# Design of An Efficient Aging-Aware Reliable Multiplier Based On Adaptive Hold Logic

### T. Shireesha[1], P.P. Nagaraja Rao[2]

[1]M. Tech in VLSI System Design, Department of ECE, Sri Venkatesa Perumal College of Engineering & Technology, Puttur, Chittoor (dt), A.P.

[2]Associate Professor & HOD, Department of ECE, Sri Venkatesa Perumal College of Engineering & Technology, Puttur, Chittoor (dt), A.P.

**[1]Shireeshat3@gmail.com,** [2]**nagaraj9s@gmail.com**

*Abstract—* **Digital multipliers are among the most critical arithmetic functional units. The overall performance of the Digital multiplier systems depends on throughput of the multiplier. The negative bias temperature instability effect occurs when a pMOS transistor is under negative bias (Vgs = −Vdd), increasing the threshold voltage of a pMOS transistor and reducing the multiplier speed. Similarly, positive bias temperature instability occurs when an nMOS transistor is under positive bias. Both effects degrade the speed of the transistor and in the long term, the system may be fail due to timing violations. Therefore, it is required to design reliable high-performance multipliers. In this paper, we implement an aging-aware multiplier design with a novel adaptive hold logic (AHL) circuit. The multiplier is able to provide the higher throughput through the variable latency and can adjust the adaptive hold logic (AHL) circuit to lessen performance degradation that is due to the aging effect. The proposed design can be applied to the column bypass multiplier.**

*Keywords—Adaptive hold logic (AHL), negative bias temperature instability (NBTI), positive bias temperature instability(PBTI), reliable multiplier, variable latency.*

## 1. INTRODUCTION

Digital multipliers are among the most critical arithmetic functional units in many applications, such as the Fourier transform, discrete cosine transforms, and digital filtering. The through put of these applications depends on multipliers, if the multipliers are too slow, the performance of entire circuits will be reduced. Furthermore, negative bias temperature instability (NBTI) occurs when a pMOS transistor is under negative bias (Vgs = −Vdd). In this situation, the interaction between inversion layer holes and hydrogen-passivated Si atoms breaks the Si–H bond generated during the oxidation process, generating H or H2 molecules. When these molecules diffuse away, interface traps are left. The accumulated interface traps between silicon and the gate oxide interface result in increased threshold voltage (Vth), reducing the circuit switching speed. When the biased voltage is removed, the reverse reaction occurs, reducing the NBTI effect. However, the reverse reaction does not eliminate all the interface traps generated during the stress phase, and Vth is increased in the long term. Hence, it is important to design a reliable high-performance multiplier. The corresponding effect on an nMOS transistor is positive bias temperature instability (PBTI), which occurs when an nMOS transistor is under positive bias. Compared with the NBTI effect, the PBTI effect is much smaller on oxide/polygate transistors, and therefore is usually ignored. However, for high-k/metal-gate nMOS transistors with significant charge trapping, the PBTI effect can no longer be ignored. In fact, it has been shown that the PBTI effect is more significant than the NBTI effect on 32-nm high-k/metal-gate processes.

A traditional method to mitigate the aging effect is overdesign including such things as guard-banding and gate over sizing; however, this approach can be very pessimistic and area and power inefficient. To avoid this problem, many NBTI-aware methodologies have been proposed. An NBTI-aware technology mapping technique was proposed into guarantee the performance of the circuit during its lifetime. In an NBTI-aware sleep transistor was designed to reduce the aging effects on pMOS sleep-transistors, and the lifetime stability of the power-gated circuits under consideration was improved. Wu and Marculescu proposed a point logic restructuring and pin reordering method, which is based on detecting functional symmetries and transistor stacking effects. They also proposed an NBTI optimization method that considered path sensitization. In dynamic voltage scaling and body-basing techniques were proposed to reduce power or extend circuit life. These techniques, however, require circuit modification or do not provide optimization of specific circuits. Traditional circuits use critical path delay as the overall circuit clock cycle in order to perform correctly. However, the probability that the critical paths are activated is low. In most cases, the path delay is shorter than the critical path. For these noncritical paths, using the critical path delay as the overall cycle period will result in significant timing waste. Hence, the Variable-latency design

was proposed to reduce the timing waste of traditional circuits.

