# Novel Decoupling Capacitor Designs for sub- 90nm CMOS Technology

Xiongfei Meng, Karim Arabi†, and Resve Saleh

*SoC Research Laboratory, Department of Electrical and Computer Engineering,*
*University of British Columbia, 2356 Main Mall, Vancouver, BC, V6T 1Z4, Canada*
*†PMC-Sierra, Inc. 100-2700 Production Way, Burnaby, BC, V5A 4X1, Canada*
*E-mail: xmeng@ece.ubc.ca, Karim_Arabi@pmc-sierra.com, res@ece.ubc.ca*

## Abstract

*On-chip decoupling capacitors are generally used to reduce power supply noise. Traditional decoupling capacitor designs using NMOS devices may no longer be suitable for 90nm CMOS technology due to increased concerns on thin-oxide gate leakage and electrostatic discharge (ESD) reliability. A cross coupled design for standard cells has recently been proposed to address the ESD issue. In this paper, three modifications of the cross coupled design are introduced and the tradeoffs among ESD performance, transient response and gate leakage are analyzed. As shown here, the modifications offer designers greater flexibility in decoupling capacitor design for 90nm and below.*

## 1. Introduction

With increasing clock frequency and decreasing supply voltage as technology scales, maintaining the quality of power supply becomes a critical issue [1]. Typically, decoupling capacitors (decaps) are used to keep the power supply within a certain percentage (e.g., 10%) of the nominal supply voltage [2]. Decaps hold a reservoir of charge and are placed close to the power pads and near any large drivers. When large drivers switch, the decaps provide instantaneous current to the drivers to reduce IR drop and Ldi/dt effects [3], and hence keep the supply voltage relatively constant. A standard decap is usually made from NMOS transistors in a CMOS process [3].

At the 90nm technology node, the oxide thickness of a transistor is reduced to roughly 2.0nm. The thin oxide causes two new problems: possible electrostatic discharge (ESD) induced oxide breakdown and gate tunneling leakage [2][4]. Potential ESD oxide breakdown increases the likelihood that an integrated circuit (IC) will be permanently damaged during an ESD event, and hence raises a reliability concern. Higher gate tunneling leakage increases the total static power consumption of the chip. As technology scales further down, with a thinner oxide, the result is an even higher ESD risk and more gate leakage. The standard decap design experiences these two problems and therefore becomes rather inappropriate for 90nm and below.

A new cross coupled standard-cell design approach [5] addresses the issue of ESD performance. The design provides certain ESD input protection to the decap, but does not offer any savings in gate leakage. This paper suggests three modifications to the cross coupled cells that trade off ESD, transient response and leakage. Different modifications improve different design aspects, with certain drawbacks on others. The modifications are made to be suitable for different technology nodes and processes, and the overall effect is to provide designers with decap design flexibility for 90nm and below.

The rest of the paper is organized as follows. In Section 2, the necessary background on decap design, ESD problem, and gate leakage is briefly discussed. The cross coupled design is then analyzed and verified by SPICE simulations in Section 3. Modifications are proposed and compared with the standard and the cross coupled designs in Section 4. Section 5 suggests future directions and provides conclusions.

## 2. Background

A standard decap is usually implemented using an NMOS transistor with the gate connected to $V_{DD}$ and both source and drain connected to $V_{SS}$ (see Figure 1), or a PMOS device with opposite connections. The effective capacitance at low frequencies can be written as:

$$C_{eff} = C_{OX}WL + C_{OL}W \qquad (1)$$

where $C_{OX}$ is the oxide capacitance per unit area, $C_{OL}$ is the overlap and fringing capacitance per unit width, and W and L are the width and length of the transistor, respectively [3].

A standard decap also exhibits parasitic resistance of the channel that imposes certain delay on the transient response of the decap [6]. The decap's effective resistance at low frequencies is given by:

$$R_{eff} = \frac{L}{6\mu C_{OX}W(V_{GS} - V_T)} \qquad (2)$$

where $\mu$ is the mobility, $V_{GS}$ (or $V_{GD}$ since source and drain are tied) is the voltage across the oxide, and $V_T$ is the threshold voltage. From (2), $R_{eff}$ is proportional to the channel length L. That is, for faster transient response, a decap design should maintain L in a reasonably small range to keep $R_{eff}$ small. To capture the transient behavior, a decap can be modeled as a lumped RC circuit, as shown in Figure 1. Both $R_{eff}$ and $C_{eff}$ can be considered

constant at low or moderate operating frequencies, but they are degraded at high frequencies [6].
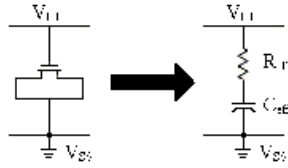


