

A Novel ASIC Design Flow using Weight-Tunable Binary Neurons as Standard Cells

Ankit Wagle, Gian Singh, *Member, IEEE*, Sunil Khatri, *Senior Member, IEEE*, Sarma Vrudhula, *Fellow, IEEE*

Abstract—In this paper, we describe a design of a mixed-signal circuit for a binary neuron (a.k.a perceptron, threshold logic gate) and a methodology for automatically embedding such cells in ASICs. The binary neuron, referred to as an FTL (flash threshold logic) uses floating gate or flash transistors whose threshold voltages serve as a proxy for the weights of the neuron. Algorithms for mapping the weights to the flash transistor threshold voltages are presented. The threshold voltages are determined to maximize both the robustness of the cell and its speed. The performance, power, and area of a single FTL cell are shown to be significantly smaller (79.4%), consume less power (61.6%), and operate faster (40.3%) compared to conventional CMOS logic equivalents. Also included are the architecture and the algorithms to program the flash devices of an FTL. The FTL cells are implemented as standard cells, and are designed to allow commercial synthesis and P&R tools to automatically use them in synthesis of ASICs. Substantial reductions in area and power without sacrificing performance are demonstrated on several ASIC benchmarks by the automatic embedding of FTL cells. The paper also demonstrates how FTL cells can be used for fixing timing errors after fabrication.

Index Terms—Artificial Neuron, Perceptron, Neural Circuits, Threshold Logic, Floating Gate, Flash, Low Power, High Performance

I. INTRODUCTION AND MOTIVATION

In this paper, we introduce a new *programmable ASIC primitive*, referred to as a *flash threshold logic* (FTL) cell, and show how it can be used to substantially improve the area and power consumption of an ASIC, without increasing its delay. An FTL cell and its use in an ASIC is different from any other type of ASIC component previously reported. It is a mixed-signal circuit that is designed as a *standard cell*, so that it is fully compatible with conventional CMOS logic synthesis, technology mapping, and place-and-route tools.

An FTL cell of n inputs realizes any *threshold* function of n or fewer variables. A threshold function $f(x_1, \dots, x_n)$ [1] is a unate Boolean function whose on-set and off-set are *linearly separable*, i.e. there exists a vector of weights $\mathbf{W} = (w_1, w_2, \dots, w_n)^T$ and a threshold T such that

$$f(x_1, x_2, \dots, x_n) = 1 \Leftrightarrow \sum_{i=1}^n w_i x_i \geq T, \quad (1)$$

Authors A. Wagle, G. Singh, and S. Vrudhula are with the School of Computing and Augmented Intelligence, Arizona State University, Tempe, AZ, USA (e-mail: (awagle1,gsingh58,vrudhula)@asu.edu)

Author S. Khatri is with the Dept. of Electrical and Computer Engineering, Texas A&M University, College Station TX, USA (e-mail: sunilkhatri@tamu.edu)

The research was supported in part by NSF PFI award #1701241, #1361926.

Copyright © 2022 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

¹W.L.O.G, weights can be assumed to be positive integers [2], and for a given truth table of a threshold function, there is a weight vector whose sum is minimum [2].

where \sum here denotes the arithmetic sum. A threshold function can be equivalently represented by $(\mathbf{W}, T) = (w_1, w_2, \dots, w_n; T)$.

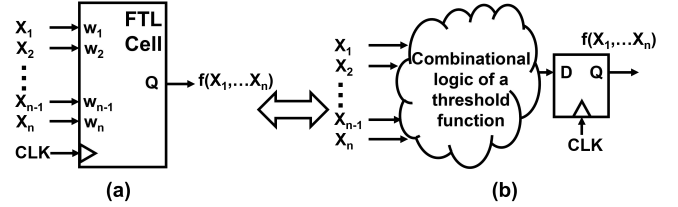


Fig. 1: (a) FTL Schematic, (b) Functional equivalent

Figure 1(a) shows a block diagram of an FTL cell, in which the weights \mathbf{W} are shown as internal parameters of the cell and Figure 1(b) shows its functional equivalent. The schematic is meant to convey that the input-output behavior of an FTL cell may be viewed as an *edge-triggered, multi-input* flip-flop, whose output is the value of a threshold function, that is internally latched at the rising edge of the clock signal CLK.

A distinctive characteristic of the FTL cell design is that the actual threshold function realized by an FTL instance within an ASIC is *programmed after the circuit is manufactured*. An FTL-based ASIC integrates *flash* or *floating gate* [3] transistors along with conventional MOSFETs within the FTL cell. Thus, unlike many of the *emerging technologies* [4], [5], [6], [7], an FTL cell employs mature IC technologies (CMOS and Flash) that are routinely integrated today and commercially manufactured with high yields.

A. FTL in ASIC Design – Overview

Before proceeding to the details of FTL design, it will be instructive to understand how it can be used in an ASIC [8], and how it can improve performance, power, and area.

Consider the logic netlist shown in Figure 2(a) which has two registered outputs F and G . Suppose that the transitive fan-in (TFI) cones of F and G are traversed and two subcircuits labeled A and B (see Figure 2(b)) are found, such that A and B are threshold functions of their inputs. The remaining subcircuit is labeled as C . Suppose that subcircuits A and B are replaced by FTL cells, which are to be later programmed to realize A and B . This replacement is shown in Figure 2(c), where the FTL cells are shown as black boxes. Now, subcircuit C would be re-synthesized to account for the changes in the delay of FTL cells and the new loads that the inputs of the FTL cells present to the outputs of C . There are two reasons why the circuit in Figure 2(c) would have improved power, performance and area.

- 1) Subcircuits A and B and the two flip-flops are each replaced by an FTL cell which has much fewer transistors, resulting in a significant reduction in area and power.

