

64 Bit×64 Bit Multiprecision Multiplier for Operands Scheduler with Dynamic Voltage Scaling



B.Ravi Teja
M.Tech Student,
Department of ECE,
CMR Institute of Technology,
Hyderabad, Telangana, India.



Muni Praveena Rela
Associate Professor,
Department of ECE,
CMR Institute of Technology,
Hyderabad, Telangana, India.

Abstract:

In this paper, we present a multiprecision (MP) reconfigurable multiplier that incorporates variable precision, parallel processing (PP), razor-based dynamic voltage scaling (DVS), and dedicated MP operands scheduling to provide optimum performance for a variety of operating conditions. All of the building blocks of the proposed reconfigurable multiplier can either work as independent smaller-precision multipliers or work in parallel to perform higher-precision multiplications. Given the user's requirements (e.g., throughput), a dynamic voltage/ frequency scaling management unit configures the multiplier to operate at the proper precision and frequency. Adapting to the run-time workload of the targeted application, razor flip-flops together with a dithering voltage unit then configure the multiplier to achieve the lowest power consumption. The single-switch dithering voltage unit and razor flip-flops help to reduce the voltage safety margins and overhead typically associated to DVS to the lowest level. The large silicon area and power overhead typically associated to reconfigurability features are removed. Finally, the proposed novel MP multiplier can further benefit from an operands scheduler that rearranges the input data, hence to determine the optimum voltage and frequency operating conditions for minimum power consumption. Experimental results show that the proposed MP design features a 28.2% and 15.8% reduction in circuit area and power consumption compared with conventional fixed-width multiplier. When combining this MP design with error-tolerant razor-based DVS, PP, and the proposed novel operands scheduler, 77.7%–86.3% total power reduction is achieved with a total silicon area overhead as low as 11.1%. This paper successfully demonstrates that MP architecture can allow more aggressive frequency/ supply voltage scaling for improved power efficiency.

Index Terms:

Computer arithmetic, dynamic voltage scaling, low power design, multi-precision multiplier.

I.INTRODUCTION:

Consumers demand for increasingly portable yet high-performance multimedia and communication products imposes stringent constraints on the power consumption of individual internal components [1]–[4]. Of these, multipliers perform one of the most frequently encountered arithmetic operations in digital signal processors (DSPs) [4]. For embedded applications, it has become essential to design more power-aware multipliers [4]–[13]. Given their fairly complex structure and interconnections, multipliers can exhibit a large number of unbalanced paths, resulting in substantial glitch generation and propagation [8], [11]. This spurious switching activity can be mitigated by balancing internal paths through a combination of architectural and transistor-level optimization techniques [8], [11]. In addition to equalizing internal path delays, dynamic power reduction can also be achieved by monitoring the effective dynamic range of the input operands so as to disable unused sections of the multiplier [6], [12] and/or truncate the output product at the cost of reduced precision [13]. This is possible because, in most sensor applications, the actual inputs do not always occupy the entire magnitude of its word-length. For example, in artificial neural network applications, the weight precision used during the learning phase is approximately twice that of the retrieval phase [14]. Besides, operations in lower precisions are the most frequently required. In contrast, most of today's full-custom DSPs and application-specific integrated circuits (ASICs) are designed for a fixed maximum word-length so as to accommodate the worst case scenario.

