

## Taking Silicon to the Limit: Challenges and Opportunities

by Tsu-Jae King

Silicon-based CMOS transistors can be scaled well into the sub-10 nm regime. However, new materials and processes, in conjunction with advanced transistor structures, will be needed for nanometer-scale MOSFETs to meet performance specifications in the International Technology Roadmap for Semiconductors (ITRS). This paper discusses challenges for achieving target performance metrics at the end of the Roadmap, and approaches to overcoming them.

The steady miniaturization of the metal-oxide-semiconductor field-effect transistor (MOSFET) with each new generation of complementary-MOS (CMOS) technology has yielded continual improvements in integrated-circuit performance (speed) and cost per function over the past several decades, to usher in the Information Age. Continued transistor scaling will not be as straightforward in the future as it has been in the past, however, because fundamental materials and process limits are rapidly being approached.<sup>1</sup> New materials and processes, as well as non-classical transistor structures, will be necessary in order to extend CMOS technology to the last node of the ITRS.<sup>2</sup> Minimization of leakage current, parasitic resistance, and capacitance to minimize power consumption and maximize circuit performance, and reduction in device-to-device variability to increase yield (and thereby lower cost), will be key challenges for sustaining the rapid growth of the industry to usher in the age of ambient intelligence and ubiquitous computing. This paper discusses recent CMOS technology developments and remaining work needed to address these challenges.

### Advanced Transistor Structures and Materials

In order to scale the classical bulk-Si MOSFET structure (Fig. 1a) down to the 10 nm physical gate length ( $L_g$ ) regime, heavy halo and channel doping (greater than  $1 \times 10^{18} \text{ cm}^{-3}$ ) will be required to suppress leakage current and short-channel effects.<sup>3</sup> As a result, field-effect carrier mobilities will be degraded, resulting in incommensurate improvements in transistor drive current with  $L_g$  scaling.<sup>4</sup> Thin-body transistor structures (Figs. 1b and 1c)<sup>5</sup> rely not on heavy channel doping but on a sufficiently thin body/channel region ( $T_{Si} < L_g$ ) to limit leakage current. The use of a lightly doped or undoped channel provides

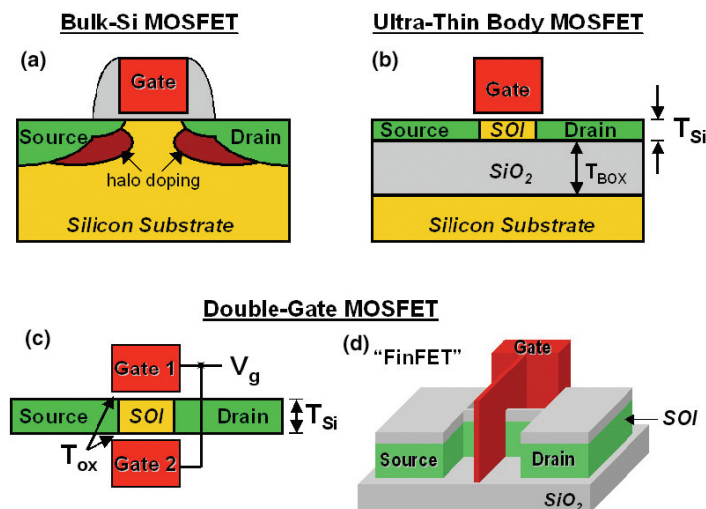


Fig. 1. Schematic diagrams of MOSFET structures: (a) classical bulk-Si, (b) ultrathin-body (UTB), (c) double-gate (DG), and (d) FinFET.

immunity to variations in threshold voltage ( $V_T$ ) resulting from statistical dopant fluctuations in the channel, as well as enhanced carrier mobility for higher transistor drive current because of the lower transverse electric field in the inversion layer.<sup>4</sup> Therefore, thin-body MOSFETs offer improved circuit performance as compared to the bulk-Si transistor structure (Fig. 2).<sup>6</sup> To provide a means for adjusting  $V_T$  without channel doping during the manufacturing process, a tunable-work-function gate technology is necessary. For thin-body CMOSFETs, the required range of gate work function ( $\Phi_M$ ) tunability is from 4.5 eV to 5.0 eV.<sup>7</sup>