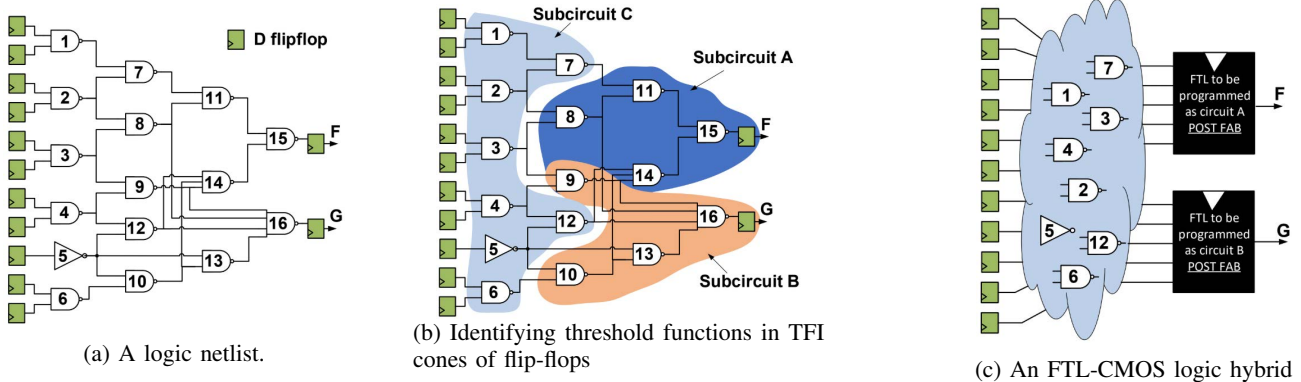


Fig. 2: Use of FTL in ASIC design.

- 2) The clock-to-Q delay of FTL cells turns out to be much less than the total delay (combinational logic delay plus clock-to-Q delay of DFF) of subcircuits *A* and *B*, which results in creating a substantial amount of slack (required time minus arrival time) on the outputs of subcircuit *C*. This in turn will allow synthesis and technology mapping tools to reduce the logic area of subcircuit *C*. The extent of the improvement depends on how the logic is *absorbed* into the FTL cell.

Note that the reason why the FTL cells are shown as black boxes in Figure 2(c) is to convey the fact that their functions are not known at the time of fabrication because the flash transistors are *programmed* after the chip is manufactured.

B. Main Contributions

- 1) An FTL cell is a mixed-signal circuit that is implemented as a standard cell. The new design incorporates flash transistors, which allow its function to be *programmed* after fabrication.
- 2) The set of threshold voltages of the flash transistors in an FTL cell serve as a proxy for the weights $[W, T]$. The weights can be realized with great fidelity because the flash transistors can be programmed with high precision [3]. However, the relationship between the weights and threshold voltages is a non-linear and multi-valued mapping that depends on the complex electrical and layout characteristics of the MOSFETs and flash transistors. We present a new algorithm called the *modified perceptron learning algorithm* (mPLA) [9] that works in concert with HSPICE and *learns* a mapping between the weights and threshold voltages. The mPLA algorithm also maximizes the noise tolerance and robustness of the FTL cell in the presence of process and environmental variations.
- 3) We present an efficient architecture and methodology for programming the threshold voltages of each flash transistor within each FTL cell that is embedded in an ASIC. We also demonstrate how the post-fabrication threshold voltage assignment capability can be used to improve functional yield and correct timing errors.
- 4) FTL cells are designed as standard cells to be compatible with existing CMOS design methodologies and tools.

We demonstrate this compatibility by using commercial CAD tools to perform synthesis, technology mapping, and place-and-route of several complex function blocks with FTL cells included in the cell library. The results show that automatic embedding of FTL cells results in substantial improvements in the area, power, and wire-length, without sacrificing performance.

The remainder of the paper is organized as follows. Section II gives a very brief overview of threshold logic and flash transistor technology. Sections III, IV, V and VI contain the main body of this work. The architecture and operation of the FTL cell are described in Section III. Section IV explains the mechanism for programming FTL cells once they are embedded in an ASIC, using a separate scan chain reserved for that purpose. Section V describes the mapping between the weights of a threshold function and the threshold voltages of the flash transistors in the FTL, considering factors such as robustness to noise, process variations and circuit delay. Section VI contains an extensive set of experimental results, demonstrating the significant improvements in PPA of FTL cells over their CMOS equivalents both at cell-level and when they are integrated into ASICs. It also includes results validating several uses of post-fabrication programming/tuning of the flash devices. Conclusions appear in Section VII.

II. BACKGROUND

A. Threshold Logic

Threshold functions are an interesting and valuable subset of Boolean functions. They were first proposed in 1943 as simple models of neurons [10], which generated a vast number of papers on neural networks – a subject that has been revived recently with the emergence of machine learning. The use of threshold logic in digital design and synthesis was extensively investigated in the 1960s and 1970s, culminating in two authoritative works [11], [1]. Since then there has been a substantial body of work on new circuit architectures and implementations of threshold logic. Surveys of designs prior to 2003 appear in [12], [13], [14], detailing over fifty different implementations.

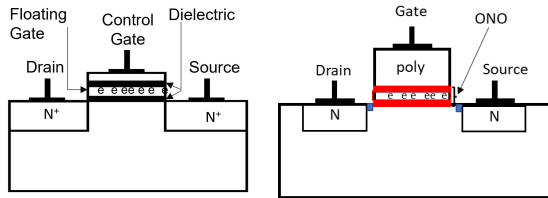
The earlier works and even later ones such as [15], [16], [17], [18], [19], [20], [5], have not been integrated into mainstream VLSI design. It is, however, still very valuable

to develop efficient implementations of threshold functions. This is due to the fact that many Boolean functions that require large AND/OR networks can be realized by smaller, fixed depth threshold networks [2] and nearly 70% of the functions in two standard cell libraries (observed in a 65nm and a 40nm library from different vendors) are threshold functions.

Recently, [8] reported an architecture of a threshold gate and showed how it can be integrated with the standard-cell ASIC design methodology using commercial tools. Unlike our approach, [8] uses only CMOS devices. In addition, they reported significant improvements in PPA of an actual silicon implementation of ASIC with threshold gates [21]. Their architecture, however, severely limits the number of threshold functions that can be implemented because the width of the transistors are made proportional to the weights. This limits the fan-in and consequently, only 12 of the 117 threshold functions of 5 or fewer inputs were implemented in [8]. In contrast, our work implements all 117 threshold functions of 5 or fewer inputs because of the use of flash transistors.