Therefore, an 8-bit multiplication computed on a 32-bit Booth multiplier would result in unnecessary switching activity and power loss. Several works investigated this word-length optimization. [1], [2] proposed an ensemble of multipliers of different precisions, with each optimized to cater for a particular scenario. Each pair of incoming operands is routed to the smallest multiplier that can compute the result to take advantage of the lower energy consumption of the smaller circuit. This ensemble of point systems is reported to consume the least power but this came at the cost of increased chip area given the used ensemble structure. To address this issue, [3], [5] proposed to share and reuse some functional modules within the ensemble. In [3], an 8-bit multiplier is reused for the 16-bit multiplication, adding scalability without large area penalty. Reference [5] extended this method by implementing pipelining to further improve the multiplier's performance. A more flexible approach is proposed in [15], with several multiplier elements grouped together to provide higher precisions and reconfigurability. Reference [7] analyzed the overhead associated to such reconfigurable multipliers. This analysis showed that around 10%–20% of extra chip area is needed for 8–16 bits multipliers. Combining multiprecision (MP) with dynamic voltage scaling (DVS) can provide a dramatic reduction in power consumption by adjusting the supply voltage according to circuit's run-time workload rather than fixing it to cater for the worst case scenario [4].

When adjusting the voltage, the actual performance of the multiplier running under scaled voltage has to be characterized to guarantee a fail-safe operation. Conventional DVS techniques consist mainly of lookup table (LUT) and on-chip critical path replica approaches [17]–[19]. The LUT approach tunes the supply voltage according to a predefined voltage–frequency relationship stored in a LUT, which is formed considering worst case conditions (process variations, power supply droops, temperature hot-spots, coupling noise, and many more). Therefore, large margins are necessarily added, which in turn significantly decrease the effectiveness of the DVS technique. The critical path replica approach typically involves an on-chip critical path replica to approximate the actual critical path. Therefore, voltage could be scaled to the extent that the replica fails to meet the timing. However, safety margins are still needed to compensate for the in-tradie delay mismatch and address fast-changing transient effects [24]. In addition, the critical path may change as a result of the varying supply voltage or process or temperature variations.

If this occurs, computations will completely fail regardless of the safety margins. The aforementioned limitations of conventional DVS techniques motivated recent research efforts into error-tolerant DVS approaches [24]–[27], which can run-time operate the circuit even at a voltage level at which timing errors occur. A recovery mechanism is then applied to detect error occurrences and restore the correct data. Because it completely removes worst case safety margins, error-tolerant DVS techniques can further aggressively reduce power consumption. In this paper, we propose a low power reconfigurable multiplier architecture that combines MP with an error-tolerant DVS approach based on razor flip-flops [25]. The main contributions of this paper can be summarized follows.

- 1) A novel MP multiplier architecture featuring, respectively, 28.2% and 15.8% reduction in silicon area and power consumption compared with its conventional 32×32 bit fixed-width multiplier counterpart. All reported multipliers trade silicon area/power consumption for MP [7]. In this paper, silicon area is optimized by applying an operation reduction technique that replaces a multiplier by adders/subtractors.
- 2) A silicon implementation of this MP multiplier integrating an error-tolerant razor-based dynamic DVS approach. The fabricated chip demonstrates run-time adaptation to the actual workload by operating at the minimum supply voltage level and minimum clock frequency while meeting throughput requirements. Prior works combining MP with DVS have only considered a limited number of off-line simulated precision-voltage pairs, with unnecessary large safety margins added to cater for critical paths [9], [10].
- 3) A novel dedicated operand scheduler that rearranges operations on input operands so as to reduce the number of transitions of the supply voltage and, in turn, minimize the overall power consumption of the multiplier. Unlike reported scheduling works, the function of the proposed.

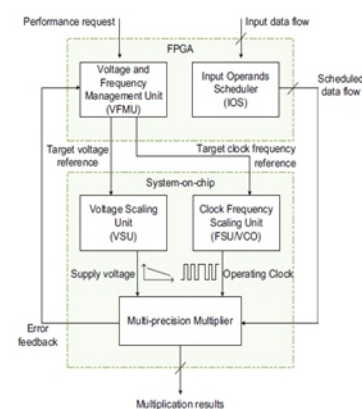


Fig 1: Overall multiplier system architecture.

