# Design Techniques to Facilitate Processor Power Delivery in 3-D Processor-DRAM Integrated Systems

Qi Wu and Tong Zhang, *Senior Member, IEEE*

*Abstract*—As a promising option to address the memory wall problem, 3-D processor-DRAM integration has recently received many attentions. Since DRAM dies should be stacked between the processor die and package substrate, we have to fabricate a large number of through-DRAM through-silicon vias (TSVs) to connect the processor die and package for power and input/output (I/O) signal delivery. Although such through-DRAM TSVs will inevitably interfere with DRAM design and induce non-negligible power consumption overhead, little prior research has been done to study how to allocate these through-DRAM TSVs on the DRAM dies and analyze their impacts. To address this open issue, this paper first presents a through-DRAM TSV allocation strategy that well fits to the regular DRAM architecture. Meanwhile, due to the longer path between power/ground pads and processor die, power delivery integrity issue may become more serious in such 3-D processor-DRAM integrated systems. Decoupling capacitor insertion is the most popular method to deal with power delivery integrity issue in high-performance integrated circuits. This paper further proposes to use 3-D stacked DRAM dies to provide decoupling capacitors for the processor die. This can well leverage the superior capacitor fabrication ability of DRAM to reduce the area penalty of decoupling capacitor insertion on the processor die. For its practical implementation, a simple uniform decoupling capacitor network design strategy is presented. To demonstrate through-DRAM TSV allocation and decoupling capacitor insertion strategy and evaluate involved tradeoffs, circuit SPICE simulations and computer system simulations are carried out to quantitatively demonstrate the effectiveness and investigate various design tradeoffs.

*Index Terms*—3-D integration, decoupling capacitor, instructions per cycle (IPC), IR drop, power delivery, processor-DRAM integrated system, through-silicon via (TSV).

## I. INTRODUCTION

AS THE computing industry enters the multi-core era, the looming memory wall problem [1] is becoming an increasingly severe issue. Continuous technology scaling can certainly integrate more SRAM and/or embedded DRAM on the processor die, but it can hardly provide enough on-chip memory capacity. Although embedded DRAM can achieve a higher storage density than SRAM, its storage density is still at least 2× lower than commodity DRAM and it may noticeably increase the processor die fabrication cost. As a promising alternative to technology scaling, the emerging 3-D integration technologies provide a viable solution to address this problem,

i.e., by stacking multiple high-capacity DRAM dies and one high-performance processor die together with massive short through-silicon vias (TSVs), 3-D processor-DRAM integrated systems can achieve drastically reduced memory access latency and increased memory access bandwidth. This has been well recognized by the computer architecture community and many recent work [2]–[7] have explored and well demonstrated its very encouraging potential.

Due to its energy-hungry nature, the high-performance processor die tends to dissipate much more heat than the 3-D stacked DRAM dies. Therefore, in high-performance 3-D processor-DRAM integrated computing systems, the processor die must directly attach to the heat spreader and heat sink. Hence, the DRAM dies have to be stacked between the package substrate and the processor die. As a result, we must fabricate a certain number of TSVs that go through the stacked DRAM dies to deliver a large amount of current and all the input/output (I/O) signals from the package to the processor die. Clearly, those through-DRAM TSVs will inevitably affect the DRAM design and incur DRAM storage capacity degradation. Moreover, the through-DRAM power TSVs themselves may also incur non-negligible power consumption overhead, particularly as we stack many DRAM dies between the package substrate and processor. Meanwhile, the power delivery path for processor die in such 3-D processor-DRAM integration becomes longer, which may result in non-negligible on-chip IR drop. Furthermore, all the DRAM dies along the power delivery path may also possibly contribute noise to the power delivery path. However, all the prior studies on 3-D processor-DRAM integrated systems did not explicitly take into account of the above important and practical issues. Although it is of great importance to fully study the impact of through-DRAM TSVs on the system design, as recently pointed out by a keynote at ISSCC'09 [8], little work has been done in the open literature to the best of our knowledge.

This paper attempts to contribute to fill this missing link by presenting two design techniques and comprehensive case studies. First, we develop a through-DRAM power and signal TSVs allocation strategy that can well fit into the regular DRAM structure. In particular, since a large amount of power TSVs may be inevitable in order to reliably deliver sufficient current to the processor, we present a uniformly distributed power TSV network design approach, which makes the fabrication of through-DRAM power TSVs do not interfere with the DRAM design itself. Using this design strategy, designers can easily adjust the tradeoff between the power consumption overhead and DRAM storage capacity degradation induced by through-DRAM power TSVs. Second, we introduce a

DRAM-based decoupling capacitor insertion strategy to further improve the quality of through-DRAM power delivery for processor. As on-chip power supply voltage continues to reduce in order to make power consumption under control, high-performance processors are subject to increasingly significant run-time voltage fluctuations or variations, which tends to demand a heavier use of decoupling capacitor insertion along the power delivery path. The longer distance between the package substrate and processor in 3-D processor-DRAM integration can make this issue more critical. Motivated by the fact that DRAM process is highly optimized for fabricating high-capacitance capacitors with very small footprint, it is very intuitive that, in 3-D processor-DRAM integrated computing systems, DRAM dies can very effectively provide a large amount of decoupling capacitors for the processor die at very small area penalty. Following this intuition, this work further presents a simple uniform decoupling capacitor network design strategy to realize DRAM-based decoupling capacitors for the 3-D stacked processor die. By providing an abundant amount of decoupling capacitors uniformly across the entire processor die, this simple method does not require accurate circuit activity analysis for decoupling capacitor allocation. Meanwhile, it does not affect the structural regularity in DRAM design and may only incur very small DRAM storage capacity degradation.

We modified the widely used memory modeling tool CACTI [9] to quantitatively evaluate the above two design strategies. The developed modeling tool can estimate the power consumption overhead incurred by through-DRAM power TSVs, and the DRAM capacity degradation incurred by through-DRAM power TSVs and DRAM-based decoupling capacitors. We also carried out extensive circuit SPICE simulations to further evaluate the decoupling efficiency. The modeling and simulation consider a wide range of different system parameters, including the number of DRAM dies, through-DRAM power TSV diameters, and TSV-metal contact resistance. Moreover, using the M5 full system simulator [10] and a wide range of benchmarks, we evaluated the performance, in terms of instructions per cycle (IPC), of such 3-D processor-DRAM integrated computing systems to further demonstrate the impact and tradeoffs of the proposed design techniques.

The remainder of this paper is organized as follows. Section II reviews the basics of 3-D integration technology and prior work on 3-D processor-DRAM integration. Section III presents a through-DRAM signal and power TSVs allocation strategy, and Section IV presents power delivery integrity analysis when using the proposed through-DRAM TSVs allocation strategy. Section V presents a DRAM-based decoupling capacitor insertion scheme to further improve the power delivery integrity, and Section VI shows full system simulation results using the M5 simulator under various design parameters. Conclusions are drawn in Section VII.

## II. BACKGROUND AND PRIOR WORK

### A. 3-D Integration Technology

3-D integration refers to a variety of technologies which provide electrical connectivity between stacked multiple active device planes. Various 3-D integration technologies are currently pursued and can be divided into the following three categories [11].

1) *3-D packaging technology*: It is enabled by wire bonding, flip-chip bonding, and thinned die-to-die bonding [12]. Its major limitation is very low inter-die interconnect density (e.g., only few hundreds of inter-die bonding wires) compared to the other emerging 3-D integration technologies.

2) *Transistor build-up 3-D technology*: It forms transistors layer by layer, on poly-silicon films, or on single-crystal silicon films. Although a drastically high vertical interconnect density can be realized, it is not readily compatible to existing fabrication process and is subject to severe process temperature constraints that tend to degrade the circuit electrical performance.

3) *Monolithic, wafer-level, back-end-of-the-line (BEOL) compatible 3-D technology*: It is enabled by wafer alignment, bonding, thinning and inter-wafer interconnections [13]. Realized by TSVs[14], inter-die interconnects can have very high density. Wafer-level BEOL-compatible 3-D integration appears to be the most promising option for high-volume production of highly integrated systems. Therefore, this work assumes the use of wafer-level BEOL-compatible 3-D integration technology and hence the availability of TSVs for inter-die interconnect.

### B. 3-D Processor-DRAM Integration

Although DRAM can realize a very high storage density, its fabrication process is not compatible with that of high-performance logic dies. As a result, in current design practice, DRAM and processor always locate on separate chips connected with chip-to-chip links. The limited bandwidth and relatively long latency of chip-to-chip links greatly contribute to the looming memory wall problem. Emerging 3-D integration technologies make it possible to stack processor die and DRAM dies together that are linked through TSVs with massive inter-die communication bandwidth and very low latency. Very intuitively, such 3-D processor-DRAM integration appears to be a very natural solution to tackle the memory wall problem, hence it has attracted many attentions in computer architecture community and very promising results have been demonstrated (e.g., see [2]–[7]).

One of major issues in 3-D integration is the heat dissipation of stacked dies. Microprocessors are generally very energy-hungry and hence tend to generate a significant amount of heat, compared with which DRAM consume much less energy and generate much less heat. Therefore, stacking a single microprocessor die with multiple DRAM dies appears to be the most plausible option, at least in the foreseeable future, in which the microprocessor die locates closest to the heat sink. Such a 3-D integration configuration has been assumed in most prior work. Therefore, in this work, we only focus on the scenario where a single processor die is stacked with multiple DRAM dies.