B. Flash Transistors

A flash transistor is functionally similar to a field effect transistor (FET), except that it is made to have an additional layer of charge between the control gate and the channel as a means to adjust the threshold voltage (V_T) of the transistor. Programming and erasing these devices correspond to increasing and decreasing V_T , and this is achieved by electrons tunneling into or out of the charge layer via Fowler-Nordheim (FN) tunneling [22]. In general, programming is performed by applying a sequence of high voltage pulses (the duration and magnitude of which determines V_T) [23] to the control gate and holding the bulk, source, and drain terminals at 0V. Erasure is achieved by holding the control gate at 0V and allowing the source and drain to float [23]. In a flash memory, the value of V_T , which is determined by sensing the current, represents the state or stored value, and this can be retained for more than ten years [23], while the bulk is driven to a high voltage. Several variants of flash transistors with different structures and materials have been investigated over the past two decades to reduce the programming voltage, improve reliability, encode multiple bits, reduce the number of fabrication steps, and improve the yield.



(a) Floating gate Transistor (b) Charge Trap Transistor
Fig. 3: Cross section of flash transistors.

Figure 3 shows a cross-section of the two common types of flash transistors: (a) the floating gate transistor (FGT), which has an additional buried and un-contacted floating gate, and (b) the charge trapping transistor (CTT) which has an oxidenitride-oxide (ONO) layer between the gate and the substrate. Floating gate technology (Figure 3(a)) is fully compatible with

and used along with CMOS, and because of its dominant role as NVM in flash drives and solid-state drives (SSD), its design and fabrication has been continuously improved over two decades. However, it still has several drawbacks, including the need for additional masks, the requirement of higher voltages for programming and erasure, and most importantly, the difficulty in scaling its dimensions below 40nm due to poor scaling of the thin oxide.

In CTTs (Figure 3(b)), the electrons are trapped in the insulating nitride layer whereas in FGTs they are in the conducting floating gate. There are several variants of CTT device such as SONOS [24], MONOS [25] and High-K Metal Gate (HKMG) [26]. All have been successfully scaled to 14nm/16nm CMOS FinFET technology. Additionally, the HKMG device can be programmed to multiple V_T levels [27] and the SONOS device can be programmed up to 128 V_T levels [28]. One important advantage of the HKMG CTTs is that they do not require any additional processes or masks and operate at logic-compatible voltages.

In summary, flash transistors, including FGTs and all variants of CTTs can co-exist with CMOS transistors on the same substrate, in a high-yield, cost-effective manufacturing flow. In this paper, their use will be for realizing logic as opposed to just storage. The use of flash transistors in logic design was described in [29], [30], [31], [32], where the authors demonstrated substantial improvements in power, performance, and area over conventional CMOS standard-cell based design and resilience to aging by reprogramming the V_T s of the flash transistors when aging-related speed degradation occurred. The main drawbacks of the approach are (1) the circuit structures are not compatible with the standard-cell based design flow that is practiced in industry and (2) they are potentially subject to the similar read/write disturb issues found in memory applications.

In this paper, we describe how flash transistors can be used to realize threshold logic gates. The flash devices are used as resistors (by varying their threshold voltages) and the resistances serve as a proxy for the weights in a threshold function. This concept first appeared in [33]. Their design was a stand-alone, custom-designed analog circuit to realize a 16-input threshold function. In contrast, threshold gates described in this paper are designed as standard cells, and automatically incorporated in large-scale ASIC design using commercial tools.

III. FLASH THRESHOLD LOGIC (FTL) CELL

Figure 4 shows the architecture of the FTL cell. It has five main components: (i) the left input network (LIN), (ii) the right input network (RIN), (iii) a sense amplifier (SA), (iv) an output latch (LA), and (v) a flash transistor programming logic (PL). The LIN and RIN consist of two sets of inputs (ℓ_1, \dots, ℓ_n) and (r_1, \dots, r_n), respectively, with each input in series with a flash transistor. In our implementation, $\ell_i = r_i = x_i$ for all i . The conductance of the LIN and RIN is determined by the state of the inputs and the threshold voltages of the flash transistors. The assignment of signals to the LIN and RIN is done to ensure sufficient difference in their conductance across all minterms.

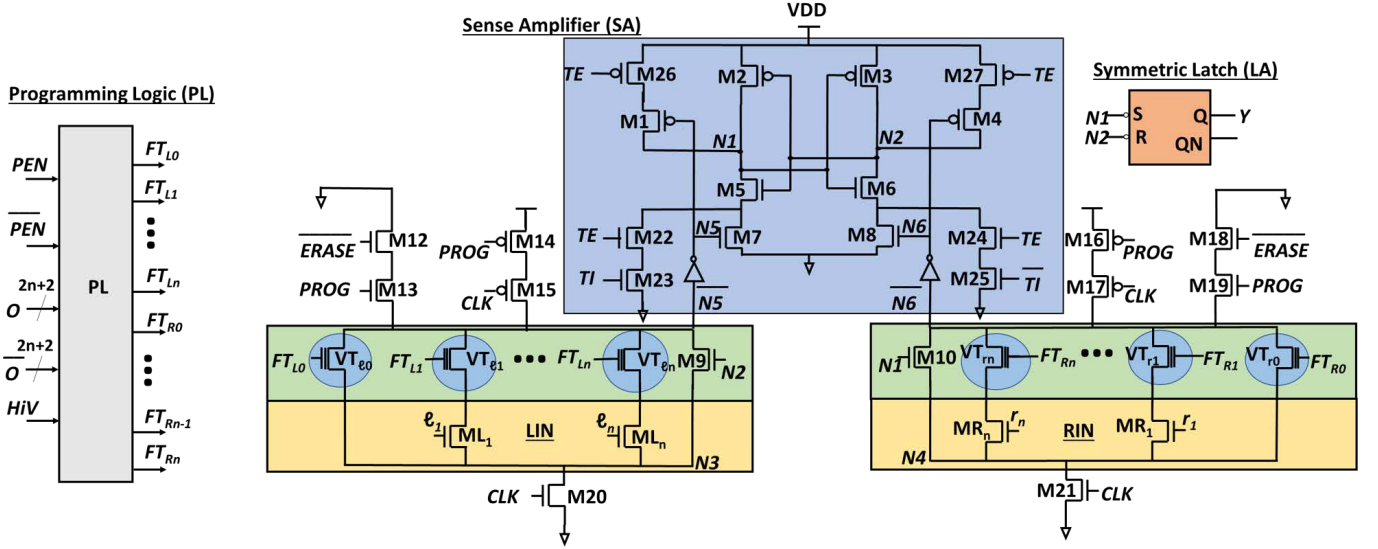


Fig. 4: FTL Cell Architecture showing LIN, RIN, Sense Amplifier (SA), Latch (LA), and Programming Logic (PL).