**Figure 1. Decap modelling as a series RC circuit**

Decaps can be used in the open areas of the chip between blocks (called "white-space" decaps) or inside the blocks composed of standard cells. If placed inside standard-cell arrays, it is more convenient to make decaps using both types of NMOS and PMOS to form a decap cell, knowing that the n-well is already implemented (Figure 2(a)) [6]. The overall impedance of two parallel RC circuits is determined as

$$(R_{eff\_n} + \frac{1}{sC_{eff\_n}}) // (R_{eff\_p} + \frac{1}{sC_{eff\_p}})$$ , and simplified as

$$\frac{R_{eff\_n}R_{eff\_p}}{R_{eff\_n} + R_{eff\_p}} + \frac{1}{s(C_{eff\_n} + C_{eff\_p})} + higher\_order\_terms$$

For first-order hand calculations, the higher-order terms are negligible. Thus, the overall effective capacitance is the sum of the two individual decoupling capacitances, and the overall effective resistance is the parallel combination of the two individual effective resistances. That is:

$$C_{eff\_overall} \approx C_{eff\_n} // C_{eff\_p} = C_{eff\_n} + C_{eff\_p} \qquad (3)$$

$$R_{eff\_overall} \approx R_{eff\_n} // R_{eff\_p} = \frac{R_{eff\_n}R_{eff\_p}}{R_{eff\_n} + R_{eff\_p}} \qquad (4)$$

Figure 2(b) illustrates a sample layout of this N+P configuration that uses two fingers to cut the overall $R_{eff}$ in half so that its transient response is maintained.
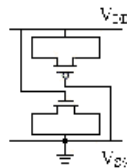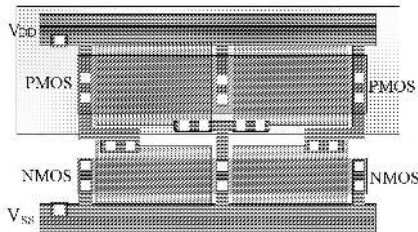


**Figure 2(a). Standard N+P decap configuration**



**Figure 2(b). A sample layout of standard N+P decap with two fingers**

## 2.1 Thin-oxide Gate Tunneling Leakage

A new design issue for decaps due to oxide thickness reduction is the gate tunneling current. The current is in the form of tunneling electrons or holes from substrate to gate or from gate to substrate through the gate oxide, depending on the voltage biasing conditions [7]. Two forms of gate tunneling exist: Fowler–Nordheim (FN) tunneling and direct tunneling. For normal operations on short-channel devices, FN tunneling is negligible, and direct tunneling is dominant [7]. In the case of direct tunneling, the gate leakage current in PMOS is much less than in NMOS, and it has been shown experimentally that PMOS gate leakage is roughly 3 times smaller than NMOS gate leakage for same size transistors [8][9]. The gate leakage simulations can be carried out by using BSIM4 SPICE models [10]. Assuming a 90nm technology with 2.0nm oxide thickness and 1.0V power supply, the gate leakage current is shown in Figure 3.
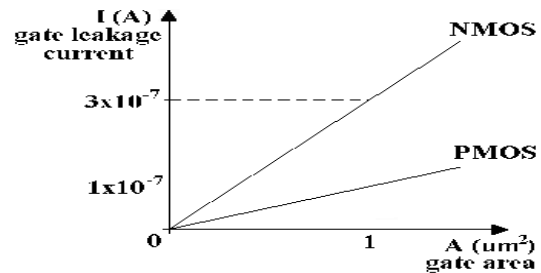


