# Leakage Power Reduction For Coarse Grained Dynamically Reconfigurable Processor Arrays With Fine Grained Power Gating Technique

Yoshiki Saito, Tomoaki Shirai, Takuro Nakamura, Takashi Nishimura, Yohei Hasegawa, Satoshi Tsutsumi, Toshihiro Kashima, Mitsutaka Nakata, Seidai Takeda, Kimiyoshi Usami, Hideharu Amano
Department of Information and Computer Science, Keio University
3-14-1 Hiyoshi, Yokohama, 223-8522 Japan (muccra@am.ics.keio.ac.jp)
Department of Technology, Shibaura Institute of Technology
3-9-14 Shibaura, Minato-ku, Tokyo, 108-8548 Japan

## Abstract

*One of the benefits of coarse grained dynamically reconfigurable processor array(DRPA) is its low dynamic power consumption by operating a number of processing elements(PE) in parallel with low clock frequency. However, in the future advanced processes, leakage power will occupy a considerable part of the total power consumption, and it may degrade the advantage of DRPAs. In order to reduce the leakage power, a fine grained Power Gating(PG) is applied to a DRPA, MuCCRA-2.32b, and leakage power and area overhead are measured. We evaluated the effect of two control modes;* Pair *and* Unit Individual *based on layout design and real applications. It appears that by applying PG for ALUs and SMUs in PEs individually, 48% of leakage power can be reduced with 9.0% of area overhead.*

## 1. Introduction

Coarse grained dynamically reconfigurable processor arrays(DRPAs) have been paid attention as an efficient off-loading engine for various types of System-on-a-Chip. Their circuit composition can be dynamically changed quickly by switching a hardware context, which is the configuration data set prepared before running application. Some DRPAs have been already commercially available in the consumer equipments[2, 5]. One of the major advantages of DRPAs is low power consumption by making the best use of parallel processing in an array of Processing Elements(PEs) with low clock frequency. It achieves better performance of that of FPGAs with much less clock frequency and smaller power consumption[4, 7]. SONY's VME[2] adopted for portable games, due to its advantages to stretch the battery life. Power reduction techniques for further decrease of dynamic power have been investigated[7].

After 90nm CMOS process, the amount of the leakage power is remarkably increased, and it will occupy a considerable portion of the total power in the future processes. Since the leakage power is relational to the chip area, the technique to reduce it will be important especially in DRPAs using a large number of PEs. Experimental results show that not all PEs in DRPAs are utilized in a hardware context. Although the dynamic power of such PEs can be reduced by clock gating and operand isolation, the leakage power consumed in such unused PEs become non-negligible overhead. The leakage power reduction techniques has not been well investigated for DPRAs except a few researches. A context dependent voltage scale control for DRPA was proposed[6], and the power aware synthesis techniques were discussed[8].

Power Gating(PG) is a representative technique to reduce the leakage power for such unused components. Although the conventional PG technique is mainly applied to a large modules like whole part of IPs and CPUs, recent fine grained PG techniques enable to shut down and wake up small components like ALUs, quickly[3].

Here, DRPAs with fine grained PG technique is proposed to cope with the future leakage power problem. The power of unused PEs or unused components in a PE is shut-off if the shut down time is longer than the break even time. The area overhead and leakage power reduction are evaluated, and the utility of applying this technique to such DRPA is considered.

## 2. MuCCRA-2.32b

Here, a model architecture for a target of fine grained PG, MuCCRA-2.32b is introduced. MuCCRA-2.32b is a small scale, multi-context DPRA designed for power analysis and investigating reduction techniques. The architecture is almost the same as MuCCRA-2[1], which is now working on a real chip, but chip dependent design optimization is eliminated for analyzing a common DRPA design.

As shown in Fig. 1, MuCCRA-2.32b has 4×4 PE array and four distributed memories(MEMs), which have 32bit×256entries at the bottom. The granularity of the architecture is 32bits, that is, all functional units treat 32-bit data, except wires for 2-bit carry signals. The carry bits are used for indicating overflow signals of sums or products, signs of subtraction results, for controlling branches, etc. Task Configuration Controller and Context Switching Controller are provided to manage reconfiguration.