There are two differential signals $N1$ and $N2$ in an FTL cell, which serve as inputs to an SR latch. When $[N1, N2] = [0, 1]$ ($[1, 0]$), the latch is set (reset) and the output $Y = 1(0)$. The magnitudes of the two sides of the inequality in the definition of a threshold function (see Equation 1) are mapped to the conductance G_L of the LIN and G_R of the RIN. Ideally, the mapping is such that $[N1, N2] = [0, 1] \Leftrightarrow G_L > G_R$ and $[N1, N2] = [1, 0] \Leftrightarrow G_L < G_R$. As stated earlier, the flash transistor threshold voltages serve as a proxy to the weights of the threshold function – the higher the weight, the lower will be the threshold voltage. For a given threshold function, this non-linear monotonic relationship is *learned* using a modified perceptron learning algorithm described in Section V.

The FTL cell has four modes: *regular*, *erase*, *programming* and *scan-testing* mode. The VT values of the flash transistors are set in the programming mode and erased in the erase mode. The logic operation of an FTL cell takes place in regular mode. **FTL Regular Mode:** ($PROG = ERASE = 0$, $TE = 0$, $HiV = 0$). Assume that the VT s of the flash transistors have been set to appropriate values corresponding to the weights of the threshold function, and their gates are being driven to 1 by setting FT_i to VDD. When $CLK = 0$, the circuit is reset. In this phase, the nodes $N5$ and $N6$ of LIN and RIN are connected to the supply, $N5 = N6 = 0$, and $N1 = N2 = 1$. Therefore, the output Y remains unchanged.

Assume now that an on-set minterm is applied to the inputs in the LIN and RIN. With properly assigned VT values to the flash transistors, suppose that $G_L > G_R$ for the given minterm. When $CLK : 0 \rightarrow 1$, both the LIN and RIN will conduct, and $N5$ and $N6$ will both transition from $0 \rightarrow 1$. Assuming $G_L > G_R$, $N5$ rises faster than $N6$, and hence $N5$ will make $M7$ active before $N6$ makes $M8$ active. This will start to discharge $N1$ before $N2$. When $N1$ falls below the threshold voltage of $M6$, it will stop further discharge of $N2$, and turn on $M3$, resulting in $N2 : 0 \rightarrow 1$. Finally, $[N1, N2] = [0, 1]$ sets the SR latch, resulting in $Y = 1$. For an off-set minterm, $G_L < G_R$, and $[N1, N2] = [1, 0]$ resulting in $Y = 0$.

FTL Program, Erase and Scan-testing mode: Figure 4 shows a circuit block labeled PL (programming logic) that generates the signals to select and program the FTL cells at the chip-level. Details of the programming architecture and protocol are given in Section IV. During flash-programming of a single FTL, the PL redirects HiV to FT_i , to program the i^{th} flash transistor.

FTL Programming Mode: ($ERASE=0$, $PROG=1$, $CLK=0$, $HiV=20V$, $TE=0$). The ERASE and PROG signals turn on M12 and M13 and turn off M14. In this state, the source of the flash transistor is floating while the drain and bulk are connected to the ground. Activating the appropriate signals in the PL unit causes high voltage pulses to be applied to the HiV line and the gate of the flash transistor to set the desired threshold voltage (VT).

FTL Erase Mode: ($ERASE=1$, $PROG=1$, $CLK=0$, $HiV=-20V$, $TE=0$). M12 is turned off by the ERASE signal. Both the source and drain of the flash transistors are floating in this state, while the bulk is connected to the ground. Using the PL unit, the gate terminals of all the flash transistors in the FTL are connected to HiV . A negative high voltage pulse at HiV in this state will tunnel the charge from the floating gate, thereby erasing the flash transistor.

FTL Scan-testing Mode: ($ERASE=0$, $PROG=0$, $CLK=0$, $HiV=0$, $TE=1$). The scan-testing mechanism in the FTL cells is implemented in the same way as described in [8]. It uses the *test enable* (TE) and *test input* (TI and \bar{TI}) signals. In this mode, TE acts as the clock with the main clock $CLK = 0$. Hence the scan-testing chains for the D flip-flops and FTL cells are kept separate. The procedure to inject data into the scan-testing chain of FTL cells is straightforward. On each TE cycle, the bit to be scanned in is applied to TI. Then $TE : 0 \rightarrow 1$, which causes the either $N1$ or $N2$ to discharge resulting in the output latch being set or reset. This process is repeated to load all FTLs with the data in a test vector. Transistors M26 and M27 block potential DC paths from VDD to VSS during testing.

Note that the problem of read and write disturb [34], [35]

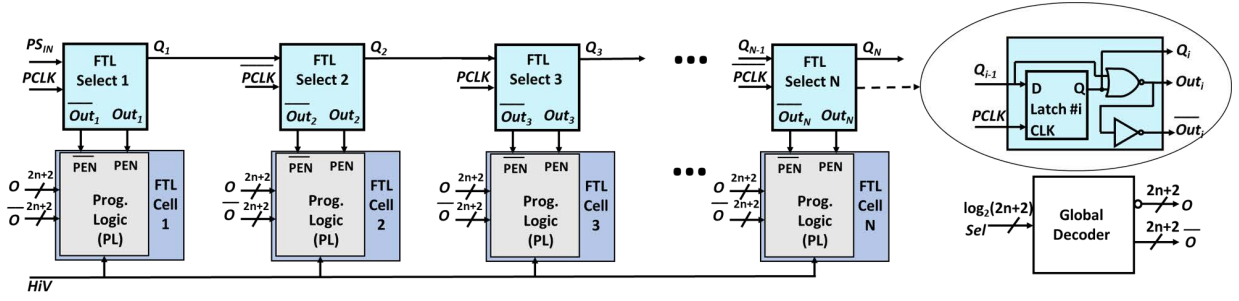


Fig. 5: Programming Scan Chain for FTL cells in an ASIC.

found in NAND flash memories does not exist with an FTL cell because there is only one flash transistor in each stack in the input network. Also the problem of *write endurance* [36] in flash memories, which refers to a limit on the number of writes (from 10K - 100K cycles), is not an issue with FTL cells, because an FTL cell would be programmed or erased only a handful of times over the life of the design.