The variable-latency design divides the circuit into two parts: 1) shorter paths and 2) longer paths. Shorter paths can execute correctly in one cycle, whereas longer paths need two cycles to execute. When shorter paths are activated frequently, the average latency of variable-latency designs is better than that of traditional designs. For example, several variable-latency adders were proposed using the speculation technique with error detection and recovery. A short path activation function algorithm was proposed into improve the accuracy of the hold logic and to optimize the performance of the variable-latency circuit. An instruction scheduling algorithm was proposed into schedule the operations on non-uniform latency functional units and improves the performance of Very Long Instruction Word processors. In a variable-latency pipelined multiplier architecture with a Booth algorithm was proposed. In process-variation tolerant architecture for arithmetic units was proposed, where the effect of process-variation is considered to increase the circuit yield. In addition, the critical paths are divided into two shorter paths that could be unequal and the clock cycle is set to the delay of the longer one. These research designs were able to reduce the timing waste of traditional circuits to improve performance, but they did not consider the aging effect and could not adjust themselves during the runtime. A variable-latency adder design that considers the aging effect was proposed. However, no variable-latency multiplier design that considers the aging effect and can adjust dynamically has been done.

This paper has been organized in the following way, we propose an aging-aware reliable multiplier design with novel adaptive hold logic (AHL) circuit. The multiplier is based on the variable-latency technique and can adjust the AHL circuit to achieve reliable operation under the influence of NBTI and PBTI effects. To be specific, the contributions of this paper are summarized as follows: 1) novel variable latency multiplier architecture with an AHL circuit. The AHL circuit can decide whether the input patterns require one or two cycles and can adjust the judging criteria to ensure that there is minimum performance degradation after considerable aging occurs; 2) comprehensive analysis and comparison of the multiplier's performance under different cycle periods to show the effectiveness of our proposed architecture; 3) an aging-aware reliable multiplier design method that is suitable for large multipliers. Although the experiment is performed in 16-and 32-bit multipliers, our proposed architecture can be easily extended to large designs; 4) the experimental results show that our proposed architecture with the $16 \times 16$ and $32 \times 32$ column-bypassing multipliers can attain up to 62.88% and 76.28% performance improvement compared with the $16 \times 16$ and $32 \times 32$ fixed-latency column-bypassing (FLCB) multipliers. In addition, our proposed architecture with $16 \times 16$ and $32 \times 32$ row-bypassing multipliers can achieve up to 80.17% and 69.40% performance improvement as compared with $16 \times 16$ and $32 \times 32$ fixed-latency row-bypassing multipliers.

## II. PRELIMINARIES

### A. Column-Bypassing Multiplier

A column-bypassing multiplier is an improvement on the normal array multiplier (AM). The AM is a fast parallel AM and is shown in Fig. 1. The multiplier array consists of $(n-1)$ rows of carry save adder (CSA), in which each row contains $(n-1)$ full adder (FA) cells. Each FA in the CSA array has two outputs: 1) the sum bit goes down and 2) the carry bit goes to the lower left FA. The last row is a ripple adder for carry propagation.
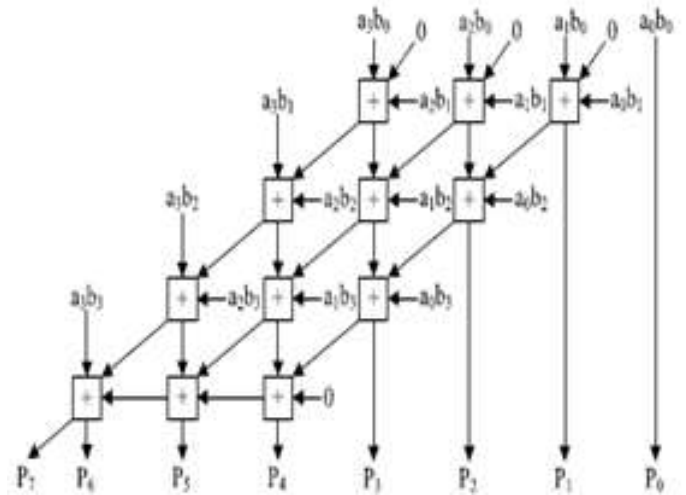


Fig. 1.4 × 4 normal AM.