**Figure 3. Gate leakage current *vs.* gate area**

The gate leakage current density J and the oxide thickness $t_{OX}$ have an empirical relationship as follows if assuming the voltage across the oxide $V_{OX}$ is fixed [8]:

$$J = 10^{(A - B \cdot t_{OX})} \qquad (5)$$

where A and B are experimental constants and are process dependent. Equation (5) implies that the gate leakage current is exponentially related to the oxide thickness. A typical J and $t_{OX}$ relationship for a fixed $V_{OX}$ is illustrated in Figure 4.
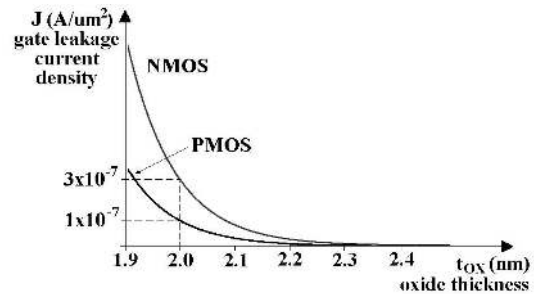


**Figure 4. Gate leakage current density *vs.* oxide thickness**

It is evident that from 90nm technology, the gate leakage from decaps is already significant [9]. The gate leakage contributes to the total static power consumption, and decaps usually occupy a large on-chip area. The use

of PMOS devices exclusively is not a viable solution for high-frequency circuits since they have a poor frequency response relative to the NMOS devices [6].

## 2.2 ESD Reliability in Decap Design

Another new consideration has arisen in the form of ESD protection due to the thin oxide in 90nm technology. ESD is the process of static discharge that can typically arise from human contact with any IC pin. Approximately 0.6uC of charge is carried on a body capacitance of 100pF, generating a potential of 2KV or higher to discharge from the contacted IC pin to ground for a duration of more than 100ns [4]. Under such an event, the peak discharge current is in the ampere range, leading to permanent damage on certain transistors in the chip if not properly protected. The damage can be in one of the two forms or the combination of the two: one is thermal burnout in devices or interconnects, while the other is oxide breakdown of devices due to the high voltage across the oxide [4]. When running simulations for an ESD event, the maximum current density $J_{max}$ of devices and interconnects is measured to check for potential thermal damage. The oxide voltage also needs to be measured to compare with the oxide breakdown voltage of a device for a given fabrication process. The oxide breakdown voltage is almost linearly proportional to the oxide thickness [4]. For instance, assuming that a 90nm process uses 2.0nm of oxide thickness, the corresponding oxide breakdown voltage is below 5V. If the thickness is doubled, the oxide breakdown voltage is also doubled [4].

An ESD event can be delivered between any two pins of an IC. To properly protect an IC from ESD damage, an ESD circuit must shunt ESD current between these two pins [4]. In the case of decaps within standard cells, the only two pins that the decaps have access to are the two local power rails, namely $V_{DD}$ and $V_{SS}$. Primary and local (sometimes called "secondary") protection elements are needed to protect the two rails by limiting the voltage difference between the two rails to a value below the oxide breakdown voltage. The primary element will shunt most of the ESD current, whereas the local element serves to limit the voltage or current at the local circuit until the primary element is fully operational [4]. A primary element can be a thick oxide transistor, a silicon controlled rectifier, an open-gate, grounded-gate or coupled-gate NMOS transistor, or a large ESD diode [4]. A local protection element can be simply a diode formed by a grounded-gate NMOS transistor [4].

The complete ESD protection scheme is illustrated in Figure 5. In addition to the primary and local elements, a resistor $R_{in}$ is required to limit the maximum current flow to the decap and to limit the voltage seen from the gate of the decap. For better ESD protection, this resistance is normally large and can be in the forms of polysilicon, diffusion, n-well, or even channel resistance [4]. The resistance is generally not implemented together with primary and local protection devices. It is likely to be inserted within standard cells where ESD damage is a concern.
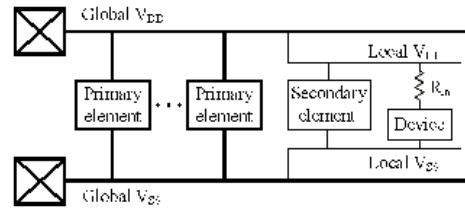


**Figure 5. Complete ESD protection scheme**

Previous decap designs (before 90nm technology) did not consider ESD performance mainly because: 1. The transistor's oxide thickness is thick and the oxide breakdown voltage is high enough so that the transistor is likely to survive during an ESD event with adequate protection circuits. 2. Insertion of the large resistance $R_{in}$ dramatically reduces the transient response of the decap. However, starting from 90nm, the oxide thickness is so thin that the designer cannot ignore the increased ESD risk. A large resistance is therefore recommended to be placed inside the decap cells to protect from potential ESD damage. As a consequence, this tradeoff between ESD performance and transient response becomes the main decap design challenge in 90nm.