IV. ARCHITECTURE FOR PROGRAMMING FTL CELLS

In this section, we describe the programming architecture used for FTL cells embedded in an ASIC. This architecture individually addresses the flash transistors of the FTL cells and then redirects HiV pulses to them. Although this architecture is a part of the ASIC, its use ends once the FTL cells are programmed. Therefore, its design must meet its functional requirements without severely impacting the overall area and wirelength of the ASIC. This is achieved by a scan-chain programming architecture.

Figure 5 shows the structure of the programming scan chain. Each stage of this chain consists of an FTL cell with its programming logic and a select cell that identifies the FTL cell to be programmed. Suppose that the FTL cells realize all threshold functions of n or fewer variables.² Then each such cell has $2n + 2$ flash transistors. Suppose further that there are N FTL cells. Initially, all Q_i s are set to 1. Then cell i is selected by making $Q_i = Q_{i-1} = 0$, while all other Q s remain at 1. Thus, clocking in the appropriate sequence using PCLK selects cell i . Since each FTL cell has $2n + 2$ flash cells, a global decoder with $\log(2n + 2)$ inputs and $2n + 2$ outputs is used. These outputs of the decoder activate the appropriate path for the HiV pulses to the inputs of the flash transistors of the selected FTL cell.

The programming architecture involves the use of high voltage nets. In the physical layout, the high voltage wires are bundled, and wire-shielding [37] is used to avoid any cross-talk due to high voltage signals to the other low voltage lines and transistors. Programming is done with a dedicated scan chain, and all the associated high voltage nets are systematically bundled and shielded. This results in reducing the total wirelength of those nets. Furthermore, since it is a linear iterative array, it easily scales to accommodate any number of cells.

Assuming that FTL cells use floating gate transistors, the program and erase modes require +20V and -20V (HiV) pulses

to be applied to their inputs (see Section II-B). Note that other flash technologies such as SONOS [24], MONOS[25] and High-K Metal Gate (HKMG) [26] require similar infrastructure for programming and erasure, but may differ in the voltage levels of the pulses.

V. COMPUTING THE RELATIONSHIP BETWEEN THE WEIGHTS AND THE V_T VALUES FOR AN FTL CELL

A. Overview

Let $\mathbf{VT}_\ell(f) = (VT_{\ell_0}(f), \dots, VT_{\ell_n}(f))$, and $\mathbf{VT}_r(f) = (VT_{r_0}(f), \dots, VT_{r_n}(f))$ denote the threshold voltages of the flash transistors in the LIN and RIN of an FTL, respectively (see Figure 4). Further, let $\mathbf{VT}(f) = (\mathbf{VT}_\ell(f), \mathbf{VT}_r(f))$. In this section, we present an algorithm to determine these voltages for an FTL to realize a given threshold function f having a weight vector $[\mathbf{W}, T]$.

To configure an FTL, the method described herein determines $\mathbf{VT}(f)$ for each f a priori, using models that include circuit parasitics and global and local process variations in the device and circuit parameters. This is to ensure that an overwhelming majority ($\gg 99\%$) of the FTL instances on a chip can be programmed by attempting at most a few pre-computed values of $\mathbf{VT}(f)$. The remaining small fraction of FTLs for which a feasible, model-based $\mathbf{VT}(f)$ could not be found, can be programmed directly on the chip.

Let $G_L(x|\mathbf{VT}(f))$ and $G_R(x|\mathbf{VT}(f))$ for $x \in B^n$, denote the conductance of the LIN and RIN as functions of a minterm x of f and the flash transistor threshold voltages. Henceforth, for clarity, we refer to $G_L(x|\mathbf{VT}(f))$ and $G_R(x|\mathbf{VT}(f))$ as G_L and G_R respectively.

The problem is to find a $\mathbf{VT}(f)$ that determines a mapping between the Boolean space B^n and the conductance space (G_L, G_R) such that, in the *ideal* case,

$$\begin{aligned} G_R &< G_L, \text{ if } f(x) = 1 \\ G_R &> G_L, \text{ if } f(x) = 0. \end{aligned} \quad (2)$$

This mapping, depicted in Figure 6, is one-to-many, since there can be many (an uncountable number) *feasible* values of $\mathbf{VT}(f)$ for a given f with a weight vector $[\mathbf{W}, T]$.

In practice, to avoid variations due to parasitics which could make the circuit behavior erroneous, we require finding a subset of the feasible set where

$$\begin{aligned} G_R &< G_L - \Delta_L, \text{ if } f(x) = 1 \\ G_R &> G_L + \Delta_R, \text{ if } f(x) = 0 \end{aligned} \quad (3)$$

²In the experimental results, $n = 5$.

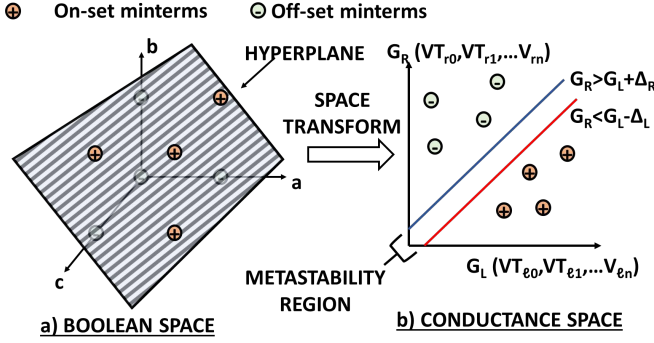


Fig. 6: Transformation from Boolean space to conductance space.

for some sufficiently large $\Delta_L \in (0, G_L)$ and $\Delta_R \in (0, G_R)$. Note that Δ_L and Δ_R are related due to the constraints imposed by the truth table.