The FAs in the AM are always active regardless of input states. A low-power column-bypassing multiplier design is proposed in which the FA operations are disabled if the corresponding bit in the multiplicand is 0. Fig. 2 shows a 4×4 column-bypassing multiplier. Supposing the inputs are $1010_2 * 1111_2$, it can be seen that for the FAs in the first and third diagonals, two of the three input bits are 0: the carry bit from its upper right FA and the partial product $a_ib_i$. Therefore, the output of the adders in both diagonals is 0, and the output sum bit is simply equal to the third bit, which is the sum output of its upper FA.
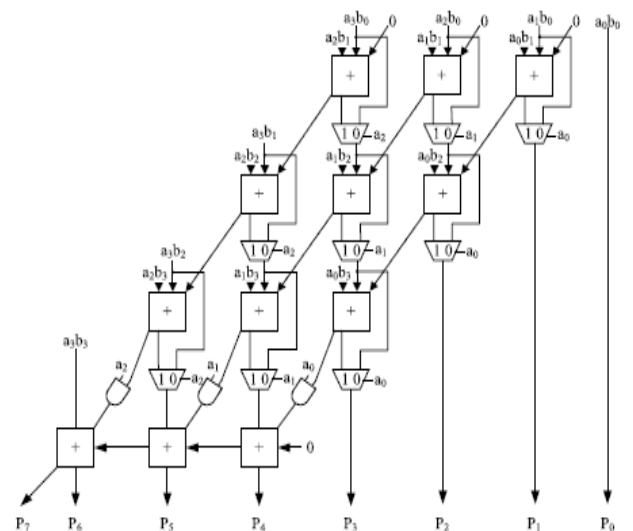


Fig. 2. 4 × 4 column-bypassing multiplier

Hence, the FA is modified to add two tristate gates and one multiplexer. The multiplicand bit $ai$ can be used as the selector of the multiplexer to decide the output of the FA, and $ai$ can also be used as the selector of the tristate gate to turn off the input path of the FA. If $ai$ is 0, the inputs of FA are disabled, and the sum bit of the current FA is equal to the sum bit from its upper FA, thus reducing the power consumption of the multiplier. If $ai$ is 1, the normal sum result is selected.

## B. Row-Bypassing Multiplier

A low-power row-bypassing multiplier is also proposed to reduce the activity power of the AM. The operation of the low-power row-bypassing multiplier is similar to that of the low-power column-bypassing multiplier, but the selector of the multiplexers and the tristate gates use the multiplicator.
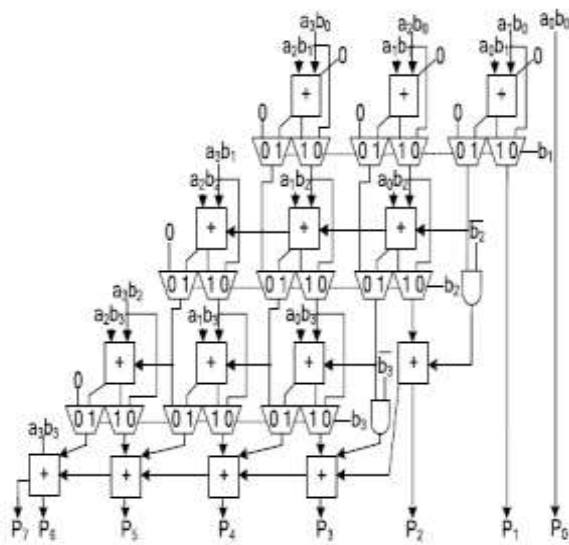