## 3. Cross Coupled Decap Design

Knowing that the standard N+P decap design may no longer be suitable for 90nm technology due to the increased ESD risk, a new cross coupled decap design has been proposed in [5] to address the issue of ESD reliability. The new cross coupled design (Figure 6) reconnects the terminals of the two transistors: the drain of the PMOS connects to the gate of the NMOS, whereas the drain of the NMOS is tied to the gate of the PMOS [5].
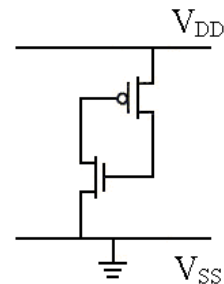


**Figure 6. Cross coupled decap schematic [5]**

The design can be modeled as a series connection of $R_{eff}$ and $C_{eff}$ at low frequencies, similar to the standard decap. The overall $C_{eff}$ is roughly the same, while the overall $R_{eff}$ increases significantly. Both transistors are still in the linear region, but the channel resistance is modified. Specifically,

$$C_{eff\_overall} \approx C_{eff\_n} // C_{eff\_p} = C_{eff\_n} + C_{eff\_p} \qquad (6)$$

$$R_{eff\_overall} \approx (R_{eff\_p} + R_{on\_n}) // (R_{eff\_n} + R_{on\_p})$$

$$\approx R_{on\_n} // R_{on\_p} = \frac{R_{on\_n} R_{on\_p}}{R_{on\_n} + R_{on\_p}} \quad (7)$$

where $C_{eff\_n}$, $C_{eff\_p}$, $R_{eff\_n}$ and $R_{eff\_p}$ are the intrinsic effective capacitances and resistances, respectively, and $R_{on\_p}$ and $R_{on\_n}$ are the channel resistance of the two transistors. Since $R_{on\_p}$ and $R_{on\_n}$ are at least one order of magnitude larger, $R_{eff\_p}$ and $R_{eff\_n}$ can be rounded off from the overall $R_{eff}$ calculation. $C_{eff\_n}$, $C_{eff\_p}$, $R_{eff\_n}$ and $R_{eff\_p}$ can be obtained from (1) and (2). Both $R_{on\_p}$ and $R_{on\_n}$ can be computed as follows [3]:

$$R_{on} \approx R_{eq} \frac{L}{W} \quad (8)$$

where $R_{eq}$ is the process-dependent square resistance (k$\Omega$). It is important to realize that (6) - (8) are first-order, low-frequency approximations only. The real transistor channel resistance by nature is nonlinear and depends strongly on applied voltages, operating frequency, and geometry [3]. The only reason for providing these formulae is to give designers useful information when making design tradeoffs.

This cross coupled design improves the ESD performance of the decap by making the overall effective resistance larger without adding additional area. The tradeoff of the design is a reduced transient response. The larger $R_{eff}$ corresponds to a longer RC delay. Assuming that the cell area is fixed and that only the terminal connections are swapped, this design provides no savings in gate leakage as compared to the standard decap.

To quantitatively measure ESD performance, RC delay in transient response, and gate leakage, a number of simulations were carried out. The layouts were created in Virtuoso™ Layout Editor, verified by Calibre™ DRC checker, and then extracted by Calibre™ XRC parasitic extraction tool. The extracted data were simulated with HSPICE™ for different simulation setups. For fairness, the same cell area was used for all the designs.

The ESD simulation requires an ESD generation model. Among all the existing models, the human body model (HBM) was adopted for simplicity. Following the standard MIL-STD-883x method 3015.7, a human body can be simulated as a series of 1.5K$\Omega$ resistance $R_{HBM}$ and 100pF capacitance $C_{HBM}$. The capacitor $C_{HBM}$ is initially charged to 2KV that needs to be discharged through some primary elements [4]. The primary element is arbitrarily chosen to be an ESD diode plus a gate-coupled NMOS device (GCNMOS) with an n-well resistor $R_{nwell}$ (~15K$\Omega$) and an NMOS bootstrap capacitor $C_b$ [4]. Two identical primary elements are used to protect the circuit placed in between the HBM generation and the elements, as shown in Figure 7 [4]. For simplicity, no secondary element is used.
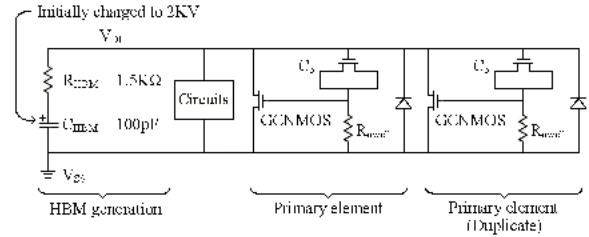


