (C_{gs}), and gate-to-drain (C_{gd}) capacitances for two lengths of P -channel MOSFET in all regions of device operation (the accumulation, linear and saturation regions) in Figs. 7 and 8. Excellent agreement is achieved for both bias-dependent intrinsic and overlap capacitances. We note that an accurate profile is critical for fitting the bias-dependent accumulation capacitance.

Finally, we applied our method to three N -channel devices with varying implant conditions as shown in Table I. In Fig. 9, we show a schematic diagram of half of an N -channel MOSFET with the n -implant angle. The implant conditions for devices A and B corresponds to a large-angle-tilt-

TABLE I. Processing conditions for N -channel MOSFET.

Processing conditions	Device type		
	A	B	C
n^- S/D (ions/cm ² , keV, tilt angle)	2.4×10^{13} , 40, 38°	3.2×10^{13} , 40, 38°	3×10^{13} , 25, 7°
n^+ S/D (ions/cm ² , keV)	5×10^{15} , 40	5×10^{15} , 40	5×10^{15} , 40

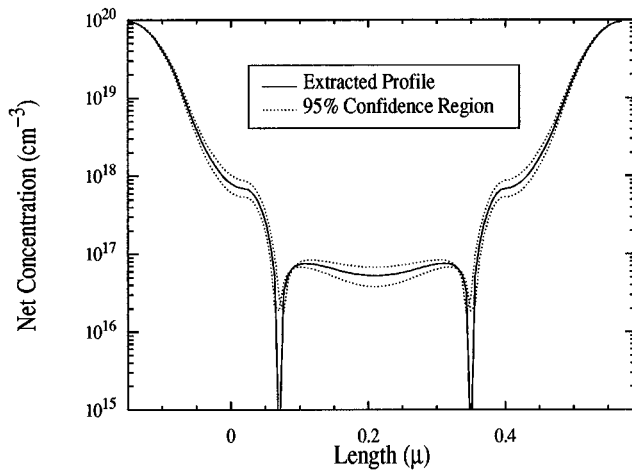


FIG. 11. Extracted N -channel net doping along the SiO_2/Si interface and the corresponding 95% confidence region.

implanted-drain (LATID) implant,¹⁶ while device C received a regular lightly doped drain implant. Figure 10 shows the extracted net doping at the SiO_2/Si interface for the three devices. It is clear that the method is able to resolve the difference in the lateral diffusion length as the implant angle/energy increases. As shown, the method is also successful in determining the characteristic shape of a LATID implant profile.

IV. DISCUSSION

The following points illustrate some of our method characteristics, limitations, and associated uncertainties.

(1) The method capability is limited to the extraction of the electrically active net doping, not the chemical concentration of species atoms. Moreover, the required electrical measurements can only be taken after the completion of processing up to the metal layers. This limits the applicability of the method during process development.

(2) Uncertainties in the device geometrical structure, such as nonplanar surfaces, have a direct effect on the extracted profile. These can be resolved by independent determination of the structural information using TEM imaging. Other types of input uncertainties, e.g., errors in the S/D SIMS profile, can also influence the extraction. We plan on performing a Monte Carlo simulation analysis to estimate the accuracy of the extracted profiles subject to the inherent errors in the inputs.

(3) The confidence intervals of the extracted TPS coefficients were computed according to the algorithms described in Ref. 17 assuming error-free inputs. The size of the confidence interval is an estimate of the standard deviation of the profile. Figure 11 shows the extracted surface net doping for an N -channel device with a 95% confidence region. It indicates that the profile extraction is reasonably accurate in determining the doping level and the lateral junction location.

(4) The extent of the device region where the profile can be determined depends on the range of measurement voltage, the doping level, and the device characteristics. For example,

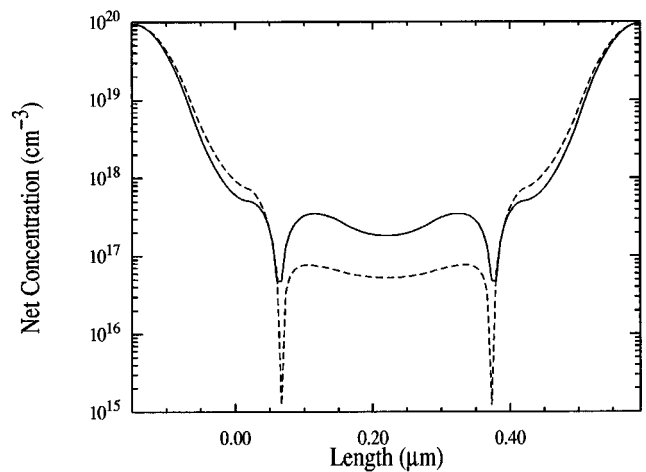


FIG. 12. N -channel net doping at the SiO_2/Si interface extracted with (solid line) and without (dashed line) Gaussian charges to model RSCE.

the capability of the gate to deplete the S/D-gate overlap region of carriers under accumulation bias, without breaking down the gate oxide, limits the resolution of the method for high concentration S/D profiles.