Our approach to find $\mathbf{VT}(f)$ consists of three steps which are implemented by Algorithms mPLA^0 , mPLA^+ , and mPLA^{++} . These are described in the following sections.

B. Algorithm mPLA^0

Algorithm mPLA^0 is a modified version of the classical perceptron learning algorithm (PLA) [9] that works in concert with HSPICE (for verifying the truth table of f) to search through the space of values of $\mathbf{VT}(f)$ until each minterm (a point in the (G_L, G_R) space) of f is correctly classified. It does this by iteratively adjusting the threshold voltages of flash transistors such that points in the conductance space that correspond to the on-set and off-set minterms satisfy the constraints in Equation 3 (see Figure 6). It calls HSPICE (line 3 of Algorithm mPLA^0) to determine whether any point falls above or below the lines corresponding to these inequalities. Since a layout extracted FTL circuit is being used, the circuit parasitics are accounted for in the HSPICE simulation. Consequently, Algorithm mPLA^0 implicitly finds values of Δ_L and Δ_R .

Algorithm mPLA^0 Modified Perceptron Learning Algorithm

Input: Truth table TT of threshold function f

Output: \mathbf{VT}_0 to program FTL cells with f

```

1: Initialize  $\mathbf{VT}_0$ 
2: for  $k = 0$  to  $k_{\max} - 1$  do
3:    $\text{OT} = \text{SPICE}(\mathbf{VT}_0)$  // Truth table from simulation
4:   if  $\text{TT}$  and  $\text{OT}$  disagree on some minterm  $m$  then
5:     if  $\text{TT}(m) = 1$  then
6:       Update  $\mathbf{VT}_0$ : decrement (increment) the  $V_T$  of every
       active transistor in LIN (RIN) that is '1' in  $m$  by  $\delta$ 
7:     else
8:       Update  $\mathbf{VT}_0$ : increment (decrement) the  $V_T$  of every
       active transistor in LIN (RIN) that is '1' in  $m$  by  $\delta$ 
9:     end if
10:  else
11:    Break
12:  end if
13: end for
```

Given the truth table (TT) of f , mPLA^0 applies all the minterms of f to the FTL cell and records the HSPICE

response in OT (output table). If $\text{TT}(m) \neq \text{OT}(m)$, for some minterm m , then the constraint in Equation 3 was not satisfied. In such a case, mPLA^0 adjusts the values of $\mathbf{VT}^0(f)$ (Algorithm mPLA^0 line 4-9) associated with the active input transistors within the interval $[\delta, V_{DD} - \delta]$, by a minimum increment δ , according to Equation 4. Here, m_ℓ (m_r) is a binary vector that identifies the active input transistors in the LIN (RIN).

$$\mathbf{VT}_\ell^{k+1} = \begin{cases} \mathbf{VT}_\ell^k - \delta m_\ell & \text{if } m \cdot W \geq T \\ \mathbf{VT}_\ell^k + \delta m_\ell & \text{if } m \cdot W < T \end{cases} \quad (4)$$

$$\mathbf{VT}_r^{k+1} = \begin{cases} \mathbf{VT}_r^k + \delta m_r & \text{if } m \cdot W \geq T \\ \mathbf{VT}_r^k - \delta m_r & \text{if } m \cdot W < T \end{cases}$$

The term δm_ℓ (or δm_r) is a vector which has a value δ at all locations in LIN (RIN) which are 1 for a minterm m , and zero elsewhere. For instance, consider the threshold function $b + c \geq a + T$. Let $m = 110$ be an on-set minterm that was incorrectly evaluated. If $\text{TT}(m) \neq \text{OT}(m)$ then $G_R > G_L - \Delta_L$. Therefore G_L needs to be increased (threshold voltages corresponding to the flash transistors of b and c will be decreased) and G_R needs to be decreased (threshold voltages corresponding to the flash transistors of a and T will be increased) for minterm m . Consequently, the threshold voltages of all the flash transistors associated with the active input transistors should be decreased (increased) by δ in the LIN (RIN). A similar change is made if m is an off-set minterm. This is what is expressed in Equation (4). $\mathbf{VT}^0(f)$ is the value returned by Algorithm mPLA^0 .

If a given set of points in B^n is linearly separable (i.e. a threshold function), then the PLA algorithm will terminate in a finite number iterations [2], [10]. Similarly, given a threshold function f , a sufficiently small δ and an FTL instance for which there exists a feasible $\mathbf{VT}(f)$, Algorithm mPLA^0 will terminate in a finite number steps (see [2] for proof of termination). For an n -input threshold function, the upper bound on the number of iterations of the PLA given in [2] becomes $k_{\max} = 2(n+1)\|\mathbf{VT}^0(f)\|^2/\delta^2$. For instance, with $n = 5$ and $\delta = .02V$, $k_{\max} = 3000\|\mathbf{VT}^0(f)\|^2$.

C. Algorithm mPLA^+ : Improving Noise Tolerance

Algorithm mPLA^0 does not consider the relative location of the points with respect to the metastability region defined by the lines $G_R = G_L - \Delta_L$ and $G_R = G_L + \Delta_R$ (see Figure 6b). Even though minterms are classified correctly, they can be arbitrarily close to the metastability region. The further a minterm is from this region, the easier (and faster and more robust) it will be for the sense amplifier to detect the difference between $N5$ and $N6$, and discharge the appropriate side ($N1$ or $N2$) first. Thus, maximizing Δ_L and Δ_R within the feasible set will maximize its noise tolerance.