Island-style interconnection structure is employed, so that each PE is surrounded by programmable routing wire segments. At each intersection of vertical and horizontal channels, a switching element(SE) is placed, connecting the global routing channels for transferring data between PEs and MEMs. There are three channels for the global routing resources, and this network will form a data path.
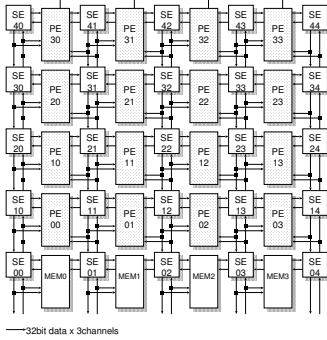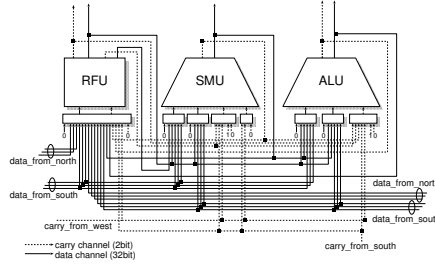
**Figure 1. PE Array**
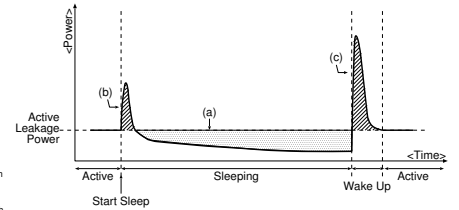


**Figure 2. PE Structure of MuCCRA-2.32b**



**Figure 3. Power Waveform While Sleeping**

In each PE, there are three units inside as shown in Fig. 2: an Arithmetic and Logic Unit(ALU), a Shift and Mask Unit(SMU) and a Register File Unit(RFU). An ALU contains an adder, a subtractor, a multiplier and a logic operator. An SMU processes shift operation, checks the equivalent between two values, and generates an immediate values. These two units output a 32-bit data and a 2-bit carry signal. RFU consists of 8entries for 34-bit flip-flops.

Each PE is connected with global routing wires and these units pick up data from them as input data. On the other hand, all output data from the units can be transferred to the other PEs or MEMs through connection block, called PICKOUT. Each SE consists of three programmable switches(SWs), handling one of the three channels. An SW transfers entering data to desired output direction.

## 3. Fine Grained PG for MuCCRA-2.32b
### 3.1. Fine Grained PG for DRPAs

Since the usage of PEs are not always high in DRPAs, the leakage power can be reduced by shutting-off the power of such unused PEs. One of the reason that this natural idea has not been applied to DRPAs yet is that, the traditional PG targets on the large semiconductor domain corresponding to entire CPUs and IP modules. Moreover, although this method has a benefit that, the target domain can be designed with the common, non-PG design method, the wake-up time is the order of micro seconds.

However, recently, high speed PG for a small component has been tried and we also proposed such a method[3]. A certain number of gates are connected with the same VGND line and multiple sleep transistors are shut-off with a sleep control signal. VGND lines are provided at the fixed position in the cell layout. Since existing ground rails in the standard cells are used as the real ground, non-PG cells such as flip-flops, clock buffers, power switch drivers and isolation cells can be placed at any location in a row. The target unit for PG can be computational modules like multipliers or adders; which are much smaller than the power domain in the traditional method. In order to distinguish from common PG, we call this technique a fine grained PG.

### 3.2. Break Even Time

By providing appropriate number of sleep transistors, the wake-up time can be reduced to a few nano seconds. However, aggressive sleep control may introduce power overhead, which has to be considered carefully. Fig. 3 is a

schematic diagram showing variability of power consumption including both leakage and dynamic over time. The module starts to sleep at the time "Start Sleep", and wakes up at "Wake Up" in the figure. The shadowed area pointed by (a) is the leakage energy reduction effect. As (b) and (c) is pointing, when a PG target domain begins to sleep or wake-up, dynamic power is required for sleep transistors, isolation cells, and buffers in sleep signal wires. Hence, the sleep time must be long enough so that the amount of the leakage energy reduction effect overcomes the overhead. The sleep time where the amount of power reduction is the same as overhead is called "Break Even Time (BET)".