Figure 7. Simulation setup for ESD analysis [4]

Since the primary elements are designed to handle large current flow, the maximum current density $J_{max}$ is assumed to be within the safe range and is not measured again for simplicity. HBM generation raises the voltage level at node $V_{DD}$, and hence turns on the primary elements to discharge. For device protection from oxide breakdown, the voltage differences across gate and source ($V_{GS}$) and across gate and drain ($V_{GD}$) of the two transistors are simulated. The $V_{GS}$ and $V_{GD}$ voltages are desired to be kept as low as possible, knowing that the oxide breakdown voltage for a typical 90nm process is below 5V. From simulation measurements, for the standard decap design, $V_{GD\_p} = V_{GS\_p} = V_{GD\_n} = V_{GS\_n} = 4.2V$. For the cross coupled design, $V_{GD\_p} = 4.0V$, $V_{GS\_n} = 3.2V$, and $V_{GS\_p} = V_{GD\_n} = 3.0V$. Clearly, the cross coupled design provides better ESD performance.

For transient response, to demonstrate how the decaps perform, a normal operating condition is set: the $V_{DD}$ node is connected to the nominal supply of 1V (for 90nm), and $V_{SS}$ is tied to a common ground. When there is no activity the current flow from $V_{DD}$ to $V_{SS}$ is solely treated as gate leakage. At 1ns, $V_{DD}$ starts to drop linearly from 1V to 0.9V, reaches 0.9V at 2ns, and then remains constant. By definition, an ideal capacitor responds to a voltage change as a current source if it is fully charged, as follows:

$$I = C_{ideal} \frac{dV}{dt} \approx C_{ideal} \frac{\Delta V}{\Delta t} \quad (9)$$

If the voltage change is a ramp, the current provided by the ideal capacitor should be a pulse. In practice, due to the presence of the effective resistance associated with the decap designs, a certain amount of RC delay exists. Good transient response should have sharp rise and fall edges (at 1ns and 2ns in this case), while it can also provide a large average current $I_{avg}$ during the time period from 1ns to 2ns. The sharpness of rise and fall is measured from the rise/fall slopes with a unit of A/s. The average capacitance $C_{avg}$ is calculated from $I_{avg}$ from (9). Figure 8 illustrates the curves for the two designs in transient analysis, and indicates that the standard decap can provide much better transient response. The two designs also have almost identical gate leakage: 53.8nA for the standard decap and 53.7nA for the cross coupled design.
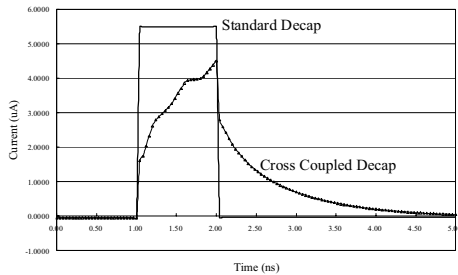
Figure 8. Simulation results for RC delay

## 4. Modified Cross Coupled Decap Designs

Three modifications are made to address different goals of decap design: ESD performance, transient response, and gate leakage. It is difficult to simultaneously make improvements on all the three goals, but trying to balance them and to make tradeoffs is certainly feasible and indeed achievable. Each modification is compared to the basic cross coupled design to show advantages and disadvantages. Again, the total cell area is fixed for all the designs.

The first modification (Mod1) attempts to improve ESD performance by making the channel lengths of the two resistors longer (Figure 9). The two fingers are combined into one. As a result, the overall $R_{eff}$ is almost doubled, while the overall $C_{eff}$ remains roughly the same. The disadvantage of this design is reduced transient response and slightly larger gate leakage since the gate area increases a little.