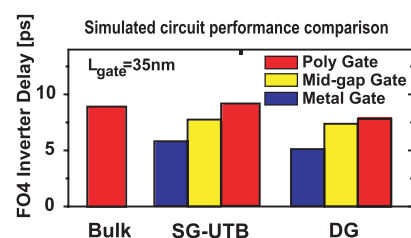


Fig. 2. Loaded-inverter delay comparison of bulk-Si vs. UTB vs. DG CMOS technologies, obtained through mixed-mode simulation using realistic device structures based on ITRS specifications.<sup>6</sup>

### Advanced Transistor Structures

**Ultrathin-body (UTB) FET**—The body thickness  $T_{Si}$  must be less than  $\sim L_g/3$  in a UTB FET in order to adequately suppress leakage current.<sup>5</sup> Because of quantum confinement effects,  $V_T$

becomes a sensitive function of  $T_{Si}$  for thicknesses below 5 nm.<sup>8</sup> Also, carrier mobilities are degraded due to enhanced interface roughness scattering for  $T_{Si} < 4 \text{ nm}$ .<sup>9</sup> For these reasons, it may be difficult to scale the UTB FET structure to below 12 nm  $L_g$ , unless techniques for achieving uniformly thin films with atomically smooth surfaces/interfaces are used. One example is the “Silicon-on-Nothing” fabrication process,<sup>10</sup> in which the ultrathin Si channel and the buried oxide are defined by epitaxy on a bulk-Si substrate, so that thickness control can be as fine as a single atomic layer.

**Double-gate (DG) FinFET**—The quasi-planar FinFET (Fig. 1d) offers the superior scalability of the DG MOSFET structure together with a process flow and layout similar to that of the conventional MOSFET.<sup>11</sup> Hence, it has been investigated by many companies.<sup>12–14</sup> FinFETs with gate lengths down to 10 nm have already been demonstrated and exhibit excellent control of short-channel effects.<sup>14,15</sup> One advantage of this vertical transistor structure is that it is relatively immune to gate line-edge roughness, a major source of variability in planar nanoscale FETs.<sup>16</sup> FinFET performance variability due to variations in fin width is a potential issue, but can be minimized by using a spacer lithography process.<sup>17</sup> In order to optimize the trade-off between parasitic series resistance and parasitic gate capacitance, a gate-underlapped structure (in which

the electrical channel length is larger than the physical gate length) will be required to achieve peak circuit performance for sub-20 nm  $L_g$ .<sup>18</sup> Parasitic series resistance and contact resistance will therefore ultimately limit FinFET performance in the nanoscale regime (Fig. 3).<sup>19</sup> The use of thick source/drain (S/D) regions, e.g. formed by selective growth of Si<sup>12</sup>, Si<sub>1-x</sub>Ge<sub>x</sub>, or Ge,<sup>20</sup> can help to alleviate this issue, particularly if low specific contact resistivity ( $\rho_c < 10^{-8} \Omega\text{-cm}^2$ ) contacts can be formed by silicidation/germanidation of the S/D fin surfaces. It should be noted that the width of the gate-sidewall spacers, which isolate the gate electrode from the raised-S/D regions, must be carefully optimized for peak circuit performance.<sup>21</sup>

Advanced transistor structures can be applied to improve the scalability of memory devices as well, to achieve very high density, non-volatile information storage. FinFET SONOS (silicon-oxide-nitride-oxide-silicon) memory devices have already been demonstrated, and exhibit excellent retention and endurance characteristics.<sup>22</sup> Because the effective gate-dielectric thickness of the ONO stack is relatively large (~10 nm), the body thickness of a SONOS memory device must be even thinner than for a logic device (a significant challenge for fabrication), in order to adequately suppress leakage current (Fig. 4).<sup>22</sup> If the two gate electrodes are electrically isolated (e.g., by applying a

chemical mechanical polishing step or a masked etch step during the fabrication process), then 2-bit storage is possible, to further increase storage density.<sup>23</sup> In order to allow each bit to be distinguished, asymmetric gates are required ( $n^+$  poly-Si for the front gate, and  $p^+$  poly-Si for the back gate). These can be obtained in a straightforward manner, via high-tilt-angle implants to dope the front and back gates separately.<sup>24</sup>