In case of applying the traditional coarse grained PG, which large PG target modules are supposed to sleep for a long time, there was no need to consider power overhead. To use of the advantage of fine grained PG, shutting leakage power off must be done as often as possible and so BET must be correctly estimated. First of all, we must investigate on the PG target modules and find how to control them. For appropriate PG target module selection, leakage power consumption, unit utilization and the BET for each module are measured. Since fine grained PG cannot be applied to flip-flops for now, we must eliminate RFUs from the target. SE is also out of our target, because the area is too small. CSC works every clock cycles so obviously, it is difficult to keep the BET. Here, our target is selected only from ALUs, SMUs and PICKOUTs, or the whole PE except RFU.

#### 3.2.1 Leakage Power
The leakage power consumption of each unit in MuCCRA-2.32b is evaluated as follows. First, MuCCRA-2.32b is synthesized with ASPLA 90nm technology library, using Synopsys Design Compiler 2006.06-SP2, and layouted with Synopsys Astro 2007.03-SP3. After RC extraction with Mentor Calibre 2007.3_18.11, the leakage power is evaluated by Synopsys HSIMplus Z-2007.03. The same tools and target processes are also used in the rest of this paper.

The evaluation is done at three different operating temperatures; 25℃, 65℃ and 80℃. As shown in Fig. 4, an ALU and an SMU consume more leakage power comparing to a PICKOUT. This comes from that a PICKOUT is mainly consisting of multiplexers, which occupy smaller area than the other components for computation.

#### 3.2.2 Module Utilization
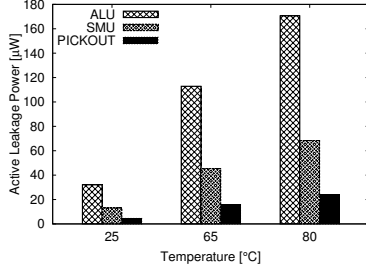To know how often units are used in practical applications, four applications: SHA-1, DCT, DWT and FFT are imple-
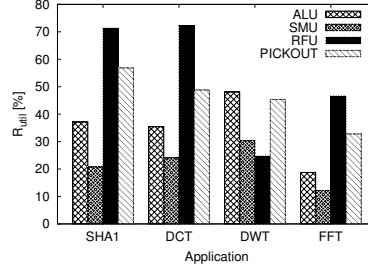
**Figure 4. Leakage Power**
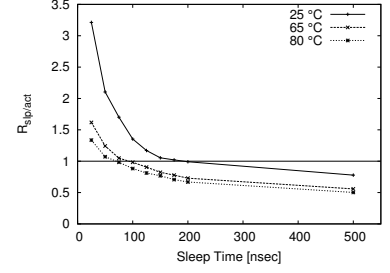


**Figure 5. Unit Utilization**



**Figure 6. BET for an ALU**

mented on MuCCRA-2.32b. The utilization rates $R_{\mathrm{util}}$ for four components in a PE, an ALU, an SMU, an RFU and a PICKOUT are calculated with the following equation,

$$R_{\mathrm{util}} = \frac{\sum_{i=0}^{cycle} U_i}{N \times cycle} \qquad (1)$$

where $i$, $N$, $cycle$ and $U$ indicate the context number, the total number of the unit, execution cycles for an application, and the number of operating units at context $i$, respectively. When more units are used inside loops, $R_{\mathrm{util}}$ increases.

Fig. 5 shows utilization ratio $R_{\mathrm{util}}$ of each component. The average of PE utilization is 65%. Through the applications, the utilization of an RFU and a PICKOUT is high compared with those of an ALU and an SMU.

The low utilization of calculating units comes from that only a component of a pair(ALU or SMU) tends to be used in most applications. Outputs from one of two units are registered to RFU, or transferred to another PE via PICKOUT, thus they are always used resulting high utilization ratio.

### 3.2.3 Evaluating BET

For evaluating the BET for each module, a real design of fine grained PG is needed. Netlist from Design Compiler is modified with hand-made script, and sleep transistors and isolation cells for preventing short-circuit current are inserted. After place and routing with Astro, the number of sleep transistors is optimized with Sequence Design's tool, Coolpower 2007.8.5, so as to make the wake-up time short. Current waveform like Fig. 3 is generated with HSIM using post-layout data. The current is measured for various sleeping time to find BET. The sleep time where the ratio $R_{\mathrm{slp/act}}$ of (a) and (b) + (c) in Fig. 3, becomes 1, is the BET. Fig. 6 indicates the change of the ratio for an ALU.