Fig. 3. 4 × 4 row-bypassing multiplier

Fig. 3 is a 4 × 4 row-bypassing multiplier. Each input is connected to an FA through a tristate gate. When the inputs are $1111_2 * 1001_2$, the two inputs in the first and second rows are 0 for FAs. Because $b1$ is 0, the multiplexers in the first row select $aib0$ as the sum bit and select 0 as the carry bit. The inputs are bypassed to FAs in the second rows, and the tristate gates turn off the input paths to the FAs. Therefore, no switching activities occur in the first-row FAs; in return, power consumption is reduced. Similarly, because $b2$ is 0, no switching activities will occur in the second-row FAs. However, the FAs must be active in the third row because the $b3$ is not zero.

## C. Variable-Latency Design

Section 1 mentioned that the variable-latency design was proposed to reduce the timing waste occurring in traditional circuits that use the critical path cycle as an execution cycle period. The basic concept is to execute a shorter path using a shorter cycle and longer path using two cycles. Since most paths execute in a cycle period that is much smaller than the critical path delay, the variable-latency design has smaller average latency.

## D. Aging Model

As mentioned in Section I, the NBTI (PBTI) effect occurs when a pMOS (nMOS) transistor is under negative (positive) bias voltage, resulting in $V$th drift. When the bias voltage is removed, the recovery process occurs, reducing the $V$th drift If a pMOS (nMOS) transistor is under constant stress, this is referred to as static NBTI (PBTI). If both stress and recovery phases exist, it is referred to as dynamic NBTI (PBTI). The $V$th drift of pMOS (nMOS) transistor due to the static NBTI (PBTI) effect can be described by dc reaction-diffusion (RD) framework. If transistors are under alternative stress and recovery phases, the dc RD model should be modified to an ac RD model.

$$\Delta V_{th}(t) \cong K_{AC} \times t^n \cong \alpha(S,f) \times K_{DC} \times t^n$$

where $\alpha$ is a function of stress frequency ($f$) and signal probability ($S$). Since the impact of frequency is relatively insignificant, the effect of signal frequency is ignored. $K$DC is a technology-dependent constant.

$$K_{DC} = A \times T_{OX} \times \sqrt{C_O} \times (V_{GS} - V_{th}) \times [1 - V_{ds}/\alpha(V_{GS} - V_{th})] \times \exp(E_{ox}/E_O) \times \exp(-E_a/kT)$$

where $A$ is a constant, and $T$OX is the oxide thickness. $E$OX is the gate electric field, which is $(V$GS–$V$th$)/T$OX; $k$ is the Boltzmann constant, and $T$ is the temperature. $E$0 and $Ea$ are technology-independent characteristics of the reaction that are equal to 1.9–2.0 MV/cm and 0.12 eV, respectively.

## III. PROPOSED AGING-AWARE MULTIPLIER

This section details the proposed aging-aware reliable multiplier design. It introduces the overall architecture and the functions of each component and also describes how to design AHL that adjusts the circuit when significant aging occurs.

## A. Proposed Architecture

The below Figure shows our proposed aging-aware multiplier architecture, which includes two $m$-bit inputs ($m$ is a positive number), one 2$m$-bit output, one column- or row-bypassing multiplier, 2$m$ 1-bit Razor flip-flops, and an AHL circuit.
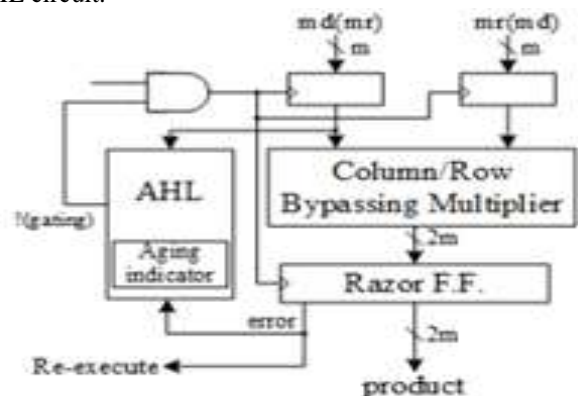


Fig.4 Proposed architecture (md means multiplicand; mr means multiplicator).