**Back-Gated (BG) UTB FET**—Power consumption will be a primary design constraint for sub-65 nm CMOS technologies, so that active leakage ( $V_T$ ) control will be necessary for optimization of energy vs. delay trade-offs in future ULSI systems. The  $V_T$  of a FinFET cannot be dynamically changed; however, if the two gate electrodes are electrically isolated so as to allow independent operation,<sup>25</sup> the FinFET can be operated as a back-gated UTB FET with the capability for dynamic  $V_T$  control. It should be noted that the optimal BG UTB FET design employs significantly different gate-oxide thicknesses for the front and back gates.<sup>26</sup> Techniques such as selective (tilted) implantation of nitrogen,<sup>27</sup> oxygen,<sup>28</sup> or argon<sup>29</sup> to simultaneously grow gate oxides of different thicknesses can be used for the FinFET structure. A planar BG UTB FET structure may ultimately be more area-efficient because the back gate would be routed in a separate layer than the front gate. The development of a cost-effective fabrication process for the planar BG UTB FET with

self-aligned gates (for optimal performance), as well as the need to achieve  $T_{Si} < 4$  nm with excellent uniformity and oxide-interface quality, remains a challenge.

**Extending the Roadmap**—The scaling limit of the Si MOSFET is well below 10 nm  $L_g$  (depending on the leakage current specification and power-supply voltage), and corresponds to the point where direct tunneling of carriers from the source to the drain in the off state becomes prohibitively large.<sup>30</sup> Practically, the MOSFET scaling limit will be determined by the degree to which the channel film thickness can be controlled in a manufacturing process. Nanofabrication techniques (e.g. self-assembly) may be useful for achieving the uniformly thin and smooth channel films required to reduce performance variations to an acceptable level, as well as for improving critical-dimension (i.e.,  $L_g$ ) control. Further scaling of  $L_g$  can also be enabled by improving the transistor design (e.g. so that the effective channel length or  $V_T$  is larger in the off state than in the on state) and/or by employing an alternative semiconductor material (one with lower dielectric permittivity to reduce drain-induced barrier lowering). Clearly, opportunities abound for innovations by technologists and device designers alike to extend transistor scaling toward atomic dimensions.

(continued on next page)

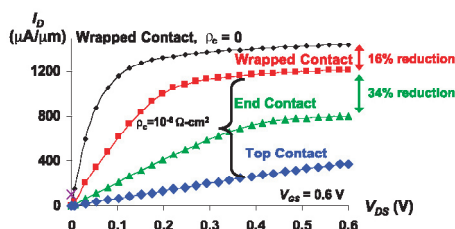
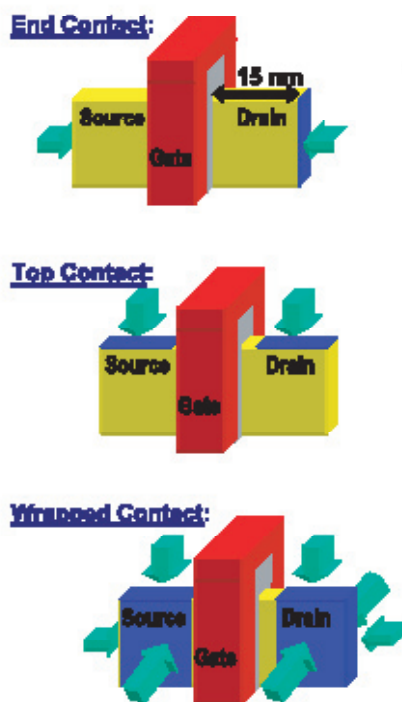


FIG. 3. Impact of S/D contact structure on FinFET drive current, obtained by 3-D device simulation.<sup>19</sup> ( $L_g = 18$  nm,  $T_{Si} = 10$  nm,  $T_{ox} = 5\text{\AA}$ , S/D profiles optimized)

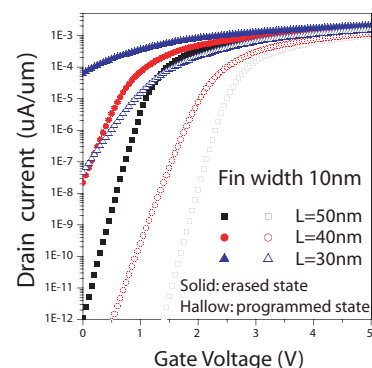


FIG. 4. Transfer characteristics for FinFET SONOS memory (ONO = 3 nm/6.1 nm/4.8 nm), from 2D device simulation.<sup>22</sup> An adequately large difference in currents for erased vs. programmed states is achieved for  $L_g < 40$  nm if a sufficiently thin (10 nm) body is used.