## III. DRAM ARCHITECTURE DESIGN IN THE PRESENCE OF THROUGH-DRAM TSVs

Fig. 1 illustrates the overall architecture of a 3-D processor-DRAM integrated system. One side of the chip are heat sink, heat spreader and thermal interface material, and the other side
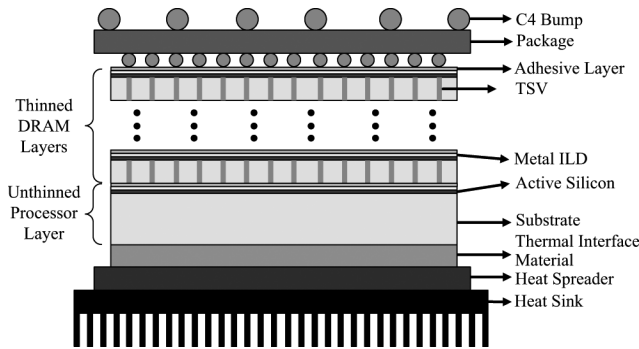
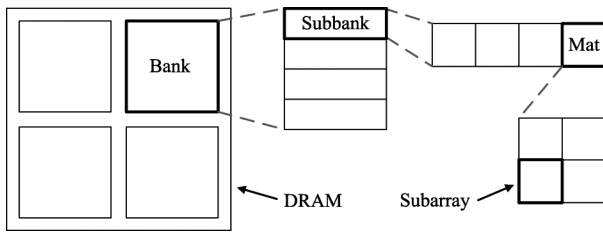Fig. 1. 3-D processor-DRAM integrated system architecture.



Fig. 2. Layout of an example DRAM array with four banks [9].



Fig. 3. Proposed design strategy to allocate through-DRAM TSVs on 3-D stacked DRAM dies.

of the chip are C4 bumps and package. Clearly, the un-thinned processor die must be directly attached to thermal interface material. All the thinned DRAM dies are stacked between the package and processor die. Besides TSVs connecting the processor die and DRAM dies, a large amount of TSVs are fabricated to connect power and signal I/O pins between the package and the processor die, which must go through all the stacked DRAM dies and are referred to as *through-DRAM* TSVs. It is obvious that through-DRAM TSVs will impact the 3-D DRAM design and degrade DRAM storage capacity. Nevertheless, how to design 3-D DRAM in the presence of through-DRAM TSVs has not been ever addressed in the open literature, and the impacts of through-DRAM TSVs on 3-D integrated systems have been largely ignored in the prior work. To address this open issue, this section presents a simple yet effective DRAM design method that can naturally accommodate those through-DRAM TSVs.

### A. Proposed Design Strategy

It is well known that, DRAM typically has a bank $\rightarrow$ sub-bank $\rightarrow$ sub-array hierarchical structure as shown in Fig. 2, i.e., one DRAM die consists of one or more banks that can be independently accessed; each bank is divided into sub-banks, and the data are read (written) from (to) one sub-bank during each memory access to one bank; each sub-bank is further divided into sub-arrays, and each sub-array contains an indivisible array of DRAM cells surrounded by supporting peripheral circuits such as word-line decoders and drivers, sense amplifiers (SAs), and output drivers, etc.

Fig. 3 illustrates the proposed strategy to allocate through-DRAM TSVs on DRAM dies and meanwhile maintain the regular DRAM architecture, where we consider the signal TSVs and power TSVs separately. The motivation is described as follows. Since each signal TSV may simply use the minimal allowable TSV size (e.g., a few $\mu$m of diameter) and microprocessors
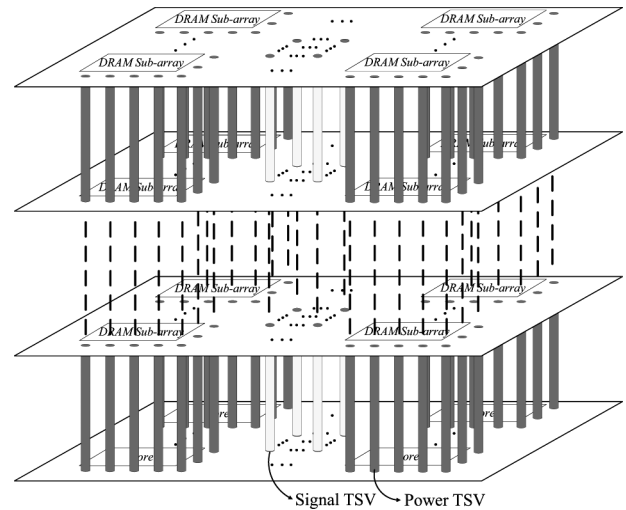
typically have a few hundred signal I/Os, through-DRAM signal TSVs tend to occupy a very small area on DRAM dies. Hence, we proposed to simply reserve a region at the center of DRAM dies dedicated for all the signal TSVs, as shown in Fig. 3. Meanwhile, we note that this strategy can readily decouple the design of the DRAM dies and the microprocessor die, i.e., we can allocate a large enough number of through-DRAM signal TSVs on DRAM dies so that the same 3-D DRAM can serve for different microprocessor dies with different amount of signal I/Os.

In comparison, through-DRAM power TSVs could have a much bigger impact and result in nontrivial system design tradeoffs. Energy-hungry microprocessor dies typically need tens or few hundreds of Ampere current, which may further increase as the supply voltage continues to reduce. To maintain a reasonable IR drop on the entire processor die (e.g., 10% or less of power supply voltage), highly distributed through-DRAM power TSVs must be used. Due to the non-negligible resistance of TSVs, through-DRAM power TSVs may incur noticeable power consumption overhead. To reduce such power consumption overhead, we need to reduce the aggregate resistance of through-DRAM power TSVs, hence we have to increase the aggregate size of those power TSVs. Moreover, as more DRAM dies are being stacked, the distance between the package and microprocessor die (i.e., the length of through-DRAM TSVs) will increase. In order to maintain the same aggregate power TSV resistance, we have to accordingly increase the size of each power TSV and/or fabricate more power TSVs.

To readily accommodate the fabrication of a large amount (e.g., thousands or tens of thousands) of through-DRAM power TSVs, we propose to arrange a regular power TSV mesh network around those individual DRAM sub-arrays, as illustrated in Fig. 4. We first partition the entire array of DRAM sub-arrays into a certain number of equal-size sub-array sets, then put power TSV channels between all the adjacent sub-array sets, where all the through-DRAM power TSVs uniformly distribute within all the channels, as shown in Fig. 4. We note that this proposed through-DRAM TSV allocation strategy does not change the existing DRAM wire profile and the number of wiring layers. We can consider that the DRAM layout is
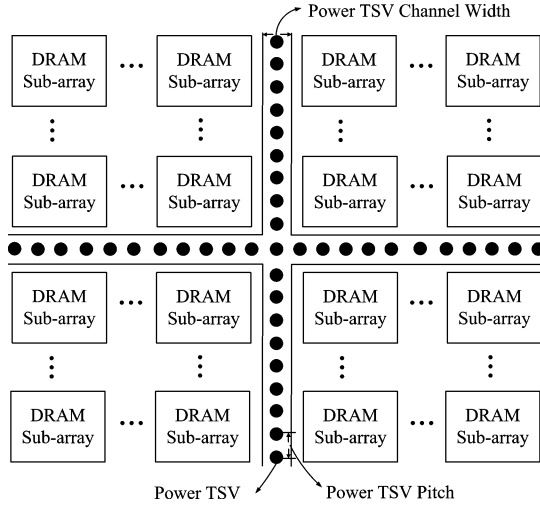
Fig. 4.   Illustration of through-DRAM power TSVs channels.

stretched to leave extra space between certain adjacent sub-arrays for through-DRAM power TSV channels. All the address and data routings in DRAM are either parallel or perpendicular to those TSV channels. Since we only consider very coarse-grained TSV pitch (e.g., tens of $\mu$m in case studies presented later), there should be enough space for address and data routings to readily go through TSV channels. By adjusting the pitch and width of power TSV channels, we could vary the aggregate size of through-DRAM power TSVs and hence tune the tradeoff between the through-DRAM power delivery quality and DRAM storage capacity degradation induced by power TSVs. Note that, although the power TSV channels as illustrated in Fig. 4 have the minimal width that can only accommodate one power TSV, a wider power TSV channel can be used if more power TSVs are required.

We note that these through-DRAM TSVs can also be used to delivery power to all the DRAM tiers. In order to isolate power delivery system for processor and DRAM and hence simplify power noise estimation, we propose to use separate through-DRAM TSVs to delivery power to processor and DRAM, e.g., 92% through-DRAM TSVs are for processor power delivery and 8% left for DRAM power delivery. Besides isolating power delivery system for processor and DRAM, this approach can also make it possible to reduce the TSV-metal contact number in processor power delivery path, which we will discuss at the end of Section III-B.