Scheduler is not task scheduling rather input operands scheduling for the proposed MP multiplier. The rest of this paper is organized as follows. Section II presents the operation and architecture of the proposed MP multiplier. Section III presents the approach used to reduce the overhead associated to MP and reconfigurability. Section IV presents the operating principle and implementation of the DVS management unit. Section V presents experimental results. Section VI presents the operands scheduler unit. Finally, a conclusion is given in Section VII.

II. SYSTEM OVERVIEW AND OPERATION:

The proposed MP multiplier system (Fig. 1) comprises five different modules that are as follows:

- 1) The MP multiplier;
- 2) The input operands scheduler (IOS) whose function is to reorder the input data stream into a buffer, hence to reduce the required power supply voltage transitions;
- 3) The frequency scaling unit implemented using a voltage-controlled oscillator (VCO). Its function is to generate the required operating frequency of the multiplier;
- 4) The voltage scaling unit (VSU) implemented using a voltage dithering technique to limit silicon area overhead. Its function is to dynamically generate the supply voltage so as to minimize power consumption;
- 5) The dynamic voltage/frequency management unit (VFMU) that receives the user requirements (e.g., throughput).

The VFMU sends control signals to the VSU and FSU to generate the required power supply voltage and clock-frequency for the MP multiplier. The MP multiplier is responsible for all computations. It is equipped with razor flip-flops that can report timing errors associated to insufficiently high voltage supply levels. The operation principle is as follows. Initially, the multiplier operates at a standard supply voltage of 3.3 V. If the razor flip-flops of the multiplier do not report any errors, this means that the supply voltage can be reduced. This is achieved through the VFMU, which sends control signals to the VSU, hence to lower the supply voltage level. When the feedback provided by the razor flip-flops indicates timing errors, the scaling of the power supply is stopped. The proposed multiplier (Fig. 2) not only combines MP and DVS but also parallel processing (PP). Our multiplier-comprises 8×8 bit reconfigurable multipliers.

These building blocks can either work as nine independent multipliers or a single 32×32 bit work in parallel to perform one, two or three 16×16 bit multiplications or a single 32×32 bit operation. PP can be used to increase the throughput or reduce the supply voltage level for low power operation.

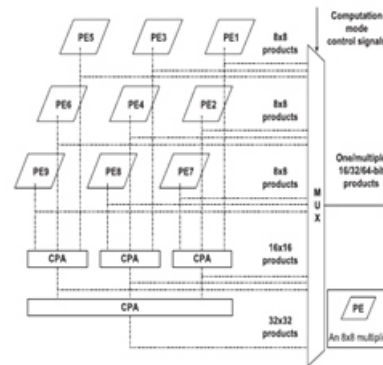


Fig 2: Possible configuration modes of proposed MP multiplier.

III. MP AND RECONFIGURABILITY OVERHEAD:

Fig.3 shows the structure of the input interface unit, which is a sub-module of the MP multiplier (Fig. 1). The role of this input interface unit (Fig. 4) is to distribute the input data between the nine independent processing elements (PEs) (Fig. 2) of the 32×32 bit MP multiplier, considering the selected operation mode. The input interface unit uses an extra MSB sign bit to enable both signed and unsigned multiplications. A 3-bit control bus indicates whether the inputs are 1/4/9 pair(s) of 8-bit operands, or 1/2/3 pair(s) of 16-bit operands, or 1 pair of 32-bit operands, respectively. Depending on the selected operating mode, the input data stream is distributed (Fig. 4) between the PEs to perform the computation. Fig. 5 shows how three 8×8 bit PEs are reused to realize a 16×16 bit multiplier. The 32×32 bit multiplier is constructed using a similar approach but requires 3×3 PEs. A 3-bit control word defines which PEs work concurrently and which PEs are disabled. Whenever the full precision (32×32 bit) is not exercised, the supply voltage and the clock frequency may be scaled down according to the actual workload. To evaluate the overhead associated to reconfigurability and MP, we define X and Y as the $2n$ -bits wide multiplicand and multiplier, respectively. XH, YH are their respective n most significant bits whereas XL, YL are their respective n least significant bits. XLYL, XHYL, XLYH, XHYH is the crosswise products. The product of X and Y can be expressed as follows:

$P = (X_H Y_H)_{2n} + (X_H Y_L + X_L Y_H)_{2n} + X_L Y_L(1)$
 Where $2n$ -bit reconfigurable multiplier can be built using adders and four n bit \times n bit multipliers to compute $X_H Y_H$, $X_H Y_L$, $X_L Y_H$, and $X_L Y_L$. this would result in overheads of 18% and 13% for the silicon area and power, respectively. However, if we define [18]

$$X_1 = X_H + X_L \quad (2)$$

$$Y_1 = Y_H + Y_L \quad (3)$$

then (1) could be rewritten as follows

$$P = (X_1 Y_1)_{2n+2} - (X_H Y_H - X_L Y_L)_{2n} + X_L Y_L(4)$$

Comparing (1) and (4), we have removed one $n \times n$ bit multiplier (for calculating $X_H Y_L$ or $X_L Y_H$) and one $2n$ -bit adder (for calculating $X_H Y_L + X_L Y_H$). The two adders are replaced with two n -bit adders (for calculating $X_H + X_L$ and $Y_H + Y_L$) and two $(2n + 2)$ -bit subtractors (for calculating $X_1 Y_1 - X_H Y_H - X_L Y_L$).

In a 32-bit multiplier, we can thus significantly reduce the design complexity by using two 34-bit subtractors to replace a 16×16 bit multiplier. We actually need two 16×16 bit multipliers (for calculating $X_H Y_H$ and $X_L Y_L$) and one 17×17 bit multiplier (for calculating $X_1 Y_1$).

To evaluate the proposed MP architecture, a conventional 32-bit fixed-width multiplier and four sub-block MP multipliers are designed using a Booth Radix-4 Wallace tree structure similar to that used for the building blocks of our MP three sub-block multiplier. These multipliers are synthesized using the synopsis design compiler with AMIS 0.35- μ m complementary metal-oxide-semiconductor (CMOS) standard cell technology library. The power simulations are performed at a clock frequency of 50 MHz and at a power supply of 3.3 V.