The clock cycles corresponding to BET of each module is summarized in Table 1 when the clock frequency of MuCCRA-2.32b is assumed to be 40MHz. The BET of an ALU and an SMU are shorter than that of a PICKOUT, since the leakage power, which can be saved, is small in a PICKOUT. Moreover, isolation cells are required at each bit of output port of the modules. A PICKOUT has 204-bit out going ports in total, while an ALU and an SMU have only 34bits. This increases the dynamic power overhead and results longer BET for a PICKOUT.

### 3.2.4 Sleep Signal Control

From Fig. 4, Fig. 5 and Table 1, it seems that a PICKOUT is disadvantageous for a target from several aspects. So, in our research, we focused ALUs and SMUs as the PG targets. There are two options: one is controlling a pair of an ALU

**Table 1. BET for Each Units in PE**

| Temperature | BET [cycles] | | |
|---|---|---|---|
| [℃] | ALU | SMU | PICKOUT |
| 25 | 8 | 18 | 67 |
| 65 | 4 | 6 | 20 |
| 80 | 3 | 4 | 14 |

and an SMU in the same PE together(*Pair*), and the other is controlling every unit individually(*Unit Individual*). The former method can reduce the overhead of sleep control, but has less sleep opportunities.

In order to control sleep signals, we propose to provide sleep control bits in the configuration data. If the time is more than the BET, it is set to shut-off the power. In order to hide the wake-up latency, the bit is reset before a clock earlier than operational context. A simple static branch prediction is enough, since the overhead of miss-prediction caused by complicated branch structures, is just a small increase of power consumption. Note that *Pair* control mode needs only 1bit for each PE, while 2bits are provided for each unit in *Unit Individual* mode.

## 4. Evaluation

### 4.1. Area Overhead

The area overhead for an ALU and an SMU are evaluated by applying fine grained PG, with a method mentioned in the previous section. as shown in Table 2. The standard cells for PG modules are not fully available for synthesizing and optimizing RTL description, due to the limited design time, and resulting larger area. Here, 16% of an ALU and 8.3% of an SMU increased compared to the case when all cells in the library are fully used. This effect is not included in the table, since making the other PG cells is just a problem of time and effort. From the technical viewpoint, it can be certainly solved. In the table, the non-PG modules are designed with the same limited cell library as PG modules. The size of PG cells is exactly the same as it of non-PG cells, thus the difference comes from substantial overhead of fine grained PG, that is, the sleep transistors and isolation cells.

The area overhead of an SMU becomes greater than an ALU. When the target module area is large, many sleep transistors are required for optimization, thus the area overhead ratio for sleep transistors tends to be almost the same. Isolation cells can be another reason of the area overhead.

**Table 2. Area Overhead for Fine Grained PG**

| Unit | Non-PG [$\mu m^2$] | PG [$\mu m^2$] | Overhead [%] |
|------|------|------|------|
| ALU | 23988.28 | 25724.06 | 7.23 |
| SMU | 7750.31 | 8863.04 | 14.4 |
| Total | 31738.59 | 34587.10 | 8.97 |

Since the number of output ports of both units are the same 34bits, the area occupied by isolation cells are the same. Thus, the overhead area is almost the same. Since the original area of an SMU is much smaller than an ALU, the relative overhead becomes large. This fact suggests that the small target module tends to cause a large relative overhead.

### 4.2. Leakage Power Reduction

With the same applications used in Section 3.2.2, the leakage power reduction is evaluated. We evaluated $N$ and $i$ that the PG module sleeps $N$ times in $i$ cycles each during execution of an application. Then, the amount of leakage power can be obtained from the graph shown in Fig. 6. By using this method, precise evaluation including every overhead of fine grained PG can be done. The leakage power reduction is computed with the following expression,

$$P_{\text{pg}} = \frac{E_{\text{non}} - \sum_{i=0} (P_i \times F_i \times t)}{T_{\text{ex}}} \quad (2)$$
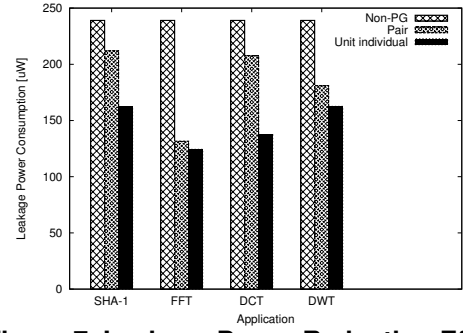
where $P_{\text{pg}}$ is the total reduction effect, $E_{\text{non}}$ means the total leakage energy consumption of a non-PG module, $P_i$ stands for the average leakage power when the PG module could sleep for $i$ cycles, $F_i$ corresponds to the frequency of $i$-cycle sleeping, $t$ is period per cycle, and $T_{\text{ex}}$ is the execution time.