In order to quantitatively evaluate the performance of the above design strategy, we need to calculate the area and power consumption overhead incurred by through-DRAM TSVs based on parameters such as thickness of each DRAM die, TSV resistivity, TSV diameter, system power consumption, and power supply voltage. The thickness of one DRAM die, denoted as $T_t$, can be expressed as

$$T_t = T_a + T_m + T_{si} + T_{sub} \tag{1}$$

where $T_a$, $T_m$, $T_{si}$, and $T_{sub}$ represent the thickness of the adhesive layer, metal layers, active silicon layer, and substrate, respectively. The resistance of through-DRAM power TSV, denoted as $R_{tsv\_p}$, consists of two parts, including the intrinsic TSV resistance and contact resistance between TSVs and top

metal layer of DRAM and processor dies. The contact resistance can be obtained from

$$R_c = \frac{\rho_c}{A_c} \tag{2}$$

where $\rho_c$ is the specific contact resistance and $A_c$ is the contact area that approximately equals to the TSV cross-section area. Since there is a contact between each through-DRAM power TSV and the top metal layer of every DRAM or processor die, $R_{tsv\_p}$ can be expressed as

$$R_{tsv\_p} = \rho \frac{N_{dram}T_t}{\pi(\frac{D_{tsv\_p}}{2})^2} + \frac{(2N_{dram} - 1)\rho_c}{\pi(\frac{D_{tsv\_p}}{2})^2} \tag{3}$$

where $\rho$ is the resistivity of the material used to form power TSV, $N_{dram}$ is the number of DRAM dies, and $D_{tsv\_p}$ is the diameter of through-DRAM power TSVs. Let $N_{tsv\_p}$ denote the total number of through-DRAM power TSVs, among which half of the TSVs connect to the power supply and half of the TSVs connect to the ground. Hence, the average current flowing through each through-DRAM power TSV, denoted as $I_{tsv\_p}$, can be estimated as

$$I_{tsv\_p} = \frac{P_{pro}}{V_{dd}\frac{N_{tsv\_p}}{2}} \tag{4}$$

where $P_{pro}$ is the power consumption of the processor die and $V_{dd}$ is processor die power supply voltage. The power consumption on all the through-DRAM power TSVs, denoted as $P_{tsv\_p}$, can be written as

$$P_{tsv\_p} = N_{tsv\_p}I_{tsv\_p}^2 R_{tsv\_p}. \tag{5}$$

Therefore, combining the above equations, we have that the power consumed by all the through-DRAM power TSVs is

$$P_{tsv\_p} = \frac{16}{\pi} \cdot \frac{\rho N_{dram}P_{pro}^2(T_a + T_m + T_{si} + T_{sub})}{N_{tsv\_p}V_{dd}^2 D_{tsv\_p}^2}$$
$$+ \frac{16}{\pi} \cdot \frac{P_{pro}^2(2N_{dram} - 1)\rho_c}{N_{tsv\_p}V_{dd}^2 D_{tsv\_p}^2}. \tag{6}$$

Hence, given the value of $\rho$, $N_{dram}$, $P_{pro}^2$, $T_a$, $T_m$, $T_{si}$, $T_{sub}$, $N_{tsv\_p}$, $D_{tsv\_p}$, and $\rho_c$, we can calculate the power overhead induced by the through-DRAM power TSVs using (6). Meanwhile, the area occupied by power and signal TSVs, denoted as $A_{tsv}$, can be expressed as

$$A_{tsv} = L_{dram\_v}W_{pch}N_{tsv\_pch\_h}$$
$$+ L_{dram\_h}W_{pch}N_{tsv\_pch\_v} + N_{tsv\_s}P_s^2 \tag{7}$$

where $L_{dram\_v}$ and $L_{dram\_h}$ are the length of DRAM along vertical direction and horizontal direction, respectively, $W_{pch}$ is the width of through-DRAM power TSV channel, $P_s$ is the pitch of signal TSV, $N_{tsv\_pch\_h}$, $N_{tsv\_pch\_v}$, and $N_{tsv\_s}$ are the number of horizontal power TSV channels, vertical power TSV channels, and signal TSVs, respectively.

### B. Performance Evaluation

To quantitatively evaluate the area penalty on 3-D DRAM after adding through-DRAM power and signal TSVs, we modified CACTI 5.3, the latest version of a memory modeling tool

TABLE I
SYSTEM PERFORMANCE WITH POWER AND SIGNAL TSVs

| $T_t$ ($\mu$m) | $D_{tsv\_p}$ ($\mu$m) | DRAM Sub-array Size (mm$^2$) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 0.0578 | | | 0.0289 | | | 0.0145 | | |
| | | $N_{tsv\_p}$ | Capacity Degradation | $P_{tsv\_p}/P_{pro}$ | $N_{tsv\_p}$ | Capacity Degradation | $P_{tsv\_p}/P_{pro}$ | $N_{tsv\_p}$ | Capacity Degradation | $P_{tsv\_p}/P_{pro}$ |
| 10 | 2 | 23468 | 0.43% | 28.93% | 46031 | 0.66% | 14.75% | 95445 | 1.07% | 7.11% |
| | 4 | 11827 | 0.62% | 14.33% | 23423 | 1.23% | 7.25% | 27008 | 3.58% | 6.28% |
| | 8 | 6022 | 1.50% | 7.04% | 12120 | 2.63% | 3.50% | 13915 | 4.98% | 3.05% |
| | 16 | 3114 | 2.96% | 3.41% | 6465 | 4.83% | 1.64% | 7476 | 7.90% | 1.42% |
| 20 | 4 | 11827 | 0.71% | 16.19% | 23423 | 1.32% | 8.18% | 27008 | 3.67% | 7.09% |
| | 8 | 6022 | 1.59% | 7.95% | 12120 | 2.72% | 3.95% | 13915 | 5.07% | 3.44% |
| | 16 | 3114 | 3.05% | 3.84% | 6465 | 4.92% | 1.85% | 7476 | 8.00% | 1.60% |
| | 32 | 1659 | 6.10% | 1.80% | 3640 | 10.01% | 0.82% | 4165 | 14.33% | 0.72% |
| 40 | 8 | 6022 | 1.76% | 9.76% | 12120 | 2.89% | 4.85% | 13915 | 5.25% | 4.22% |
| | 16 | 3114 | 3.23% | 4.72% | 6465 | 5.11% | 2.27% | 7476 | 8.19% | 1.97% |
| | 24 | 2144 | 4.74% | 3.05% | 4574 | 7.62% | 1.43% | 5263 | 11.29% | 1.24% |
| | 32 | 1659 | 6.29% | 2.21% | 3640 | 10.21% | 1.01% | 4165 | 14.54% | 0.88% |
| 80 | 16 | 3114 | 5.95% | 6.47% | 6465 | 5.48% | 3.12% | 7476 | 8.58% | 2.69% |
| | 24 | 2144 | 5.11% | 4.17% | 4574 | 8.01% | 1.96% | 5263 | 11.69% | 1.70% |
| | 32 | 1659 | 6.67% | 3.03% | 3640 | 10.62% | 1.38% | 4165 | 14.96% | 1.21% |
| | 40 | 1367 | 8.27% | 2.36% | 3076 | 13.31% | 1.05% | 3505 | 18.38% | 0.92% |

Note: The total DRAM storage capability without through-DRAM TSVs is 1.18Gb, 1Gb, and 0.83Gb when sub-array size is 0.0578mm$^2$, 0.0289mm$^2$, and 0.0145mm$^2$, respectively.

CACTI [9], to explicitly take into account of the area overhead induced by through-DRAM power and signal TSVs. In all the simulations, we assume that the through-DRAM TSVs are made of tungsten [15]. Resistivity of tungsten at 20 °C is $5.6 \times 10^{-8}\Omega - m$, and specific contact resistance is set to $2 \times 10^{-8}\Omega - cm^2$ [16]. We assume the 3-D integrated processor-DRAM system has a footprint of 101 mm$^2$. We note that current high-end multi-core microprocessors tend to have a die size of around 200 mm$^2$ among which almost half is occupied by a shared last-level on-chip cache, e.g., the INTEL Xeon X5482 45 nm quad-core processor has a die size of 214 mm$^2$ with a shared 12 MB L2 cache and 150 W power consumption [17]. Since the high-capacity last-level on-chip cache can either be migrated into the 3-D stacked DRAM or largely reduced because of the 3-D stacked DRAM, this work sets the footprint as 101 mm$^2$ and the processor power consumption as 150 W. We fix the number of through-DRAM signal TSVs as 2048. The processor power supply voltage is set to 1 V. In the foreseeable future, practical aspect ratio of TSV is 10:1 or lower [18]. Hence, in this work the maximum aspect ratio of TSV is set to 5:1 and the signal TSVs always use the maximum aspect ratio. A through-DRAM power TSV channel is created every four sub-arrays horizontally and vertically. We set that 92% of the power through-DRAM TSVs are used for processor power delivery and the others are used for DRAM power delivery.