The second modification (Mod2) attempts to reduce gate leakage while maintaining ESD performance and transient response at the same level (Figure 10). One NMOS is replaced by a PMOS with the n-well expanded to accommodate the new PMOS. The effect of this change is then increased $R_{on\_p}$ and $C_{eff\_p}$. To match ESD performance, $R_{on\_n}$ needs to be reduced. One simple change to obtain a small $R_{on\_n}$ is to reduce the channel length of the NMOS. By the same token, $C_{eff\_n}$ is also reduced. The result is comparable ESD performance and transient response if carefully designed. Knowing that the new same-area PMOS leaks 3 times less than the replaced NMOS, extra saving in gate leakage is realized.

The third modification (Mod3) (Figure 11) follows the similar approach as of Mod2. It further increases the new PMOS area and further reduces the NMOS channel length. Indeed, the minimum length NMOS is used to have the smallest possible $R_{on\_n}$ so that it dominates and makes the overall $R_{eff}$ smaller. Since the overall $R_{eff}$ is greatly decreased while the overall $C_{eff}$ is somewhat increased, the transient response dramatically improves. An even larger PMOS and smaller NMOS lead to additional savings in gate leakage as well. The only downside is reduced ESD protection capability due to the reduced overall $R_{eff}$.
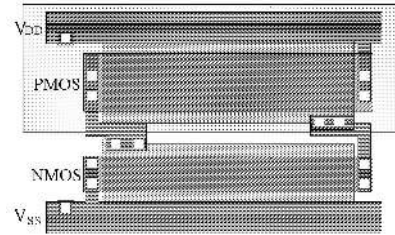


Figure 9. A sample layout of Mod 1 (Basic circuit without fingering)
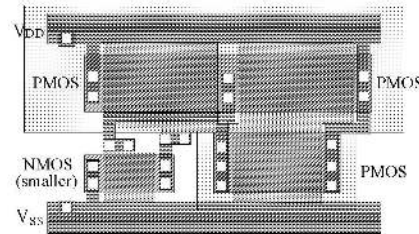


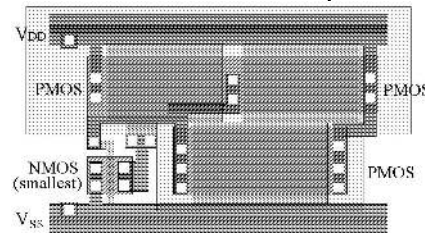Figure 10. A sample layout of Mod 2 (Replace NMOS with PMOS)



Figure 11. A sample layout of Mod 3 (Replace NMOS with PMOS, and use smallest NMOS)

Following the same simulation procedures outlined earlier, Table I lists the comparison for all the designs on ESD performance, transient response, and gate leakage. The **bold** numbers indicate the best results in the comparison. One can see clearly the improvements and the tradeoffs from the simulation results.

### Table I. Comparison on ESD performance, transient response and gate leakage

| | ESD performance with 2 primary elements | | | Transient response | | Gate leak age (nA) |
|---|---|---|---|---|---|---|
| | $V_{GD\_p}$ (V) | $V_{GS\_n}$ (V) | $V_{GS\_p} = V_{GD\_n}$ (V) | Rise slope (A/s) | Avg cap (fF) | |
| Decap | 4.2 | 4.2 | 4.2 | **2.8e5** | **54.3** | 53.8 |
| Cross Coupled | 4.0 | 3.2 | 3.0 | 8.2e4 | 33.1 | 53.7 |
| Mod1 | **3.8** | **2.9** | **2.8** | 8.7e4 | 21.4 | 59.7 |
| Mod2 | 4.0 | 3.7 | 3.4 | 7.0e4 | 35.8 | 33.6 |
| Mod3 | 4.1 | 3.9 | 3.8 | 1.1e5 | 47.5 | **31.8** |

There is no single design that suits all the possible situations. The reason for having several possible designs is to provide designers with different solutions so that

they can make a suitable choice for a specific process at a specific technology node. For 90nm technology, the standard decap still seems to be acceptable in ESD reliability, assuming the power rails have protection elements. However, Mod3 is more suitable because it has better ESD performance and saves roughly 41% on gate leakage. The only tradeoff then is a slightly reduced transient response. As technology further scales, or as a different process increases the transistor speed, the oxide thickness will probably become thinner and the oxide breakdown voltage will occur. Under that scenario, the standard design or the Mod3 will no longer be appropriate any more. For improved ESD performance, Mod2 will be suggested instead of the basic cross coupled design. The reason is that Mod2 has similar ESD numbers and similar transient response compared to the basic cross coupled design but saves approximately 40% on gate leakage. When technology scales down to a point that the oxide thickness makes the ESD reliability a more serious concern, the use of Mod1 will be advised for the best ESD performance, although its transient response will be sacrificed significantly.