The simulation was done with the temperature of 80℃, and MuCCRA-2.32b's operating frequency is fixed at 40MHz. Although 80℃ is high, we evaluated in the worst case in the real usage. From the results, we can also prospect cases when the future processes with more leakage power will be used.

Fig. 7 shows the leakage power consumption of each application. Here, the sleep bits in configuration data are manually set by programmer. In all application, fine grained PG can reduce the leakage power considering its overhead. In FFT, it becomes up to 48%. Comparing *Pair* and *Unit Individual*, the latter makes the use of more opportunities to sleep than those of *Pair*. For example, in DCT, the reduction effects of *Unit Individual* and *Pair* become 43% and 13%. In SHA-1 and FFT, the utilization of PE is not so high, and the leakage power can be reduced both with *Pair* and *Unit Individual*. From Fig. 5, the unit utilization $R_{\text{unit}}$ for SHA-1 is almost the same as that of DCT, but comparing the effect, the DCT achieves better reduction result. Looking at the mapping of the compiler, only particular units are used in DCT, whereas SHA-1 uses different units for executing the application. That is, units in the latter application have smaller chance to sleep satisfying BET. It suggests that the mapping and scheduling optimization by the compiler can improve the efficiency of fine grained PG.

## 5. Conclusion

Fine grained PG is applied to a DRPA, MuCCRA-2.32b, and leakage power and area overhead are evaluated. The

**Figure 7. Leakage Power Reduction Effect**

sleep control bits are inserted into configuration data, and are set when unused time is longer than BET, according to the pre-analysis results. We evaluated the effect of two control modes; *Pair* and *Unit Individual*, based on layout design and real applications. It appears that by applying PG for ALUs and SMUs individually, 48% of leakage power can be reduced with 9.0% of area overhead.

## References

[1] H.Amano and Y.Hasegawa and S.Tsutsumi and T.Nakamura and T.Nisimura and V.Tanbunheng and A.Parimala and T.Sano and M.Kato. MuCCRA Chips: Configurable Dynamically-Reconfigurable Processors. In *Proc. of ASSCC 2007*, pages 384–387, Nov. 2007.

[2] K.Kurose and et al. A 90nm Embedded DRAM Single Chip LSI with a 3D Graphics, H.264 Codec Engine, And a Reconfigurable Processor. In *Hot Chips 16*, 2004.

[3] K.Usami, N.Ohkubo. An Approach for Fine-grained Runtime Power Gating using Locally Extracted Sleep Signals. *ICCD*, Oct. 2006.

[4] T. Kodama, et al. Flexible Engine: A Dynamic Reconfigurable Accelerator with High Performance and Low Power Consumption. In *Proc. of Int'l Symp. on Low-Power and High-Speed Chips (COOL Chips)*, pages 393–408, Apr. 2006.

[5] T. Sugawara and K. Ide and T. Sato. Dynamically Reconfigurable Processor Implemented with IPFlex's DAPDNA Technology. *IEICE Trans. on Information & System*, E87-D(8):1997–2003, May 2004.

[6] Thomas Schweizer, et al. Exploiting Slack Time in Dynamically Reconfigurable Processor Architectures. In *Proc. of IEEE Int'l Conf. on Field Programmable Technology (FPT)*, pages 381–384, Dec. 2007.

[7] T.Nishimura and K.Hirai and Y.Saito and T.Nakamura and Y.Hasegawa and S.Tsutsusmi and V.Tunbunheng and H.Amano. Power Reduction Techniques for Dynamically Reconfigurable Processor Arrays. In *Proc. of Int'l Conf. on Field Programmable Logic and Application (FPL)*, Sept. 2008.

[8] X.Wnag and S.G.Ziavras and J.Hu. Energy-Aware System Synthesis for Reconfigurable Chip Multiprocessors. In *Proc. of ERSA 2007*, pages 61–68, July 2007.