## Tunable-Work-Function Gate Materials

Two metal-gate technologies have been demonstrated to provide a range of  $\Phi_M$  tunability suitable for application to future thin-body CMOSFETs: molybdenum (Mo) and fully-silicided NiSi. Mo is compatible with a conventional (“gate first”) fabrication process;<sup>31</sup> it can be deposited by sputtering or chemical vapor deposition, patterned by conventional reactive ion etching, and is thermodynamically stable on SiO<sub>2</sub>. Nitrogen (N) implantation followed by thermal annealing to form a Mo nitride layer at the gate-oxide interface has been shown to be an effective way to controllably lower the effective  $\Phi_M$  of a Mo gate electrode on SiO<sub>2</sub> gate dielectric, from 5 eV down to 4.4 eV (Fig. 5).<sup>32</sup> Care must be taken to avoid Mo<sup>+</sup> diffusion<sup>33</sup> and gate-oxide damage during sputter deposition;<sup>34</sup> the Mo gate film thickness and N implant energy must be carefully co-optimized to minimize damage (due to implant straggle) to the underlying gate dielectric.<sup>35</sup> A capping layer (e.g., of TiN) deposited *in situ* is beneficial to prevent out-diffusion of the implanted N during thermal annealing, to maximize the  $\Phi_M$  shift for a given implant dose and to improve uniformity,<sup>36</sup> which is critical for precise  $V_T$  control.

The work function of a fully silicided (FUSI) gate material can be adjusted via doping of the precursor Si gate material.<sup>37</sup> Researchers have recently demonstrated that  $\Phi_M$  for a FUSI NiSi gate on SiO<sub>2</sub> gate dielectric can be tuned over a significant range<sup>38</sup> (from 4.5 eV to 4.9 eV, for dopant implant doses up to  $\sim 3 \times 10^{15} \text{ cm}^{-2}$ )<sup>39</sup> and have successfully applied this gate technology to fabricate CMOS FinFETs with nearly symmetrical  $V_T$ 's.<sup>40</sup> The NiSi is formed at low temperature ( $\leq 500^\circ\text{C}$ ) and cannot withstand high annealing temperature; hence, the silicidation must be the last thermal processing step in the transistor fabrication process. The intrinsic tensile stress ( $\sim 0.8 \text{ GPa}$ ) in a NiSi gate induces tensile strain in the Si channel for narrow-width UTB FETs, which enhances both electron and hole carrier mobilities and hence drive current.<sup>41</sup>

It should be noted that  $\Phi_M$  for a metal gate electrode can vary significantly with the gate-dielectric material.<sup>42</sup> In addition, process integration challenges will change with the gate-dielectric material: Ni atoms in a FUSI NiSi gate can easily penetrate HfO<sub>2</sub> during the silicidation process, leading to yield and reliability problems;<sup>43</sup> N implanted into a Mo gate can easily penetrate HfO<sub>2</sub> during subsequent annealing steps and degrade the oxide-silicon interface and hence transistor performance.<sup>44</sup> Therefore, a metal-gate