Given the above setup and assuming we stack 10 DRAM dies, we use CACTI to estimate the area of 3-D DRAM after adding through-DRAM TSVs and use the formulas derived above to estimate the power consumptions overhead introduced by through-DRAM TSVs with different configurations in terms of DRAM die thickness, DRAM sub-array size and through-DRAM power TSV diameter. The results are summarized in Table I, which clearly shows the tradeoff between the through-DRAM power TSVs' impacts on storage capacity and power consumption. This can be intuitively justified as follows: In order to reduce the power consumed by the power TSVs, we must reduce the aggregate power TSV resistance by increasing the power TSVs' number and/or the size of
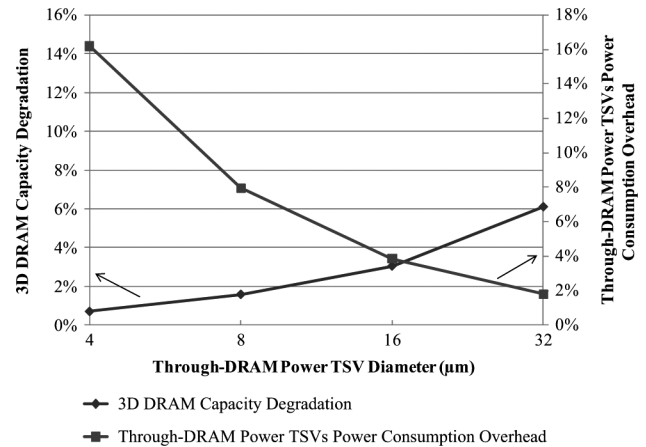


Fig. 5. Impact of individual power TSV size on the DRAM storage capacity versus power consumption overhead tradeoff. DRAM die thickness is 20 $\mu$m and the number of DRAM dies is 10.

each power TSV, which will inevitably occupy more area on the DRAM dies and hence result in a higher DRAM storage capacity degradation. Results listed in Table I quantitatively shows how different parameters could impact such storage capacity versus power overhead tradeoff. For the purpose of further illustration, Fig. 5 shows the impacts of the size of each individual power TSV. As we reduce the DRAM sub-array size, more power TSV channels will be created and hence the number of TSVs will increase, which directly leads to less power consumption overhead at the cost of a higher DRAM storage capacity degradation. This is illustrated in Fig. 6.

Another important parameter is the DRAM substrate thickness that directly determines the length of through-DRAM TSVs. To maintain the same aggregate resistance of power TSVs, as the DRAM substrate thickness increases, we have to accordingly increase the aggregate size of power TSVs, leading to a higher DRAM storage capacity degradation. As pointed out in the above, we set the maximum aspect ratio of TSV is 5:1, so the minimum TSV diameter corresponding to substrate thickness of 10, 20, 40, and 80 $\mu$m is 2, 4, 8, and 16 $\mu$m,
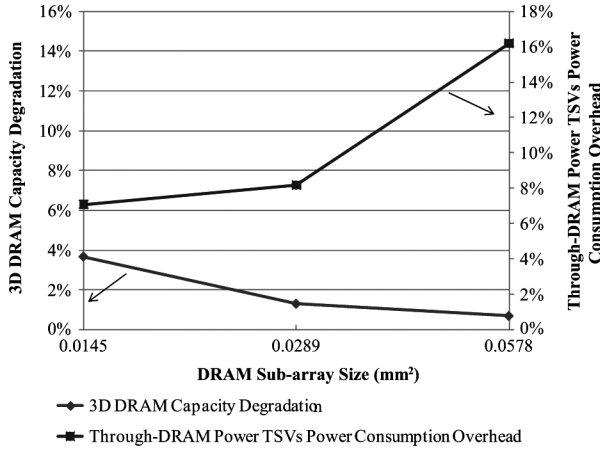
Fig. 6. Impact of individual DRAM sub-array size on the DRAM storage capacity versus power consumption overhead tradeoff. DRAM die thickness is 20 $\mu$m, diameter of power TSV is 4 $\mu$m, and the number of DRAM dies is 10.
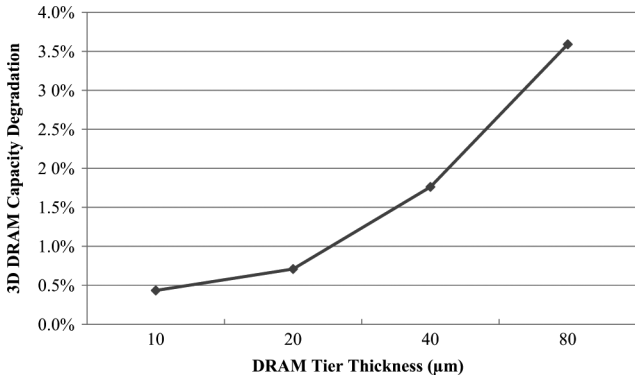


Fig. 7. 3-D DRAM capacity degradation with DRAM die thickness $(T_t)$ increasing. Through-DRAM power TSV maximum aspect ratio is set to 5:1 and the number of DRAM dies is set to 10.



Fig. 8. Impact of contact resistance on through-DRAM power TSVs power consumption. DRAM die thickness is set to 20 $\mu$m and the number of DRAM dies is set to 10.

respectively, as shown in Table I. Fig. 7 further illustrates the impact of DRAM substrate thickness.

Finally, it should be pointed out that the contact resistance $\rho_c$ also has a significant impact. Results shown in Table I are obtained by assuming the specific contact resistance is $2 \times 10^{-8} \Omega -$ cm$^2$. In practice, this value may vary and depends on particular fabrication and stacking technologies. In this study, we further vary the specific contact resistance from $4 \times 10^{-9} \Omega -$ cm$^2$ to $1 \times 10^{-7} \Omega -$ cm$^2$ to investigate the impacts of contact resistance, and the results are shown in Fig. 8. It clearly shows that the TSV power consumption overhead $P_{\text{tsv\_p}}$ will increase dramatically as the specific contact resistance increases. Meanwhile, under a bigger specific contact resistance, $P_{\text{tsv\_p}}$ will change more dramatically as we reduce the TSV size. As pointed out in the above, we set that there is a TSV-metal contact every DRAM die, which appears to be the most practical scenario in the foreseeable future. However, as suggested by the results shown in Fig. 8, if the specific contact resistance turns to be too large, we have to rely on more advanced TSV fabrication and 3-D stacking technologies that could reduce the TSV-metal contact number and/or resistance. For instance, if a small number of DRAM dies are stacked and every DRAM die is very thin, we may be able to fabricate TSVs that can go all the way through these DRAM dies without etch stop and mean-
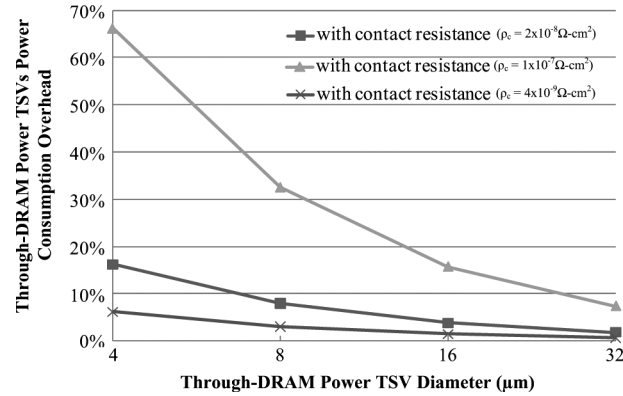
while maintain a reasonable TSV aspect ratio. In this case, we can first align and bond all the DRAM dies together, then etch a via through all DRAM dies without any TSV-metal contact, followed by via filling. If we have to stack many DRAM dies, to maintain a reasonable TSV aspect ratio we can partition all the DRAM dies into several groups, within each group we fabricate TSVs without an etch stop, and between adjacent groups we use one etch stop. Therefore, although we cannot remove all the TSV-metal contacts, we can still effectively reduce the TSV-metal contact number. Finally, we note that in the absence of TSV-metal contact at each DRAM tier, those through-DRAM power TSVs can no longer provide power to the DRAM dies due to the insulation barrier around each TSV. Hence, a separate power delivery system must be implemented to provide power to those DRAM dies, e.g., as we pointed out earlier that we could reserve 8% TSVs for DRAM power delivery by intentionally introducing TSV-metal contact at each DRAM tier, while make the other 92% TSVs go all the way through DRAM tiers without any contact.

## IV. POWER DELIVERY NETWORK INTEGRITY ANALYSIS

As pointed out above, there is a tradeoff between the through-DRAM power TSVs' impacts on storage capacity and power consumption. When we want to reduce the power consumed on through-DRAM power TSVs, we need to sacrifice more DRAM die area to layout those through-DRAM power TSVs. Beyond this tradeoff, the on-chip IR drop will also be affected by the through-DRAM power TSVs. As the through-DRAM power TSV channel resolution becomes coarse, the region surrounded by adjacent through-DRAM power TSV channels will become bigger. As a result, the IR drop within this region tends to increase, especially at the center of this region. In this section, we present further studies on power delivery integrity analysis when the above through-DRAM power TSV implementation strategy is being used.