In the proposed architecture, the column- and row-bypassing multipliers can be examined by the number of zeros in either the multiplicand or multiplicator to predict whether the operation requires one cycle or two cycles to complete. When input patterns are random, the number of zeros and ones in the multiplicator and multiplicand follows a normal distribution. According to the bypassing selection in the column or row-bypassing multiplier, the input signal of the AHL in the architecture with the column-bypassing multiplier is the multiplicand, whereas that of the row-bypassing multiplier is the multiplicator. Razor flip-flops can be used to detect whether timing violations occur before the next input pattern arrives.

Fig.5 shows the details of Razor flip-flops. A 1-bit Razor flip-flop contains a main flip-flop, shadow latch, XOR gate, and mux. The main flip-flop catches the execution result for the combination circuit using a normal clock signal, and the shadow latch catches the execution result using a delayed clock signal, which is slower than the normal clock signal. If the latched bit of the shadow latch is different from that of the main flip-flop, this means the path delay of the current operation exceeds the cycle period, and the main flip-flop catches an incorrect result. If errors occur, the Razor flip-flop will set the error signal to 1 to notify the system to re-execute the operation and notify the AHL circuit that an error has occurred. We use Razor flip-flops to detect whether an operation that is considered to be a one-cycle pattern can really finish in a cycle. If not, the operation is re-executed with two cycles. Although the re-execution may seem costly, the overall cost is low because the re-execution frequency is low.
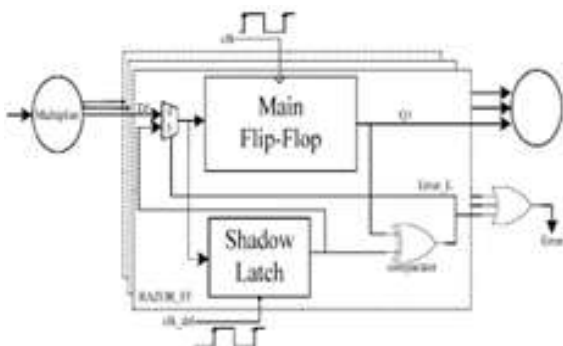


Fig.5. Razor flip flops

The AHL circuit is the key component in the aging-ware variable-latency multiplier. Fig. 6 shows the details of the AHL circuit. The AHL circuit contains an aging indicator, two judging blocks, one mux, and one D flip-flop. The aging indicator indicates whether the circuit has suffered significant performance degradation due to the aging effect. The aging indicator is implemented in a simple counter that counts the number of errors over a certain amount of operations and is reset to zero at the end of those operations. If the cycle period is too short, the column- or row-bypassing multiplier is not able to complete these operations successfully, causing timing violations. These timing violations will be caught by the Razor flip-flops, which generate error signals. If errors happen frequently and exceed a predefined threshold, it means the circuit has suffered significant timing degradation due to the aging effect, and the aging indicator will output signal 1;

otherwise, it will output 0 to indicate the aging effect is still not significant, and no actions are needed.
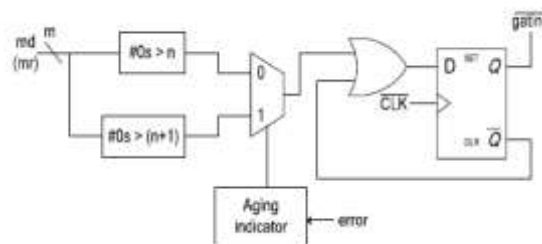


Fig. 6. Diagram of AHL (md means multiplicand; mr means multiplicator).