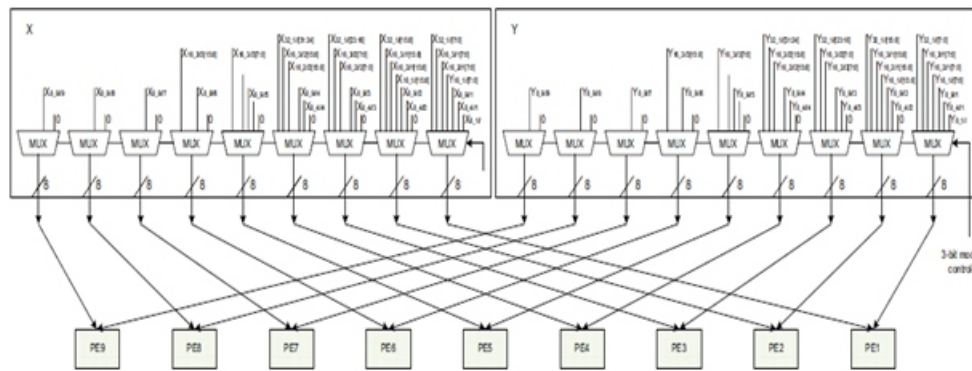


Fig.3: Structure of input interface unit

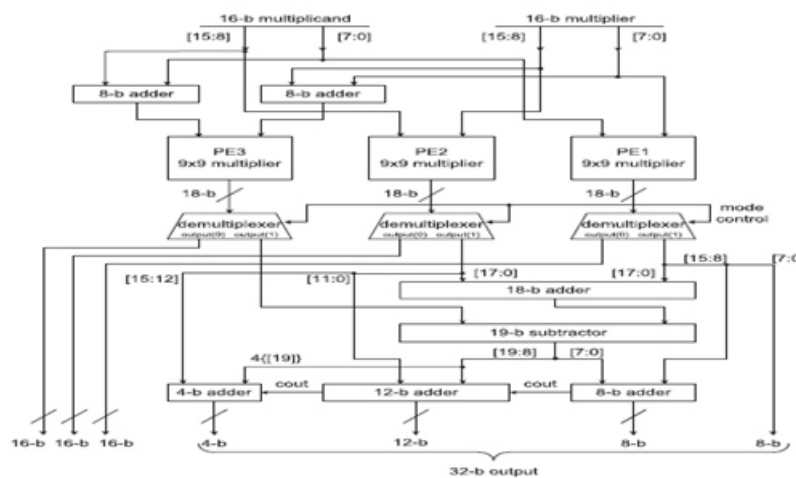


Fig.4: Three PEs combined to form 16 x 16 bit multiplier.

IV. DYNAMIC VOLTAGE AND FREQUENCY SCALING MANAGEMENT

A. DVS Unit:

In our implementation (Fig. 1), a dynamic power supply and a VCO are employed to achieve real-time dynamic voltage and frequency scaling under various operating conditions. In [28], near-optimal dynamic voltage scaling can be achieved when using voltage dithering, which exhibits faster response time than conventional voltage regulators. Voltage dithering uses power switches to connect different supply voltages to the load, depending on the time slots. Therefore, an intermediate average voltage is achieved. This conventional voltage dithering technique has some limitations. If the power switches are toggled with overlapping periods, switches can be turned on simultaneously, giving rise to a large transient current. To mitigate this, nonoverlapping clocks could be used to control power switches. However, this may result in system instability as there are instances where all supply voltages are disconnected from the load.

The requirement for multiple supplies can also result in system overhead. To address these issues, we implemented a single-supply voltage dithering scheme [Fig. 6(a)], which operates as follows. When the supply voltage (V_n) of the multiplier drops below the predefined reference voltage (V_{ref}), the comparator output (V_a) toggles. Therefore, the VFMU turns on the power switch via V_{ctrl} , for a predefined duration $T_c = 5 \mu s$. The chosen value for the off-chip storage capacitor C_s is $4.7 \mu F$. This value is chosen to achieve a voltage ripple magnitude of 50 mV [Fig. 6(b)] with a charging current set to 50 mA, hence to limit the resistive power loss of the dithering unit to less than 1% of the total power consumption. The value of C_s is a tradeoff between ripple magnitude, tracking speed, and area/power overheads. Fig. 6(b) shows experimental results for the voltage control loop.

B. Dynamic Frequency Scaling Unit:

In the proposed 32×32 bit MP multiplier, dynamic frequency tuning is used to meet throughput requirements. It is based on a VCO implemented as a seven-stage current starved ring oscillator. The VCO output frequency can be tuned from 5 to 50 MHz using four control bits (5 MHz/step). This frequency range is selected to meet the requirements of general purpose DSP applications.

The reported multiplier can operate as a 32-bit multiplier or as nine independent 8-bit multipliers. For the chosen 5–50 MHz operating range, our multiplier boasts up to $9 \times 50 = 450$ MIPS. The simulated power consumption for the VCO ranges from 85 (5 MHz) to $149 \mu W$ (50 MHz), which is negligible compared with the power consumed by the multiplier. Fig. 7 shows experimental measurements showing the transient response for the worst case frequency switching (from 50 to 5 MHz). Clock frequency can settle within one clock cycle as required.