### A. Power Delivery Circuit Model

In order to analyze the power delivery integrity, we need to first build the power delivery model for the entire 3-D processor-DRAM integrated system. As pointed out earlier, we set that the 3-D DRAM dies locate on top of the processor, as illustrated in Fig. 1. The on-chip grid structured power delivery networks
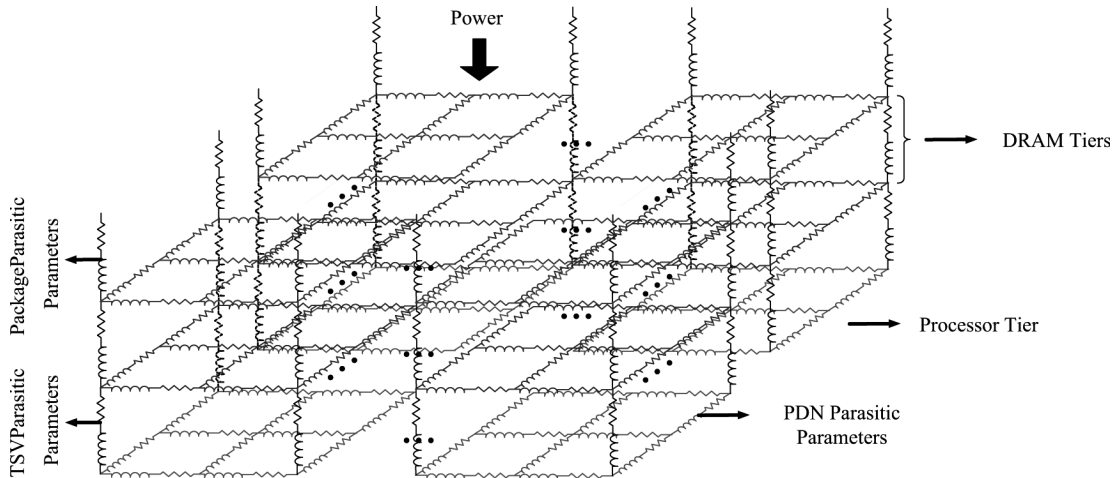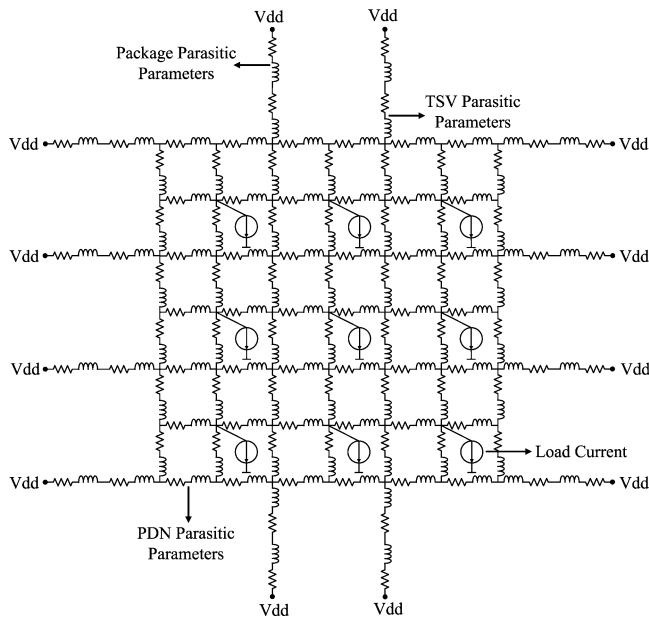
Fig. 9.   Power delivery network model.



Fig. 10.   PDN unit model for bottom processor die.



Fig. 11.   Current load model.

simulation. The width of the power and ground lines is set to 2 $\mu$m and the line spacing is set to 10 $\mu$m, i.e., the line pitch is set to 12 $\mu$m. The unit length inductance is set to 0.05 pH. The package parasitic parameters are set to 0.2 nH and 0.01 $\Omega$. The resistance of each power/ground TSV through each DRAM die is estimated using the method described in Section III-B. The supply voltage is 1 V. The PDN resolution in one PDN unit is decided by the PDN unit size. For example, if the PDN unit size is 680 $\mu$m $\times$ 680 $\mu$m which equals to the area of sixteen DRAM sub-arrays, given the power/ground line pitch of 12 $\mu$m, we have that the PDN unit resolution is $57 \times 57$ (i.e., $680/12 = 57$). Following the discussions in Section III-B, we keep the footprint size and power consumption of the processor die as set to 101 mm$^2$ and 150 W (70% for dynamic power and 30% for static power, i.e., leakage power), respectively. The power consumption of DRAM is set to 3 W. We assume 16 current load evenly distributed within each PDN unit. Every current load is modeled as a triangular current source with 100 ps rise time, 150 ps fall time, and 300 ps cycle time (2.5 GHz), as shown in Fig. 11 [21]. For the PDN under evaluation, we set the total current as two times higher than the average PDN unit current across the processor die.

Figs. 12 and 13 show the simulated maximum IR drop results under different DRAM sub-array sizes and different power TSV sizes. The maximum IR drop steadily decreases as DRAM sub-array size decreases, as shown in Fig. 12. When DRAM sub-array size decreases, the size of each PDN unit, which is surrounded by adjacent through-DRAM power TSV channels, becomes smaller, which in turn leads to less maximum IR drop. Meaning while, as shown in Fig. 6, the DRAM capacity degradation tends to increase as DRAM sub-array size decreases.

(PDN) consists of an array of uniformly spaced metal wires [19], [20]. Fed from the package through power I/Os, power enters from the top-most die and travels to the lower dies through TSVs and PDN, as shown in Fig. 9. We call the PDN between two adjacent through-DRAM power TSV channels as a PDN unit. Since we only concern the power delivery for the load on the bottom-most processor die, we further derive the PDN unit for processor die as shown in Fig. 10.

### B. Simulation Results

In the following simulations, we use two metal layer PDN unit as an example to show the design tradeoffs. The typical sheet resistance ($\Omega/\square$) in 45 nm node for top, medium and bottom metal layer is 0.031, 0.196, and 0.224, respectively. In real chips, the PDN can be composed by all the metal layers, but it should be mainly built using top two metal layers. For purpose of simplicity, we use two-metal-layer PDN unit and set the sheet resistance as $(0.224 \times 0.2 + 0.031 \times 0.8) = 0.07\Omega/\square$ in our
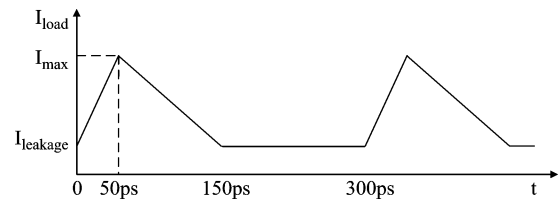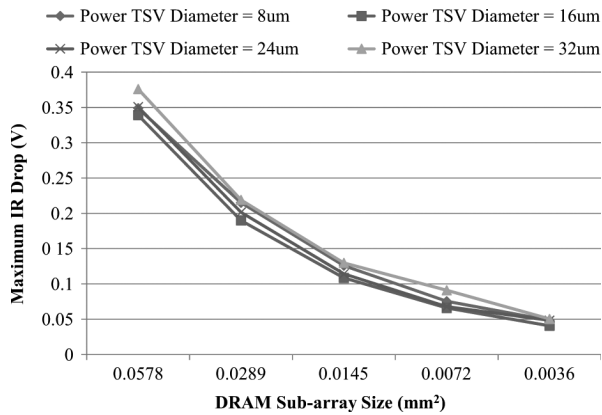
Fig. 12. Maximum IR drop under different DRAM sub-array size (without decoupling capacitors). DRAM die thickness is set to 40 $\mu$m and the number of DRAM dies is set to 10.
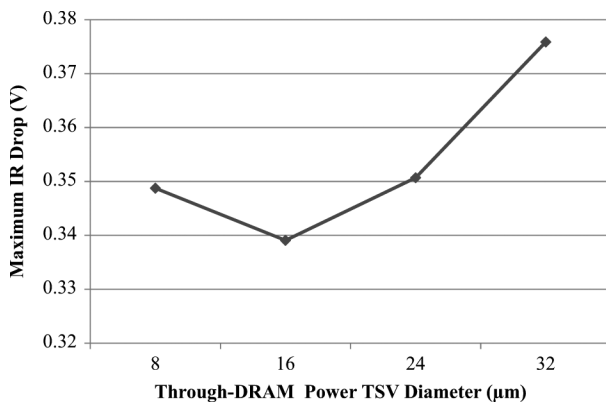


Fig. 13. Maximum IR drop under different through-DRAM power TSV diameters (without decoupling capacitors). DRAM die thickness is set to 40 $\mu$m and the number of DRAM dies is set to 10. DRAM sub-array size is set to 0.0578 mm$^2$.

Therefore, the results suggest a clear tradeoff between IR drop and DRAM capacity degradation. However, the effects of through-DRAM power TSV diameter on the maximum IR drop are more complex. As shown in Fig. 13, when the diameter of through-DRAM power TSVs increases from 8 $\mu$m to 32 $\mu$m, the maximum IR drop will decrease first and then increase. As we increase the diameter of through-DRAM power TSVs, the resistance of through-DRAM power TSVs will accordingly reduce, which can reduce IR drop. However, the number of through-DRAM power TSVs around each PDN unit will also reduce as power TSVs become bigger, which tends to result in larger IR drop within each PDN unit. Therefore, there should be an optimal power TSV diameter, which is 16 $\mu$m in our simulation as shown in Fig. 13. Simulation results for DRAM tier PDN shows the same trend except that the maximum IR drop is always well below 10% of the power supply voltage, which suggests that 8% of power through-DRAM TSVs appears to be enough for DRAM tiers power delivery.

## V. REALIZATION OF DECOUPLING CAPACITORS IN 3-D DRAM

As on-chip power supply voltage continues to reduce in order to make power consumption under control, high-performance integrated circuits, particularly microprocessors, are subject to
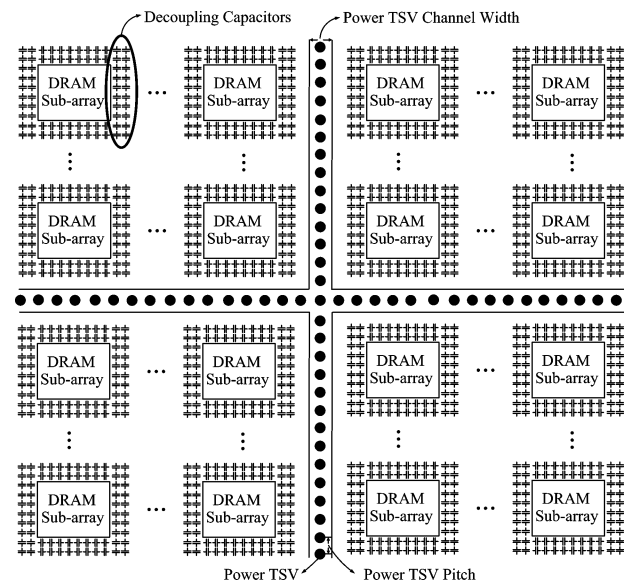


Fig. 14. Proposed decoupling capacitors network in 3-D DRAM.

increasingly significant run-time voltage fluctuations or variations. This has greatly contributed to the emerging design integrity crisis [21]. Decoupling capacitor insertion is one of the most effective means to reduce voltage variations [19], [22]. However, in current design practice, decoupling capacitors may occupy a non-negligible amount of chip area. More importantly, it is nontrivial to decide where and how much decoupling capacitors should be inserted across the chip. This tends to demand accurate knowledge of circuit run-time activities across the die, which may not always be available in the design time. Intuitively, the longer distance between the package substrate and processor can make this issue more critical in 3-D processor-DRAM integration.