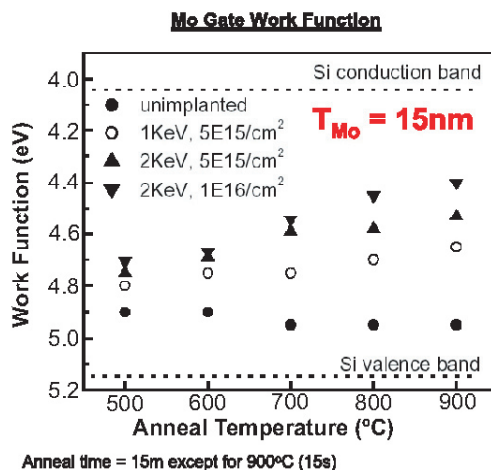


Fig. 5. Effect of N<sup>+</sup> implant and annealing conditions on the effective  $\Phi_M$  of a Mo gate on SiO<sub>2</sub>.<sup>31</sup>

technology must be tailored specifically to the gate-dielectric material. For HfO<sub>2</sub>, FUSI HfSi is a promising gate material (with  $\Phi_M$  tunable in the range 4.23 eV to 4.87 eV) that is stable at high temperatures and, therefore, compatible with a conventional (“gate-first”) planar CMOS fabrication process.<sup>45</sup>

## Performance Enhancement Approaches

Alternative approaches to transistor scaling for continued improvements in system performance and reductions in cost and power consumed per function will ultimately be needed to sustain the rapid growth of the semiconductor industry through the first half of this century. Process technology innovations which improve transistor drive current ( $I_{on}$ ) without sacrificing off-state leakage ( $I_{off}$ ) can further improve the performance vs. power trade-off. Device innovations can provide circuit designers with better “building blocks” to enable more efficient designs. Examples of these are discussed below.

**Carrier-mobility enhancement**—Techniques for increasing the average velocity of carriers in the channel—without significantly impacting cost and device reliability—will be necessary in order for the industry to maintain its historic 17%-per-year performance improvement rate.<sup>2</sup> Approaches to enhancing carrier mobility include the use of a strained capping layer,<sup>46</sup> a strained gate electrode,<sup>47</sup> or strained S/D regions (using epitaxial Si<sub>1-x</sub>Ge<sub>x</sub><sup>48</sup> or silicide<sup>49</sup>), and optimization of the channel surface crystal orientation and current flow direction;<sup>50</sup> indeed, some of these methods are already used in state-of-the-art CMOS products today. In the future, these techniques must be adapted to advanced transistor structures, with manageable device-parameter dependencies (i.e., variation with

transistor channel length and width, and with S/D length) for effective circuit design.

**Metallic-source/drain technology**—In the sub-10 nm  $L_g$  regime, intrinsic variation due to random discrete dopants in the S/D regions will cause large variations in  $I_{on}$  and  $I_{off}$ .<sup>16</sup> The use of metallic-S/D regions rather than doped-S/D regions can eliminate this issue. Sub-20 nm  $L_g$  CMOSFETs with silicide S/D have been successfully fabricated using the UTB structure to achieve low leakage current.<sup>51</sup> The primary challenge for metallic-S/D technology is achieving sufficiently low ( $\leq 0.1 \text{ eV}$ )<sup>52</sup> Schottky barrier height  $\Phi_b$  to meet the Roadmap  $I_{on}$  specifications. Techniques proposed for lowering  $\Phi_b$  include passivation of extrinsic interface states,<sup>53</sup> straining the Si channel,<sup>54</sup> and using very heavily doped “tips” formed by silicidation-induced impurity segregation<sup>55</sup> (which is susceptible to random discrete dopant effects).

**Negative differential resistance devices**—Negative differential resistance (NDR) devices are a prime example of alternative semiconductor devices that can potentially be used to reduce the power consumption and cost of integrated circuits (ICs). The defining feature of an NDR device is that the current flowing between two of its terminals actually decreases as the voltage difference between those two terminals increases over a range of voltages. A key figure of merit for NDR devices is the “peak-to-valley” current ratio (PVCR). The higher the PVCR, the better: A high “peak” current is needed for fast and reliable circuit operation, while a low “valley” current is needed to minimize power consumption. Silicon-based NDR devices typically exhibit a PVCR no greater than 10 at room temperature,<sup>56</sup> decreas-

(continued on page 42)



ing significantly as the temperature is increased.