The details of the operation of the AHL circuit are as follows: when an input pattern arrives, both judging blocks will decide whether the pattern requires one cycle or two cycles to complete and pass both results to the multiplexer. The multiplexer selects one of either result based on the output of the aging indicator. Then an OR operation is performed between the result of the multiplexer, and the .Q signal is used to determine the input of the D flip-flop. When the pattern requires one cycle, the output of the multiplexer is 1. The !(gating) signal will become 1, and the input flip flops will latch new data in the next cycle. On the other hand, when the output of the multiplexer is 0, which means the input pattern requires two cycles to complete, the OR gate will output 0 to the D flip-flop. Therefore, the !(gating) signal will be 0 to disable the clock signal of the input flip-flops in the next cycle. Note that only a cycle of the input flip-flop will be disabled because the D flip-flop will latch 1 in the next cycle.

The overall flow of our proposed architecture is as follows: when input patterns arrive, the column- or row-bypassing multiplier, and the AHL circuit execute simultaneously. According to the number of zeros in the multiplicand (multiplicator), the AHL circuit decides if the input patterns require one or two cycles. If the input pattern requires two cycles to complete, the AHL will output 0 to disable the clock signal of the flip-flops. Otherwise, the AHL will output 1 for normal operations. When the column- or row-bypassing multiplier finishes the operation, the result will be passed to the Razor flip-flops. The Razor flip-flops check whether there is the path delay timing violation. If timing violations occur, it means the cycle period is not long enough for the current operation to complete and that the execution result of the multiplier is incorrect. Thus, the Razor flip-flops will output an error to inform the system that the current operation needs to be re-executed using two cycles to ensure the operation is correct. In this situation, the extra re-execution cycles caused by timing violation incurs a penalty to overall average latency. However, our proposed AHL circuit can accurately predict whether the input patterns require one or two cycles in most cases. Only a few input patterns may cause a timing variation when the AHL circuit judges incorrectly. In this case, the extra re-execution cycles did not produce significant timing degradation.

IV.SIMULATION RESULTS

The below figure shows the Block diagram of the aging-aware multiplier.



Fig7: Block diagram of the aging-aware multiplier

The below figure shows the RTL Schematic of the 64 x 64 aging-aware multiplier.



Fig8: RTL Schematic of the aging-aware multiplier

The below figure shows the Technology Schematic of the 64x 64 aging-aware multiplier.



Fig9: Technology Schematic of the aging-aware multiplier.

The below figure shows the Simulation output waveform of the 64 x 64 aging-aware multiplier.



Fig10: Simulation output waveform of the aging-aware multiplier.

## V. CONCLUSION

This paper proposed an aging-aware variable-latency multiplier design with the AHL. The multiplier is able to adjust the AHL to mitigate performance degradation due to increased delay. our proposed variable latency multipliers have less performance degradation because variable latency multipliers have less timing waste, but traditional multipliers need to consider the degradation caused by both the BTI effect and electromigration and use the worst case delay as the cycle period. The experimental results show that our proposed architecture with 4x4 multiplication with CLA as last stage instead of Normal RCA adder it will decrease the delay and improve the performance compared with previous designs.

## REFERENCES

[1] Y. Cao. (2013). Predictive Technology Model (PTM) and NBTI Model [Online]. Available: http://www.eas.asu.edu/~ptm

[2] S. Zafar et al., "A comparative study of NBTI and PBTI (charge trapping) in SiO2/HfO2 stacks with FUSI, TiN, Re gates," in Proc. IEEE Symp. VLSI Technol. Dig. Tech. Papers, 2006, pp. 23–25.

[3] S. Zafar, A. Kumar, E. Gusev, and E. Cartier, "Threshold voltage instabilities in high-k gate dielectric stacks," IEEE Trans. Device Mater. Rel., vol. 5, no. 1, pp. 45–64, Mar. 2005.

[4] H.-I. Yang, S.-C. Yang, W. Hwang, and C.-T. Chuang, "Impacts of NBTI/PBTI on timing control circuits and degradation tolerant design in nanoscale CMOS SRAM," IEEE Trans. Circuit Syst., vol. 58, no. 6,pp. 1239–1251, Jun. 2011.

[5] R. Vattikonda, W. Wang, and Y. Cao, "Modeling and miimization of pMOS NBTI effect for robust naometer design," in Proc. ACM/IEEE DAC, Jun. 2004, pp. 1047–1052.