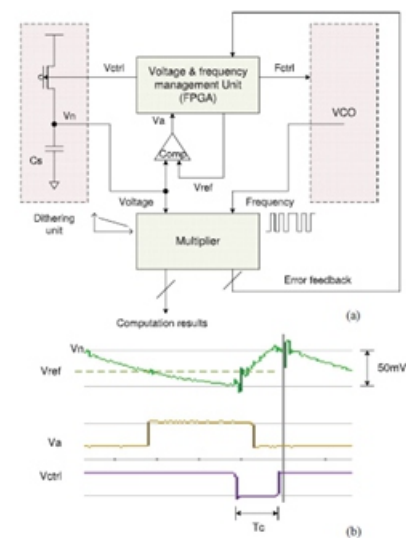


Fig. 6. (a) Proposed single-header voltage dithering unit and voltage and frequency tuning loops. (b) Experimental timing results from voltage dithering unit.

V. INPUT OPERANDS SCHEDULER:

Here, we present three different algorithms to reduce this overall power consumption. Each of these algorithms constitutes a different approach to process the mixed-precision data held in the operands buffer (Fig. 5). The performance of each algorithm is evaluated using a mixed precision data set of 120 000 randomly operands, with a third corresponding to each precision (8, 16, 32 and 64-bit). In the following, the specified throughput T_p for the proposed 64×64 bit multiplier is $64 F$ (Mbits/s), where F is the multiplier's operating frequency.

Algorithm A:

In the first algorithm, the multiplier throughput $T_p = 64 F$ is kept constant by fixing the operating frequencies (f_{32} , f_{16} , or f_8) of each precision-data group (32, 16, or 8-bit) to $f_{32} = F$, $f_{16} = F/2$, $f_8 = F/4$ (10)

Where F is the multiplier's operating frequency. This is because the throughput in 8×8 bit multiplication mode is four times that of the 32×32 bit multiplication mode and double that of the 16×16 bit multiplication mode, as a result of the multiplier PP. The minimum supply voltage (V_{min32} , V_{min16} or V_{min8}) associated to each operating frequency (f_{32} , f_{16} or f_8) is determined through a V_{min} - f LUT. Algorithm A shows its limitations when 32-bit operands are processed initially all N_{32} operands of the data block are processed,

the supply voltage (V_n) needs to decrease rapidly from point A (V_{min32}) to point B (V_{min16}) at which all N_{16} 16-bit operands of the data block should be processed. If N_{16} is too small, most 16-bit operands will be actually processed in Sections A and B, that is at a voltage possibly much higher than the minimal V_{min16} level. Similarly 8-bit operands of the data block could be processed in Sections C and D, B-C, or even A-B for the worst case. This contributes to increasing $P_{comp_overhead}$.