IC manufacturers have actively investigated devices that exhibit significant NDR behavior since the invention of the Esaki diode.<sup>57</sup> This is because such devices used together with conventional transistors result in more efficient circuits because fewer elements are needed to implement a function.<sup>58</sup> Many innovative circuit designs and applications have been proposed in the literature (see Ref. 58 for an overview) to take advantage of NDR devices, including: compact static memory (SRAM),<sup>59</sup> self-latching logic, analog-to-digital conversion, shift registers, oscillator elements, and multi-valued logic. To date, technological obstacles have hindered their widespread use in silicon-based ICs, however. This is because high-performance NDR devices typically require highly specialized (*i.e.*, expensive) fabrication processes and/or exotic materials, so that they cannot be easily integrated with conventional CMOS devices. Thus, the development of a high-performance (high-PVCR), CMOS-compatible NDR device would constitute a breakthrough advancement in IC technology and have a significant impact on the industry.

## Summary

As compared to the classical bulk-Si MOSFET structure, thin-body transistor structures achieve a better trade-off between performance and power consumption, and also can provide immunity to random variations associated with the discreteness of dopant atoms (if a tunable- $\Phi_M$  gate material and metallic S/D are used) and line-edge roughness effects. Therefore, they will likely be used to scale  $L_g$  to below 10 nm. Techniques for achieving uniformly ultrathin (<5 nm) channel films with atomically smooth surfaces will be needed to extend transistor scaling further, toward 1 nm  $L_g$ . Precise control of interfacial properties (not only at Si-dielectric interfaces, but also at metal-dielectric interfaces and Si-metal interfaces) will be critical for achieving high performance with good uniformity in nanometer-scale transistors. The need has never been greater for innovations in process technology, materials, and device design to sustain the Si revolution. ■

## Acknowledgments

I have learned much from my colleagues at UC Berkeley, especially Chenming Hu, Jeffrey Bokor, Bora Nikolić, Hideki Takeuchi, Leland Chang, Yang-Kyu Choi, Jakub Kedzierski, Pushkar Ranade, Min She, Sriram Balasubramanian, Daewon Ha, Hei Kam, Kyoungsub Shin, and Hiu Yung Wong. Research funding from DARPA

(Advanced Microelectronics Program), SRC (Advanced Devices and Technology Program), the MARCO Focus Center for Advanced Materials, Structures, and Devices, and the MARCO Focus Center for Circuits, Systems, and Software is gratefully acknowledged.