For our interested 3-D processor-DRAM integrated computing systems, because DRAM fabrication process is highly optimized for fabricating high-capacitance capacitors with small area, it is intuitive that we should use the stacked DRAM to provide decoupling capacitors for the processor die. To embed decoupling capacitors in 3-D DRAM, one may expect to simply put decoupling capacitors exactly on top of the hot spots of the processor. However, such an ad hoc approach tends to suffer from two drawbacks: 1) DRAM has a very dense and regular structure. Such *ad hoc* decoupling capacitor allocation will introduce a certain degree of structural non-regularity, which can greatly complicate the DRAM design and degrade DRAM storage density. 2) Using such an ad hoc approach demands the design of processor and DRAM closely coupled, which will largely complicate the overall system design and make the 3-D DRAM design less reusable.

To solve the above problems, we propose an uniform decoupling capacitors network design method as shown in Fig. 14. One or more circles of decoupling capacitors are formed around each individual DRAM sub-array. Each individual decoupling capacitor is fabricated as the capacitor being used in each DRAM cell. The exact amount of decoupling capacitors is determined according to system constrains, such as clock frequency, worst-case peak current load, and maximum voltage ripple tolerance. The decoupling capacitors on each side of
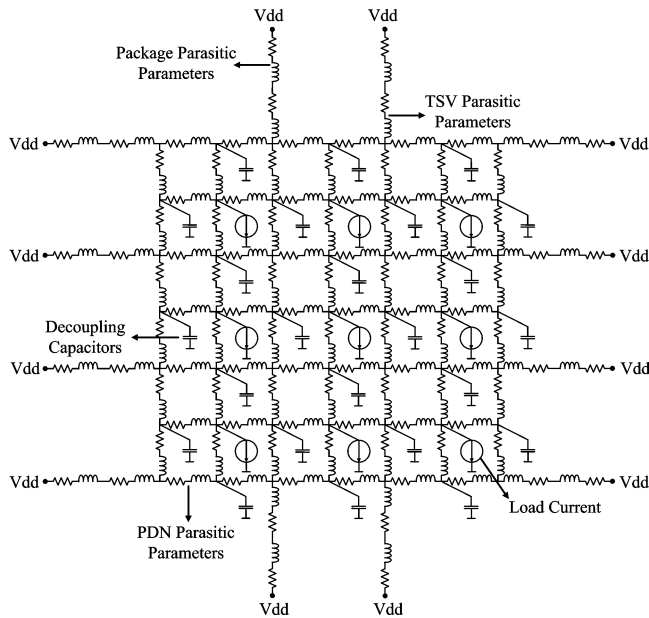
Fig. 15. Power delivery network model for processor die with decoupling capacitors.



Fig. 16. Circuit model of power delivery network.

one DRAM sub-array are connected together in parallel as a large side decoupling capacitor. The side decoupling capacitors vertically aligned across all the DRAM dies are further connected in parallel through TSVs and form a super decoupling capacitor, which is fed to the processor. Clearly, this simple uniform decoupling capacitor network design well fits to the inherently dense and regular DRAM structure, and meanwhile largely decouples the design of processor and DRAM. Once we decide the locations of all the TSVs that connect processor and DRAM, the processor and DRAM design can be carried out independently, which can reduce the overall system design complexity and improve component reusability. Finally, since all the decoupling capacitors along each side of DRAM sub-arrays are connected together to form a super decoupling capacitor, only small number of TSVs, which is linearly proportional to the number of DRAM sub-arrays, is required.

We carried out simulations to evaluate the performance of the above uniform decoupling capacitor network design in 3-D processor-DRAM integrated systems. In the following, we first present the circuit model under which the circuit SPICE simulations are carried out, then present the simulation results.

### A. Power Delivery Circuit Model

The PDN unit model without decoupling capacitors is shown in Fig. 10. After adding decoupling capacitors around every DRAM sub-array, the PDN unit model with decoupling capacitors is shown as Fig. 15, which can be simplified to the circuit model as shown in Fig. 16. The parasitic parameters in the rectangles from top to bottom represent the TSV parasitic parameters, decoupling capacitors parasitic parameters, local PDN parasitic parameters, global PDN parasitic parameters, and package parasitic parameters respectively. Functionally, the circuits can be divided into two parts, including the following:

1) decoupling capacitors discharging circuits through which the decoupling capacitors can release charges to the load when there is a big current on the load;
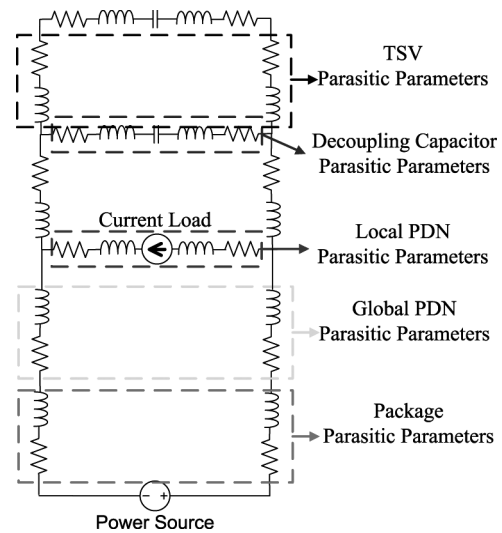
2) decoupling capacitors recharging circuits through which the decoupling capacitors restore its charges through power supply before the next clock cycle.

Both the distance between power supply and decoupling capacitors and the distance between decoupling capacitors and current load affect the effectiveness of the decoupling capacitors. As the distance between decoupling capacitors and current load increases, the voltage drop on the parasitic resistance and inductance of the PDN will become bigger, which will degrade the decoupling effectiveness. The time required to recharge the decoupling capacitors depends on the distance between power supply and decoupling capacitors. If the recharging time is longer than the current load period, decoupling capacitors cannot be fully recharged and may fail to reach its full decoupling capability.

### B. Simulation Results

In this work, we assume that two circles of decoupling capacitors are added around each DRAM sub-array. We still use CACTI to model each DRAM die while considering the area overhead incurred by the extra circles of decoupling capacitors around each individual DRAM sub-array. The footprint of the 3-D processor-DRAM integrated system is still fixed as 101 mm$^2$ as we set in the Section V-B, under which each DRAM die can support 1 Gb with 2048 I/Os based on CACTI estimation results at 45 nm node. All the other simulation parameters and constraints are the same as those in Section V-B. Given these configurations, we use the circuit model shown in Fig. 15 to run SPICE simulations to estimate the voltage of every node in this circuit model, find the minimum value of them and hence estimate the maximum IR drop. Table II lists estimated maximum IR drop under different configurations.

As shown in Table II, we considered a set of combinations of different size of individual DRAM sub-array hence different decoupling capacitance around each sub-array, different specific contact resistance and different number of DRAM dies. The results quantitatively show the trade-off between decoupling efficiency and DRAM storage capacity: As we reduce the size of each DRAM sub-array, the uniform decoupling capacitor network will become denser and the PDN unit size will be-

TABLE II
DECOUPLING CAPACITORS PERFORMANCE

| Sub-array Area (mm²) | Decap. per sub-array (pF) | # of dies | Maximum IR Drop (V) | | | DRAM Storage Capacity Reduction |
|---|---|---|---|---|---|---|
| | | | Cont. Resis. 1 | Cont. Resis. 2 | Cont. Resis. 3 | |
| 0.0578 | 702 | 2 | 0.194 | 0.243 | 0.273 | 0.40% |
| | | 5 | 0.230 | 0.258 | 0.302 | |
| | | 10 | 0.249 | 0.285 | 0.357 | |
| 0.0289 | 410 | 2 | 0.098 | 0.117 | 0.153 | 0.58% |
| | | 5 | 0.110 | 0.129 | 0.188 | |
| | | 10 | 0.123 | 0.144 | 0.219 | |
| 0.0145 | 264 | 2 | 0.048 | 0.057 | 0.084 | 0.71% |
| | | 5 | 0.054 | 0.063 | 0.102 | |
| | | 10 | 0.060 | 0.075 | 0.129 | |
| 0.0072 | 205 | 2 | 0.024 | 0.027 | 0.042 | 0.93% |
| | | 5 | 0.027 | 0.030 | 0.057 | |
| | | 10 | 0.03 | 0.039 | 0.075 | |
| 0.0036 | 133 | 2 | 0.011 | 0.012 | 0.021 | 1.19% |
| | | 5 | 0.012 | 0.018 | 0.033 | |
| | | 10 | 0.015 | 0.021 | 0.045 | |