(5) Modeling assumptions in solving Poisson's equations are also a source of uncertainty. For instance, there are two main approaches available in the literature to account for the anomalous threshold voltage variation in short-channel devices. The reverse short-channel effect (RSCE) can be explained by either the existence of oxide charges or dopant profile variation.^{18,19} The chosen method to model this phenomenon influence the extracted profile near the interface. Figure 12 compares the net doping along the surface in two cases. In the first one, oxide Gaussian charges were introduced in the simulation to model the RSCE. In the second, the oxide charge value was fixed to that determined from long-channel device data. As seen, the two profiles are clearly different. We are presently investigating the inclusion of other types of data to eliminate this uncertainty. In particular, the use of subthreshold current data in the extraction might clarify the issue.

V. CONCLUSIONS

We have described a nondestructive method for MOSFET 2D profile determination. The method does not require any difficult sample preparation and uses measurements easily performed during process characterization. It fills an important gap in the available process metrology tools. The preliminary results are very encouraging; however, further development of the technique is still needed. In particular, the following areas deserve more attention:

(1) Development of an algorithm for selecting the appropriate number of TPS knots and their locations.

(2) A quantitative study of the resolution and accuracy of the method is still needed.

(3) A comparison between the profiles extracted using the inverse modeling technique and profiles determined by direct

measurements, when possible, could serve as a cross validation check.

(4) Finally, the extension of the method to include other sources of electrical data, such as subthreshold current measurements, might improve the resolution beyond what is possible when relying solely on capacitance data.

- ¹S. H. Goodwin-Johansson, R. Subrahmanyam, C. E. Floyd, and H. Z. Massoud, *IEEE Trans. Comput.-Aided Design* **CAD-8**, 323 (1989).
- ²R. Subrahmanyam, H. Z. Massoud, and R. B. Fair, *Appl. Phys. Lett.* **52**, 2145 (1988).
- ³S. Kordic, E. Van Leonen, D. Dijkkamp, A. Hoeven, and H. Moraal, *IEDM Technol. Dig.* **1989**, 277 (1989).
- ⁴A. Tarantola, *Inverse Problem Theory* (Elsevier, Amsterdam, 1987).
- ⁵G. J. L. Ouwerling, Ph.D. dissertation, The Delft University of Technology (1989).
- ⁶S. Selberherr, *Analysis and Simulation of Semiconductor Devices* (Springer, Wien, 1984).
- ⁷N. Khalil, J. Faricelli, D. Bell, and S. Selberherr, *Dig. Symp. VLSI Technol.* **1994**, 131 (1994).
- ⁸N. Khalil, J. Faricelli, D. Bell, and S. Selberherr, *IEEE Electron Device Lett.* **EDL-16**, 17 (1995).
- ⁹R. Rios and N. D. Arora, *IEDM Technol. Dig.* **1994**, 616 (1994).
- ¹⁰P. Habaš and J. Faricelli, *IEEE Trans. Electron Devices* **ED-39**, 1496 (1992).
- ¹¹S. Selberherr, A. Schütz, and H. W. Pötzl, *IEEE Trans. Electron Devices* **ED-27**, 1540 (1980).
- ¹²C. De Boor, *A Practical Guide to Splines* (Springer, New York, 1978).
- ¹³A. T. Watson, P. C. Richmond, P. D. Krieg, and T. M. Tao, *SPE Reservoir Eng.* **3**, 953 (1988).
- ¹⁴K. Levenberg, *Q. Appl. Math.* **2**, 164 (1944).
- ¹⁵D. W. Marquardt, *J. Soc. Indus. Math Appl. Math. (SIAM)* **11**, 431 (1963).
- ¹⁶T. Hori, *IEDM Technol. Dig.* **1989**, 777 (1989).
- ¹⁷M. Sharma and N. D. Arora, *IEEE Trans. Comput.-Aided Design* **CAD-12**, 982 (1993).
- ¹⁸H. Jacobs, A. V. Schwerin, D. Scharfetter, and F. Lau, *IEDM Technol. Dig.* **1993**, 307 (1993).
- ¹⁹C. S. Rafferty, H. H. Vuong, S. A. Eshraghi, M. D. Giles, M. R. Pinto, and S. J. Hillenius, *IEDM Technol. Dig.* **1993**, 311 (1993).