## References

- D. J. Frank, *et al.*, *Proc. IEEE*, **89**, 259 (2001).
- Semiconductor Industry Association, *International Technology Roadmap for Semiconductors*, 2003 Edition, <http://public.itrs.net/Files/2003ITRS/Home2003.htm>.
- A. Hokazono, *et al.*, *Tech. Dig. - Int. Electron Devices Meet.*, p. 639 (2002).
- D. Antoniadis, *Technical Digest-Very Large Scale Integration (VLSI) Symposium, IEEE*, p. 2 (2002).
- Y.-K. Choi, *et al.*, *IEEE Electron Device Lett.*, **21**, 254 (2000).
- S. Tang, *et al.*, *Tech. Dig. - IEEE Int. Solid-State Circuits Conf.*, p. 118 (2000).
- L. Chang, *et al.*, *Tech. Dig. - Int. Electron Devices Meet.*, p. 719 (2000).
- Y.-K. Choi, *et al.*, *IEEE 58th Device Research Conf.*, p. 85 (2001).
- K. Uchida, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 47 (2002).
- M. Jurczak, *IEEE Trans. Electron Devices*, **47**, 2179 (2000).
- Lindert, *et al.*, *IEEE Electron Device Lett.*, **22**, 487 (2001).
- J. Kedzierski, *et al.*, *Tech. Dig. - Int. Electron Devices Meet.*, p. 437 (2001).
- F.-L. Yang, *et al.*, *Tech. Dig. - VLSI Symp.*, p. 104 (2002).
- B. Yu, *et al.*, *Tech. Dig. - Int. Electron Devices Meet.*, p. 251 (2002).
- F.-L. Yang, *et al.*, *Tech. Dig. - VLSI Symp.*, p. 196 (2004).
- A. R. Brown, *et al.*, *IEEE Trans. Nanotechnology*, **1**, 195 (2002).
- Y.-K. Choi, *et al.*, *IEEE Trans. Electron Devices*, **49**, 436 (2002).
- S. Balasubramanian, *et al.*, *Proc. Silicon Nanoelectronics Workshop*, p. 16 (2003).
- H. Kam, *et al.*, *Proc. Silicon Nanoelectronics Workshop*, p. 9 (2004).
- N. Lindert, *et al.*, *IEEE Int. SOI Conf.*, p. 111 (2001).
- S. Zimin, *et al.*, *Proc. 5th Int. Conf. Solid-State and Integrated Circuit Technology*, p. 188 (1998).
- P. Xuan, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 609 (2003).
- K. H. Yuen, *et al.*, *IEEE Electron Device Lett.*, **24**, 518 (2003).
- L. Mathew, *et al.*, *IEEE Int. SOI Conf.*, p. 109 (2003).
- L. Mathew, *et al.*, *IEEE Int. Conf. on Integrated Circuit Design and Technology*, p. 97 (2004).
- S. Balasubramanian, *et al.*, *IEEE Int. SOI Conf.*, p. 27 (2004).
- B. Doyle, *IEEE Electron Device Lett.*, **16**, 301 (1995).
- Y.-C. King, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 585 (1998).
- M. Togo, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 347 (1998).
- J. Wang, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 707 (2002).
- D. M. Brown, *et al.*, *Solid-State Electron.*, **11**, 1105 (1968).
- P. Ranade, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 363 (2002).
- J. D. McBrayer, R. M. Swanson, and T. W. Sigmon, *J. Electrochem. Soc.*, **133**, 1242 (1986).
- T. Amazawa, *et al.*, *J. Electrochem. Soc.*, **145**, 1297 (1998).
- T.-J. King, *Semiconductor Interface Specialists Conf.*, p. S 2.1 (2002).
- K. Shin, *et al.*, in preparation.
- M. Kakumu, *et al.*, *Tech. Dig. - VLSI Symp.*, p. 30 (1984).
- M. Qin, *et al.*, *J. Electrochem. Soc.*, **148**, G271 (2001).
- J. Kedzierski, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 13.3.1 (2003).
- J. Kedzierski, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 247 (2002).
- Z. Krivokapic, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 445 (2003).
- Y. Yeo, *et al.*, *Tech. Dig. - VLSI Symp.*, p. 49 (2001).
- J. Schaeffer, *et al.*, Paper D4.1, presented at the 2004 MRS Spring Meeting, San Francisco, CA, April 12-16, 2004.
- D. Ha, *et al.*, *Tech. Dig. - Int. Electron Device Meet.* (2004).
- C. S. Park, *et al.*, *IEEE Electron Device Lett.*, **25**, 372 (2004).
- F. Ootsuka, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 575 (2000).
- K. Ota, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 27 (2002).
- T. Ghani, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 978 (2003).
- C.-H. Ge, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 73 (2003).
- M. Yang, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 453 (2003).
- J. Kedzierski, *et al.*, *Tech. Dig. - Int. Electron Device Meet.*, p. 57 (2000).
- D. Connelly, *et al.*, *IEEE Trans. Electron Devices*, **50**, 411 (2003).
- M. Tao, *et al.*, *Appl. Phys. Lett.*, **82**, 2593 (2003).
- A. Yagishita, *et al.*, Extended Abstracts of the International Conference on Solid-State Devices and Materials, p. 708 (2003).
- A. Kinoshita, *et al.*, *Tech. Dig. - VLSI Symp.*, p. 168 (2004).
- N. Jin, *et al.*, *IEEE Trans. Electron Devices*, **50**, 1876 (2003).
- L. Esaki, *Phys. Rev.*, **109**, 603 (1958).
- P. Mazumder, *et al.*, *Proc. IEEE*, **86**, 664 (1998).
- J. P. A. van der Wagt, *Proc. IEEE*, **87**, 571 (1999).

## About the Author

Tsu-Jae King is a professor of Electrical Engineering and Computer Sciences at the University of California, Berkeley, where her research interests include nanoscale semiconductor devices and technology. She is presently on industrial leave of absence at Synopsys, Inc. (Mountain View, CA) leading the development of NDR device technology for compact, low-power memory applications. She can be reached at [tking@eecs.berkeley.edu](mailto:tking@eecs.berkeley.edu) or [tsujae@synopsys.com](mailto:tsujae@synopsys.com).