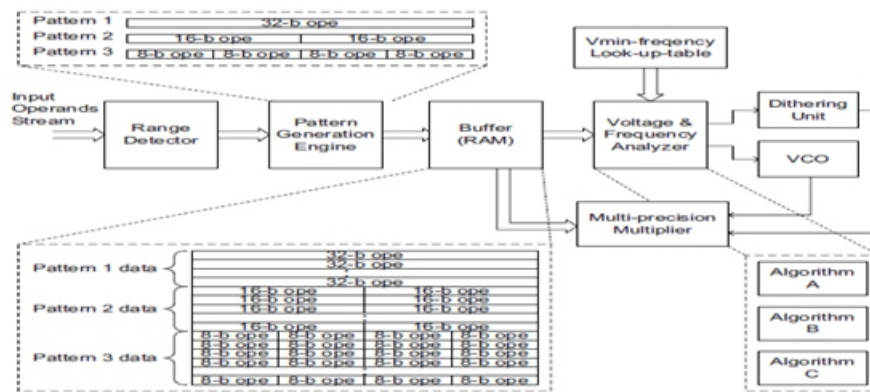


Fig.5: Block diagram of IOS

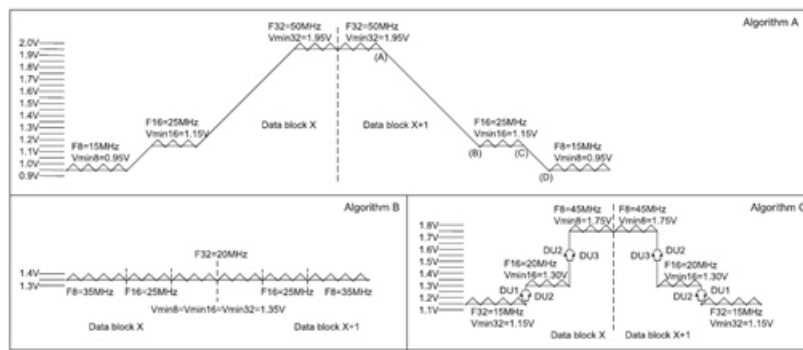


Fig.6: Operation principles of operand scheduling algorithms A, B, and C. Data Block X and Data Block X+1 refer to two-consecutive operand datablocks subsequently stored into the RAM, respectively.

Algorithm B:

This algorithm removes all transitions of the power supply voltage by making V_{min32} , V_{min16} , and V_{min8} equal and adjusting f_{32} , f_{16} , and f_8 such that the overall throughput is kept unchanged. We thus need to have the following:

$$\frac{64N_{32} + 128N_{16} + 256N_8}{\frac{N_{32}}{f_{32}} + \frac{N_{16}}{f_{16}} + \frac{N_8}{f_8}} = 64F \quad (11)$$

From a LUT, we can obtain the V_{min} - f relationship as follows:

$$V_{min32} = \psi_{32}(f_{32}) \quad (12)$$

$$V_{min16} = \psi_{16}(f_{16}) \quad (13)$$

$$V_{min8} = \psi_8(f_8) \quad (14)$$

As algorithm B keeps the supply voltage constant

$$\psi_{32}(f_{32}) = \psi_{16}(f_{16}) = \psi_8(f_8) = V \quad (15)$$

the operating frequencies f_{32} , f_{16} , and f_8 can be determined by using (11) and (15). For example, when F is set to 50 MHz, the values for V , f_{32} , f_{16} , and f_8 are found to be 1.35 V, 20 MHz, 25 MHz, and 35 MHz, respectively.

Algorithm C:

Although Algorithm B removes power supply voltage transitions by setting a single-voltage level V , there may be better power saving combinations of power supply voltages and operating frequencies: (V_{min32}, f_{32}) , (V_{min16}, f_{16}) , and (V_{min8}, f_8) . The aim of algorithm C is to find such an optimum for reduced power consumption. To limit complexity, we will only seek to minimize the dynamic power dissipated as a result of the computation

$$P = CV^2f$$

$$= C_m 32V^2 \min_{32} f_{32} + C_m 16V^2 \min_{16} f_{16} + m 8V^2 \min_{8} f_8$$

$$= \chi(f_{32}, f_{16})$$

Given that the V_{min} - f relationships are known (12)–(14), one could find the minimum of the above equation for the specified throughput (11). For example, when F is set to 50 MHz, the values for (V_{min32}, f_{32}) , (V_{min16}, f_{16}) , (V_{min8}, f_8) are found to be (1.15 V, 15 MHz), (1.30 V, 20 MHz), and (1.75 V, 45 MHz), respectively. When considering DVS, razor, RAM, and dedicated scheduling circuitry, algorithm B exhibits the least power consumption, with an overall power reduction of 86.3%, compared with the standard 32×32 bit fixed-width multiplier. However, it requires two additional dithering units to generate all three discrete power supply levels V_{min32} , V_{min16} , and V_{min8} and thus remove transitions among these different supply levels. This increases the total silicon area overhead to 27.1%. Therefore, algorithm B provides the most attractive tradeoff with 81.5% reduction and a silicon area overhead of just 11.9%.

VI SIMULATION RESULTS:

The simulation of the proposed design is carried out by using Xilinx software. The simulated waveforms are shown in below figure.