The recommendations above are good for moderate or low frequency chips. If the targeting frequency is extremely high, even Mod3 may not be able to provide desired amount of current within an excessively small period of time. Under such case, the use of thick-oxide decaps is suggested around the standard-cell blocks. For 90nm technology, the thick oxide is 3x thicker than the thin oxide, resulting in almost zero gate leakage and 3x ESD breakdown voltage. The disadvantage is the effective capacitance reduced to 1/3. Hence, the area needed for a fixed capacitance is 3x for thick-oxide decaps. The thick-oxide decaps must be properly placed around the periphery of the block. The fabrication cost for using thick-oxide devices may also be slightly higher.

## 5. Conclusions and Future Directions

Decaps are commonly used to maintain the power supply voltage within certain noise margins. For 90nm technology, a cross coupled design is suggested to replace the standard decap for better ESD performance, but it suffers from reduced transient response. Moreover, the cross coupled design provides no savings in gate leakage that begins to account for a significant portion of the static power consumption in 90nm CMOS. This paper proposed three modifications of the basic cross coupled design to make tradeoffs among ESD performance, transient response, and gate leakage. Among the three, Mod2 is designed to replace the cross coupled design; Mod1 has the best ESD performance; Mod3 provides better transient response and the least gate leakage. The designer is given the opportunity and the flexibility to choose the most suitable design for any specific situation.

As technology further scales to 65nm or below, the ultra thin oxide will increase the ESD risk and the amount of gate leakage, and will eventually limit the use of the cross coupled design and its modifications. The anticipation at this stage would be the use of high-k dielectrics as the oxide material so that the electrical thickness and the physical thickness can be differentiated to completely eliminate the concerns on ESD reliability and gate leakage. Another approach would be to utilize high-voltage thick-oxide transistors as decaps, as discussed earlier. In any case, solutions that properly balance ESD, gate leakage, transient response and area will be required.

## 6. Acknowledgments

## 7. References

[1] H. H. Chen and S. E. Schuster, "On-chip decoupling capacitor optimization for high-performance VLSI design", in *Proc. of International Symp. on VLSI Technology, Systems, and Applications*, 1995, pp. 99-103.

[2] H. H. Chen, J. S. Neely, M. F. Wang, and G. Co, "On-chip decoupling capacitor optimization for noise and leakage reduction", in *Proc. of Symp. on Integrated Circuits and Systems Design*, 2003, pp. 319-326.

[3] D. A. Hodges, H. G. Jackson, and R. A. Saleh, *Analysis and Design of Digital Integrated Circuits in Deep Submicron Technology*, 3rd Ed, McGraw-Hill, 2004.

[4] A. Amerasekera and C. Duvvury, *ESD in Silicon Integrated Circuits*, 2nd Ed, John Wiley & Sons, 2002.

[5] *TSMC 90nm CLN90G Process SAGE-X v3.0 Standard Cell Library Databook*, Release 1.0, Artisan Components Inc., 2004.

[6] J. Chia, "Design, Layout and Placement of on-chip decoupling capacitors in IP blocks", *M.A.Sc Thesis*, the University of British Columbia, 2004.

[7] K. Roy, S. Mukhopadhyay, and H. Mahmoodi-Meimand, "Leakage Current Mechanisms and Leakage Reduction Techniques in Deep-Submicrometer CMOS Circuits," in *Proc. IEEE*, vol. 91, no. 2, pp. 305-327, Feb. 2003.

[8] W. C. Lee and C. Hu, "Modeling Gate and Substrate Currents due to Conduction- and Valence-Band Electron and Hole Tunneling," in *Dig. Tech. Papers Symp. VLSI Technology*, 2000, pp. 198-199.

[9] F. Hamzaoglu and M. Stan, "Circuit-level techniques to control gate leakage for sub-100 nm CMOS," in *Proc. Int. Symp. Low Power Design*, 2002, pp. 60–63.

[10] K. Cao, W. -C. Lee, W. Liu, X. Jin, P. Su, S. K. H. Fung, J. X. An, B. Yu, and C. Hu, "BSIM4 gate leakage model including source drain partition," in *Tech. Dig. IEDM*, 2000, pp. 815-818.