Algorithm mPLA^+ repeatedly calls mPLA^0 to maximize Δ_L and Δ_R . It does this by introducing a *hypothetical* capacitance C_1 on node $N5$ (which is introduced in HSPICE) when classifying an on-set minterm, and determining the maximum value of C_1 for which Algorithm mPLA^0 converges. This modification *handicaps* node $N5$ and directs the algorithm

to find a solution, that will result in an increased gap between G_L and G_R . Similarly, we add a capacitance C_0 on node $\overline{N6}$, when classifying an off-set minterm. Since the values of Δ_L and Δ_R are linearly proportional to C_1 and C_0 respectively, the separation between the lines $G_R = G_L - \Delta_L$ and $G_R = G_L + \Delta_R$ increases, which in turn forces the training algorithm to produce a threshold voltage assignment $\mathbf{VT}^+(f)$ in a more robust (and also faster) FTL cell. Note that C_1 and C_0 are only used during HSPICE simulations, and are not part of the actual FTL cell.

Figure 7 shows the results of running Algorithms mPLA^0 and mPLA^+ on a test function ([1]) $f_{115}(a, b, c, d, e) : (W, T) = [4, 1, 1, 1, 1; 5] = a(b + c + d + e)$. It is plot of the minterms in the conductivity space that was obtained by using HSPICE after programming the FTL using $\mathbf{VT}^0(f)$ and $\mathbf{VT}^+(f)$. The largest values of C_1 and C_0 for which a feasible solution was obtained was $0.1fF$. The plot shows that training with the hypothetical capacitance values separates the two closest on-set and off-set minterms in the conductivity space by more than five times. Furthermore, the delay of an FTL programmed with $\mathbf{VT}^+(f)$ will be smaller than the one that is programmed with $\mathbf{VT}^0(f)$.

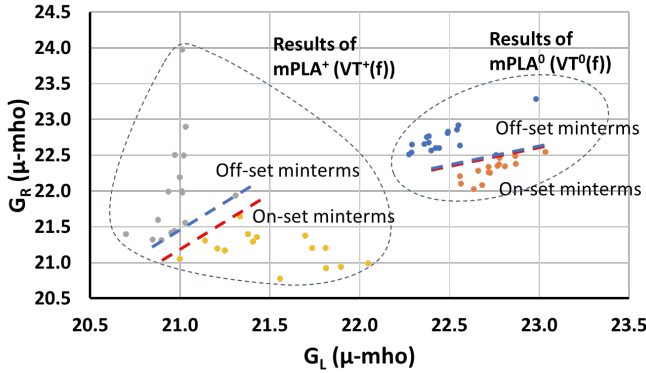


Fig. 7: Conductance G_L and G_R of FTL cell programmed for $f = [4, 1, 1, 1, 1; 5]$ using mPLA^0 and mPLA^+ ($[T, 0.9V, 25^\circ C]$).

D. Algorithm mPLA^{++} : Optimizing Yield

The threshold voltages $\mathbf{VT}^+(f)$ computed by the mPLA^+ are aimed at achieving maximal separation between the on-set and off-set minterms based on model-based estimates of parasitics. This has the twin advantages of increasing the noise margin and reducing the delay. Despite this, inevitable manufacturing variations can still result in reducing the difference between G_L and G_R associated with $\mathbf{VT}^+(f)$ of the two closest minterms, which may result in the incorrect evaluation of the intended threshold function. In this section we present a *predictive* technique to *pre-compute* a small *set* of $\mathbf{VT}s(f)$ for each threshold function f which would cover a very high percentage of manufactured variations.

Among the N manufactured FTL cells programmed to realize function f using $\mathbf{VT}^+(f)$, suppose that N_e were erroneous and let $\{\text{FTL}_1(f), \text{FTL}_2(f), \dots, \text{FTL}_{N_e}(f)\}$ be the erroneous instances. The problem is to find individual threshold voltage assignments for each of these N_e instances so that each

will correctly realize f . Our approach is motivated by two observations.

First, each erroneous function in $\{f_i^e, 1 \leq i \leq N_e\}$ is itself a threshold function. This is simply due to the fact that by construction, an FTL cell only computes threshold functions (see Figure 1). Second, our experiments show that a large number of different instances of an FTL cell, which are reprogrammed with $\mathbf{VT}^+(f)$ and are to realize the same function f , realize the same erroneous function f^e . This suggests that all the erroneous FTL cell instances can be grouped into a few equivalence classes, called *error-types*, with two FTLs belonging to the same error-type if they realize the same erroneous function.

Given a threshold function f , Algorithm mPLA^{++} first generates a *set* of N_{MC} Monte Carlo (MC) instances of an FTL cell and identifies the N_e erroneous instances (i.e. those when programmed with $\mathbf{VT}^+(f)$ do not realize f). The N_e erroneous instances are grouped into M_f error-types. Let $f_i^e, 1 \leq i \leq M_f$, denote the logic functions of the *distinct* error-types observed in a sample of N FTLs. Algorithm mPLA^{++} selects one MC instance from each error-type class and computes one $\mathbf{VT}^+(f)$ assignment for that instance using mPLA^+ . It returns a set of threshold assignments,

$$\mathbf{VT}^{++}(f) = \{\mathbf{VT}^+(f_1^e), \mathbf{VT}^+(f_2^e), \dots, \mathbf{VT}^+(f_{M_f}^e)\}, \quad (5)$$

one for each error-type for each function f .

Algorithm mPLA^{++} Modified Perceptron Learning Algorithm accounting for process variations

Input: TT of f , N_{MC}

Output: $\mathbf{VT}^{++}(f)$ to program FTL cells with f

- 1: Execute mPLA^+ to compute $\mathbf{VT}^+(f)$
- 2: Using MC sampling of the parameter space, generate $N_{MC}(f)$ instances of an FTL cell, and program them with $\mathbf{VT}^+(f)$.
- 3: Among the set of N_{MC} instances, let N_e be the number of instances, which when programmed with $\mathbf{VT}^+(f)$, realize a function other than f , and among these N_e , let M_f be the number of erroneous functions that are distinct.
- 4: Execute the mPLA^+ on one MC instance from each of the M_f erroneous functions to obtain

$$\{\mathbf{VT}^{++}(f)\} = \{\mathbf{VT}^+(f_1^e), \dots, \mathbf{VT}^+(f_{M_f}^e)\}$$

Results presented in Section VI show that using the $\mathbf{VT}^+(f_i^e)$ computed for one FTL instance from i^{th} error-type ($1 \leq i \leq M_f$) resulted in *all* the instances of the same error-type correctly realizing f . This works because instances that have the same error-type share similar parasitic variations. Thus, all instances of our sample of FTL cells were correctly programmed using one distinct $\mathbf{VT}^+(f_i^e)$ for each error-type.