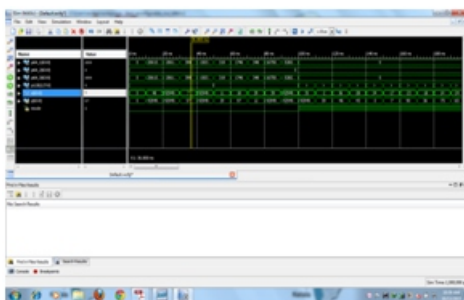


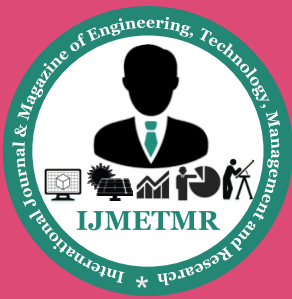
Fig.7: Simulation results of the proposed 64 x 64 bit design in signed decimal

VII CONCLUSION:

We proposed a novel MP multiplier architecture featuring, respectively, 28.2% and 15.8% reduction in silicon area and power consumption compared with its 32×32 bit conventional fixed-width multiplier counterpart. When integrating this MP multiplier architecture with an error-tolerant razor-based DVS approach and the proposed novel operands scheduler, 77.7%–86.3% total power reduction was achieved with a total silicon area overhead as low as 11.1%. The fabricated chip demonstrated run-time adaptation to the actual workload by operating at the minimum supply voltage level and minimum clock frequency while meeting throughput requirements. The proposed novel dedicated operand scheduler rearranges operations on input operands, hence to reduce the number of transitions of the supply voltage and, in turn, minimized the overall power consumption of the multiplier. The proposed MP razor-based DVS multiplier provided a solution toward achieving full computational flexibility and low power consumption for various general purpose low-power applications.

REFERENCES:

- [1] R. Min, M. Bhardwaj, S.-H. Cho, N. Ickes, E. Shih, A. Sinha, A. Wang, and A. Chandrakasan, "Energy-centric enabling technologies for wireless sensor networks," *IEEE Wirel. Commun.*, vol. 9, no. 4, pp. 28–39, Aug. 2002.
- [2] M. Bhardwaj, R. Min, and A. Chandrakasan, "Quantifying and enhancing power awareness of VLSI systems," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 9, no. 6, pp. 757–772, Dec. 2001.
- [3] A. Wang and A. Chandrakasan, "Energy-aware architectures for a real-valued FFT implementation," in *Proc. IEEE Int. Symp. Low Power Electron. Design*, Aug. 2003, pp. 360–365.
- [4] T. Kuroda, "Low power CMOS digital design for multimedia processors," in *Proc. Int. Conf. VLSI CAD*, Oct. 1999, pp. 359–367.
- [5] H. Lee, "A power-aware scalable pipelined booth multiplier," in *Proc. IEEE Int. SOC Conf.*, Sep. 2004, pp. 123–126.



[6] S.-R. Kuang and J.-P. Wang, "Design of power-efficient configurable booth multiplier," *IEEE Trans. Circuits Syst. I, Reg. Papers*, vol. 57, no. 3, pp. 568–580, Mar. 2010.

[7] O. A. Pfander, R. Hacker, and H.-J. Pflöderer, "A multiplexer-based concept for reconfigurable multiplier arrays," in *Proc. Int. Conf. Field-Programm. Logic Appl.*, vol. 3203, Sep. 2004, pp. 938–942.

[8] F. Carbognani, F. Buerger, N. Felber, H. Kaeslin, and W. Fichtner, "Transmission gates combined with level-restoring CMOS gates reduce glitches in low-power low-frequency multipliers," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 16, no. 7, pp. 830–836, Jul. 2008.

Authors Profile :

Bottu Ravi Teja is currently pursuing his M.Tech specialization in VLSI system design in CMR Institute of Technology which is affiliated to JNTUH in Hyderabad.

Muni Praveena Rela is currently working as Associate Professor in Department Of ECE in CMR Institute of Technology which is affiliated to JNTUH in Hyderabad.