Cont. Resis. $1 = 4 \times 10^{-9} \Omega$-cm², Cont. Resis. $2 = 2 \times 10^{-8} \Omega$-cm², Cont. Resis. $3 = 1 \times 10^{-7} \Omega$-cm²
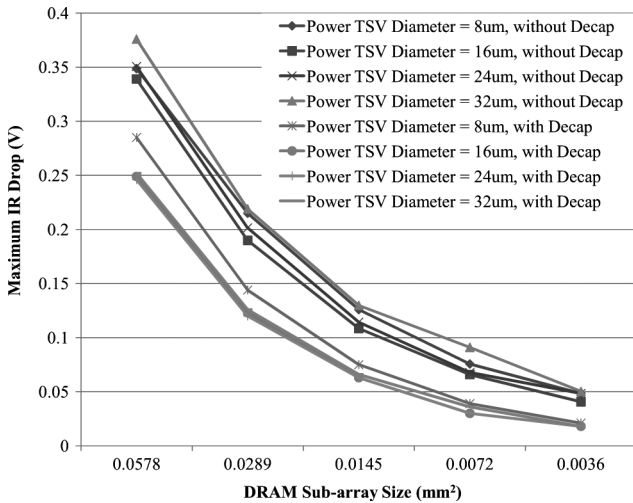


Fig. 17. Maximum IR drop under different DRAM sub-array size (with and without decoupling capacitors). DRAM die thickness is set to 40 $\mu$m and the number of DRAM dies is set to 10.



Fig. 18. Maximum IR drop vs. the number of DRAM dies with different DRAM sub-array size. DRAM die thickness is set to 40 $\mu$m. The contact resistance is set to $2 \times 10^{-8} \Omega - $cm². The diameter of the through-DRAM power TSV is set to 8 $\mu$m.



Fig. 19. Effect of TSV specific contact resistance on decoupling efficiency. DRAM die thickness is set to 40 $\mu$m and the number of DRAM dies is set to 10. The diameter of the through-DRAM power TSV is set to 8 $\mu$m.

come smaller, which leads to smaller IR drop under the same current load. Clearly, a denser decoupling capacitor network tends to occupy more silicon area and hence reduce the overall DRAM storage capacity. To further illustrate the results, Fig. 17 shows the maximum IR drop in one PDN unit with and without decoupling capacitors under different DRAM sub-array size. With decoupling capacitors, the maximum IR drop can be reduced up to 61.3% compared to the scenarios without decoupling capacitors.

Moreover, the number of DRAM dies may also affect the decoupling efficiency. Results in Table II show that the maximum IR drop will accordingly increase as the number of DRAM dies increases. Although more DRAM dies provide more decoupling capacitors, it will meanwhile increase the recharging path from the power supply to decoupling capacitors. Since the recharging path play a more important role, the decoupling efficiency tends to degrade as we stack more DRAM dies. Fig. 18 further illustrates the maximum IR drop versus the number of DRAM dies under different DRAM sub-array sizes.

We note that the effectiveness of decoupling capacitors is also affected by the parasitic parameters of power TSVs. It is clear
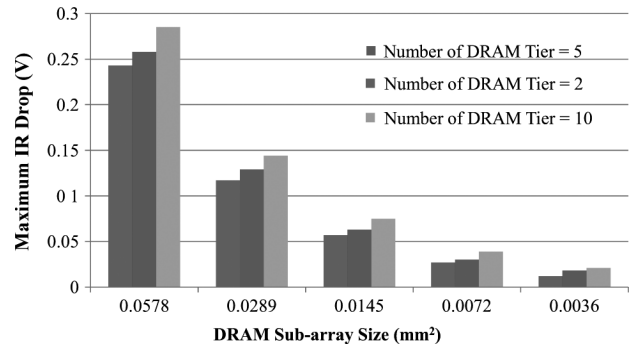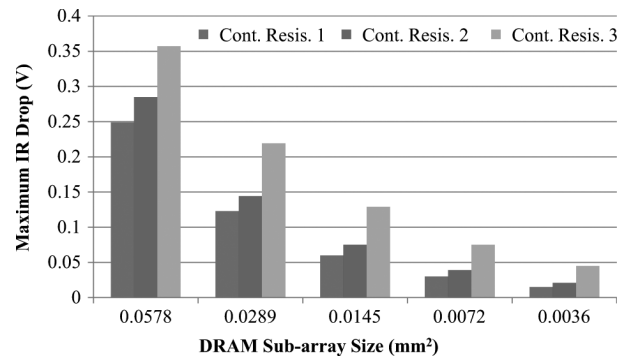
that the decoupling effectiveness drops as the TSV specific contact resistance becomes bigger. For the purpose of evaluation, we carry out SPICE simulations over a range of TSV specific contact resistance values as shown in Fig. 19.

## VI. FULL COMPUTING SYSTEM SIMULATIONS

As we can see from the above discussions, under the same footprint size 3-D DRAM will have different capacity given different design parameters. In a consequence, the 3-D processor-DRAM integrated system will have different memory resources, which can lead to different computing system performance that
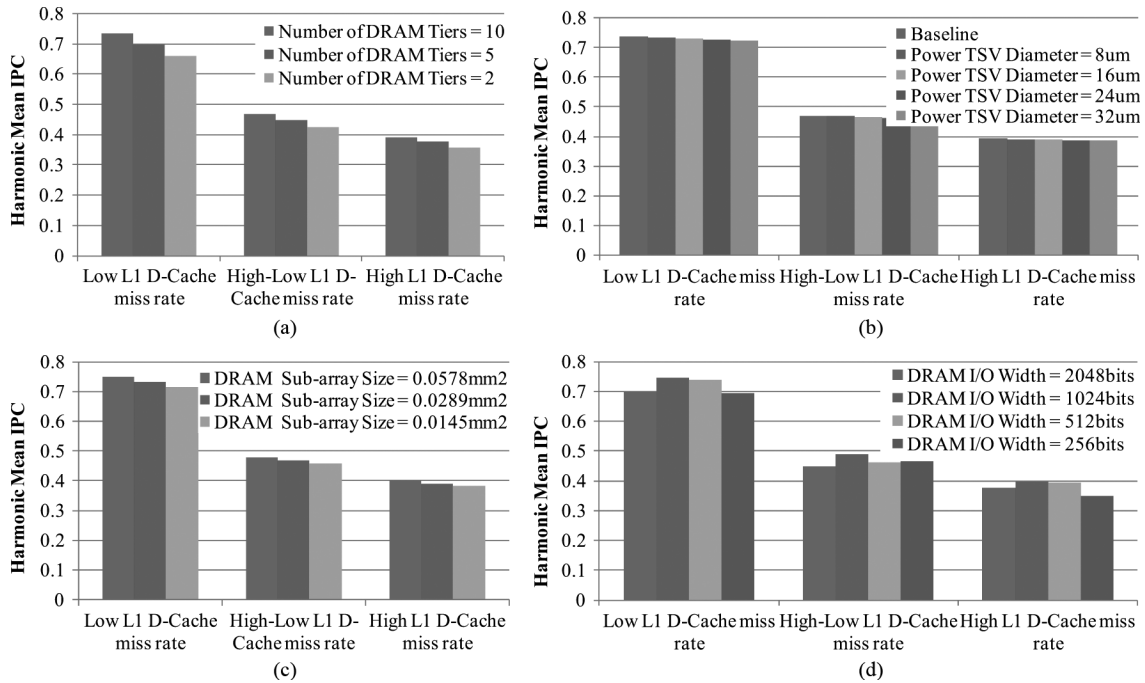
Fig. 20. Simulated computing system performance under different configurations: (a) harmonic mean IPC with different number of DRAM tiers; (b) harmonic mean IPC with different power TSV diameter; (c) harmonic mean IPC with different DRAM sub-array size; (d) harmonic mean IPC with different DRAM I/O width.

TABLE III
BENCHMARK CONFIGURATION IN FULL SYSTEM SIMULATIONS

| Benchmark Set | Four SPEC2000 Benchmarks |
|---|---|
| High-Low L1 D-Cache miss rate | art, ammp, gcc, stream |
| High L1 D-Cache miss rate | bzip2, mcf, mesa, gzip |
| Low L1 D-Cache miss rate | stream, equake, tigr, twolf |

TABLE IV
CONFIGURATION PARAMETERS OF M5 SIMULATOR

| M5 parameters | Values |
|---|---|
| # of cores | 4 |
| Core freq. | 2.5GHz |
| Core type | Alpha 21264, out-of-order |
| L1 cache | 16KB for both data and inst., 4 way, 64B block, private, hit latency: 1 cpu cycle |
| 3D DRAM | total capacity & latency: varies miss penalty: 200 cpu cycles (80ns) |

is typically measured in terms of instruction per cycle (IPC). In order to investigate the impacts of different DRAM design parameters on computing system performance, we carried out system performance simulations using the M5 full system simulator [10] developed by the University of Michigan. We assume the processor die contains four 2.5 GHZ cores and choose three different benchmark sets according to different L1 data cache miss rate, each set contains four different SPEC2000 benchmarks as listed in Table III. Detailed configuration parameters for M5 are listed in Table IV. We note that we remove on-chip L2 cache and assume the 3-D DRAM serves as the next-level memory after L1 cache. The total 3-D DRAM storage capacity and 3-D DRAM access latency are obtained from CACTI modeling based upon various design parameters including DRAM

tier number, TSV diameter, sub-array size, and memory bandwidth. We note that the 3-D DRAM access latency tends not to vary largely and is around 20 ns (i.e., 50 CPU cycles) in this case study.
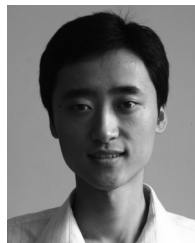
The simulated harmonic mean IPC (HMIPC) performance is shown in Fig. 20. Fig. 20(a) shows the HMIPC under different number of DRAM tiers. HMIPC decreased as the number of DRAM tiers decreased, since less DRAM tiers provided less memory resources. On the other hand, decoupling efficiency will drop as we increase the number of DRAM tiers. Therefore, there should be a tradeoff between the decoupling efficiency and IPC performance. Fig. 20(b) shows the HMIPC under different power TSV diameter. As we have seen from Fig. 5, 3-D DRAM capacity degradation will become bigger as the through-DRAM power TSV diameter increasing, which leads to HMIPC drop slightly as the through-DRAM power TSV diameter increases. Fig. 20(c) shows the HMIPC under different DRAM sub-array sizes. With the same DRAM die footprint, bigger DRAM sub-array size leads to a higher storage density, and hence higher capacity, which can lead to a slightly higher HMIPC. Fig. 20(d) shows the HMIPC under different DRAM I/O width. Although a higher bandwidth can reduce the on-chip cache miss penalty and hence can improve HMIPC, it meanwhile can result in storage density degradation. Therefore, there should be a design tradeoff and optimal configuration of DRAM bandwidth. As shown in Fig. 20(d), as we increase the bandwidth from 256 bits, the HMIPC first improves because the reduced on-chip cache miss penalty overweighs the reduced DRAM storage capacity; however if the bandwidth becomes too large (e.g., 2048 bits), the HMIPC begins to degrade because the reduced SRAM storage capacity begins to overweigh the reduced on-chip cache miss penalty.

## VII. Conclusion

This paper addresses the issue of DRAM design in the presence of through-DRAM TSVs in 3-D processor-DRAM integrated computing systems. We present a simple method to allocate through-DRAM power and signal TSVs, which can well fit to the regular DRAM architecture and hence minimize the interference to the DRAM design. Those through-DRAM TSVs inevitably incur a tradeoff between DRAM storage capacity degradation and power consumption overhead, which has been studied through memory modeling and circuit simulations over a wide range of parameters. To improve the through-DRAM power delivery integrity, this paper further proposes a method to implement decoupling capacitors for the processor die. The objective is to provide abundant decoupling capacitors at minimal cost by leveraging the proximity between processor die and DRAM dies and the superior capacitor fabrication ability of DRAM. A simple uniform decoupling capacitor network design strategy is presented, which directly adds extra capacitors around each individual DRAM sub-array as decoupling capacitors. This method can maintain the inherent regularity of DRAM structure and largely decouple the design of processor and DRAM dies. Further circuit simulations and computing system simulations are carried out to demonstrate the effectiveness of this method and the involved design tradeoffs.

## References

[1] W. A. Wulf and S. A. McKee, "Hitting the memorywall: Implications of the obvious," *Comput. Arch. News*, vol. 23, pp. 20–24, Mar. 1995.

[2] C. C. Liu, I. Ganusov, M. Burtscher, and S. Tiwari, "Bridging the processor-memory performance gap with 3D IC technology," *IEEE Des. Test Comput.*, vol. 22, no. 6, pp. 556–564, Nov.-Dec. 2005.

[3] T. H. Kgil, A. Saidi, N. Binkert, R. Dreslinski, S. Reinhardt, K. Flautner, and T. Mudge, "Picoserver: Using 3D stacking technology to enable a compact energy efficient chip multiprocessor," in *Proc. 12th Symp. Arch. Support for Program. Lang. Operat. Syst.*, 2006.

[4] B. Black, M. Annavaram, N. Brekelbaum, J. Devale, L. Jiang, G. H. Loh, D. Mccauley, P. Morrow, D. W. Nelson, D. Pantuso, P. Reed, J. Rupley, S. Shankar, J. Shen, and C. Webb, "Die stacking (3D) microarchitecture," in *Proc. Annu. IEEE/ACM Int. Symp. Microarch.*, 2006, pp. 469–479.

[5] G. H. Loh, Y. Xie, and B. Black, "Processor design in 3D die-stacking technologies," *IEEE Micro*, vol. 27, no. 3, pp. 31–48, May-Jun. 2007.

[6] G. H. Loh, "3D-stacked memory architecture for multi-core processors," in *Proc. 35th ACM/IEEE Int. Conf. Comput. Arch.*, Jun. 2008, pp. 453–464.

[7] P. G. Emma and E. Kursun, "Is 3D chip technology the next growth engine for performance improvement?," *IBM J. Res. Develop.*, vol. 32, no. 6, pp. 541–552, Nov. 2008.

[8] M. Bohr, "The new era of scaling in an SoC world," in *Proc. IEEE Int. Solid-State Circuits Conf. (ISSCC)*, Feb. 2009, pp. 23–28.

[9] Hewlett-Packard Labs, Palo Alto, CA, "CACTI: An integrated cache and memory access time, cycle time, area, leakage, and dynamic power model," 2008. [Online]. Available: http://www.hpl.hp.com/research/cacti/

[10] University of Michigan, Ann Arbor, "M5: A modular platform for computer system architecture research, encompassing system-level architecture as well as processor microarchitecture," 2006. [Online]. Available: http://www.m5sim.org/

[11] J.-Q. Lu, "3-D hyperintegration and packaging technologies for micro-nano systems," *Proc. IEEE*, vol. 97, no. 1, pp. 18–30, Jan. 2009.

[12] T. A. C. M. Claasen, "An industry perspective on current and future state of the art in system-on-chip (SoC) technology," *Proc. IEEE*, vol. 94, no. 6, pp. 1121–1137, Jun. 2006.

[13] J.-Q. Lu, T. S. Cale, and R. J. Gutmann, "Wafer-level three-dimensional hyper-integration technology using dielectric adhesive wafer bonding," in *Materials for Information Technology: Devices, Interconnects and Packaging*, E. Zschech, C. Whelan, and T. Mikolajick, Eds. London, U.K.: Springer-Verlag, 2005, pp. 386–397.

[14] Z. Xu, A. Beece, K. Rose, and J.-Q. Lu, "Structures and electrical performance of through strata vias (TSVs)," in *Proc. Int. Wafer-Level Packag. Conf.*, Oct. 2009, pp. 24–26.

[15] P. S. Andry, C. K. Tsang, B. C. Webb, E. J. Sprogis, S. L. Wright, B. Dang, and D. G. Manzer, "Fabrication and characterization of robust through-silicon vias for silicon-carrier applications," *IBM J. Res. Develop.*, vol. 52, p. 571, Nov. 2008.

[16] K. N. Chen, A. Fan, C. S. Tan, and R. Reif, "Contact resistance measurement of bonded copper interconnects for three-dimensional integration technology," *IEEE Electron Device Lett.*, vol. 25, no. 1, pp. 20–22, Jan. 2004.

[17] INTEL Corporation, Santa Clara, CA, "Xeon processors," 2009. [Online]. Available: http://www.intel.com

[18] Semiconductor Industry Association, "The international technology roadmap for semiconductors (ITRS)," 2007. [Online]. Available: http://www.itrs.net/Links/2007ITRS/Home2007.htm

[19] M. Popovich, A. V. Mezhiba, and E. G. Friedman, *Power Distribution Networks with On-Chip Decoupling Capacitors*. New York: Springer, 2007.

[20] G. Huang, M. Bakir, A. Naeemi, H. Chen, and J. D. Meindl, "Power delivery for 3D chip stacks: Physical modeling and design implication," in *Proc. Elect. Perform. Electron. Packag.*, Atlanta, GA, Oct. 2007, pp. 205–208.

[21] Q. K. Zhu, *Power Distribution Network Design for VLSI*. Hoboken, NJ: Wiley-IEEE, 2004.

[22] Z. Qi, H. Li, S. Tan, L. Wu, Y. Cai, and X. Hong, "Fast decap allocation algorithm for robust on-chip power delivery," in *Proc. Int. Symp. Quality Electron. Des. (ISQED)*, 2005, pp. 542–547.

**Qi Wu** received the B.S. and M.S. degrees in electrical engineering from Xian Jiaotong University, Xian, China, in 2003 and 2006, respectively. He is currently pursuing the Ph.D. degree from the Electrical, Computer, and Systems Engineering Department, Rensselaer Polytechnic Institute, Troy, NY.

His current research interests include design of 3-D integrated circuits, high performance computing systems, memory and storage system hierarchy, and VLSI architecture and circuits for image and video processing.

**Tong Zhang** (M'02–SM'08) received the B.S. and M.S. degrees in electrical engineering from the Xian Jiaotong University, Xian, China, in 1995 and 1998, respectively, and the Ph.D. degree in electrical engineering from the University of Minnesota, Minneapolis, in 2002.

Currently, he is an Associate Professor with the Electrical, Computer and Systems Engineering Department, Rensselaer Polytechnic Institute, Troy, NY. His current research interests include circuits and systems for data storage, signal processing, and computing.

Dr. Zhang currently serves as an Associate Editor for the IEEE Transactions on Circuits and Systems—Part II: Brief Papers and the IEEE Transactions on Signal Processing.