There is no guarantee that the set of erroneous functions found in a sample set N_{MC} will capture all manufacturing outcomes. This means that there may be some manufactured FTLs that could not be correctly programmed using any of the threshold voltage vectors computed by Algorithm mPLA^{++} . For these remaining FTLs, our approach is to apply Algorithm mPLA^0 directly on the chip. In each iteration of mPLA^0 , the step that adjusts the threshold voltages of flash transistors

is replaced by the application of an appropriate number of positive or negative pulses to the FTL cell on the chip using the programming scan chain. This capability of correcting the function of a cell after fabrication to increase yield is a signature attribute of the proposed design methodology.

VI. EXPERIMENTAL RESULTS

A. Experiment Setup

An FTL cell with $n = 5$ (see Figure 4 in Section III) was designed and a complete layout (including the programming devices) was created using the TSMC 40nm LP library. It was laid out as a double height cell requiring 24 tracks. The flash transistor models were obtained from [30] and were suitably modified to reflect the characteristics and variations of the TSMC 40nm LP library. The design rules for the flash transistors were obtained from ITRS. The standard cell area of the FTL was $15.6 \mu m^2$.

There are a total of 117 distinct positive-form threshold functions of five or fewer variables. A numbered list of these is given in [1] and can also be accessed at [38]. The one cell that was designed was copied 117 times, and each was trained to realize one of the 117 threshold functions. In this section, we use the same numbering scheme as in [1] to identify the functions. The FTL cell trained to implement the threshold function numbered n in [1] will be referred to as FTL_n , and the corresponding CMOS implementation will be denoted as $CMOS_n$. The threshold function itself will be denoted as f_n . **Note:** In all the bar charts shown in this section, the numbers on the x-axis identify the threshold function. Function f_0 is a buffer and is omitted because this would correspond to a DFF, which by itself would never be replaced by an FTL in an ASIC. The first function shown is f_1 , which is a two-input AND.

B. Training Iterations

mPLA⁺ was used to train the FTL cell for robustness (see Section V-C) for all 117 functions. Figure 8 shows the number of iterations needed for training of each of the 117 functions. The actual number of iterations was about 10X lower than the theoretical upper bound, presented in Section V-B.

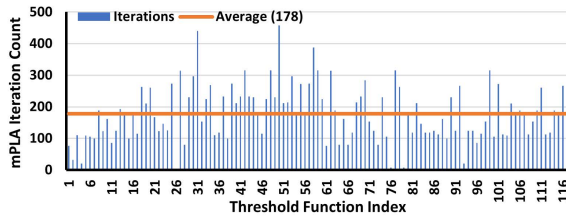


Fig. 8: Iteration count for mPLA⁺ for all 117 functions of 5 or fewer variables.

C. Individual Cell Area, Delay and Power Comparison

All 117 threshold functions of five or fewer variables were implemented using FTL cells. These functions were also synthesized by Cadence Genus [39] and placed and routed using Cadence Innovus [39], using the conventional TSMC 40nm LP standard cells. Two sets of experiments were performed

to compare the CMOS equivalent designs to the FTL cells: (1) delay optimal and (2) area optimal synthesis. The results comparing the total delay (sum of logic delay, setup-time, and clock-to-Q delay), area, and power of these circuits and the corresponding FTL implementations are shown in Figures 9(a) and 9(b), respectively.

The results show that FTL cells have the advantage of speed. Optimizing their CMOS equivalents to meet the delay of the corresponding FTL cells forces the synthesis algorithms to use high drive strength cells (larger area) for the combinational logic and larger DFFs. As the FTL implementations are faster than the fastest CMOS equivalent implementation, delay optimal synthesis results in an across-the-board improvement in all FTL cells in delay, area, and power.

When synthesizing individual cells for minimum area, FTL cells are still uniformly faster. However, the synthesis algorithm now uses the smallest combinational cells and DFFs in the CMOS equivalents. In this case, although the CMOS implementations of simpler functions are much smaller than their FTL equivalents (see Figure 9(b)), the FTL versions are still smaller for 74 out of 117 functions because the *logic absorbed* by the FTL cell results in greater area savings than the smaller drive strength cells used in the CMOS equivalents.

The dynamic power of every FTL implementation is higher than its CMOS equivalent for area optimal synthesis. The reasons for this are (1) an FTL cell resets and then evaluates its function on every clock cycle and (2) the much smaller switched capacitance of the low-drive strengths of the combinational logic in the CMOS equivalents. Figure 9(a) shows that FTL cells have a much lower power-delay product (i.e. energy) when their CMOS equivalents are synthesized for minimum delay. Figure 9(c) shows that this is also true for the majority of the CMOS equivalents when they are synthesized for minimum area. Hence, FTL cells are, in general, more energy efficient.

Figures 9(d) and 9(e) show a comparison of the leakage power of FTL cells and their CMOS equivalents. The leakage of FTL cells is practically independent of the function, and in the case of delay optimal synthesis, it is far lower than every CMOS equivalent circuit. Exactly the opposite is true for the area optimal synthesis due to reduced sizes of the combinational cells and DFFs. In these plots the circuit indices are ordered by increasing area. The area trend lines show that the leakage increases with area for the CMOS implementations.

D. Delay Distributions

This experiment compares the distributions of delays of FTL and CMOS implementations. We show the results for the threshold function f_{35} with a weight vector $[W; T] = [3, 3, 2, 1, 1; 8]$. The PVT corner setting was $[TT, 0.9V, 25^\circ C]$. 100K Monte Carlo instances were generated for both FTL_{35} and $CMOS_{35}$. Each of the 100K FTL instances was verified against the truth table for functional correctness, for both FTL_{35} and $CMOS_{35}$. Figure 10 shows the histogram of delays for both circuits. These demonstrate the delay advantage of the FTL cell over its CMOS equivalent, even in the presence of process variations. The difference in standard deviation

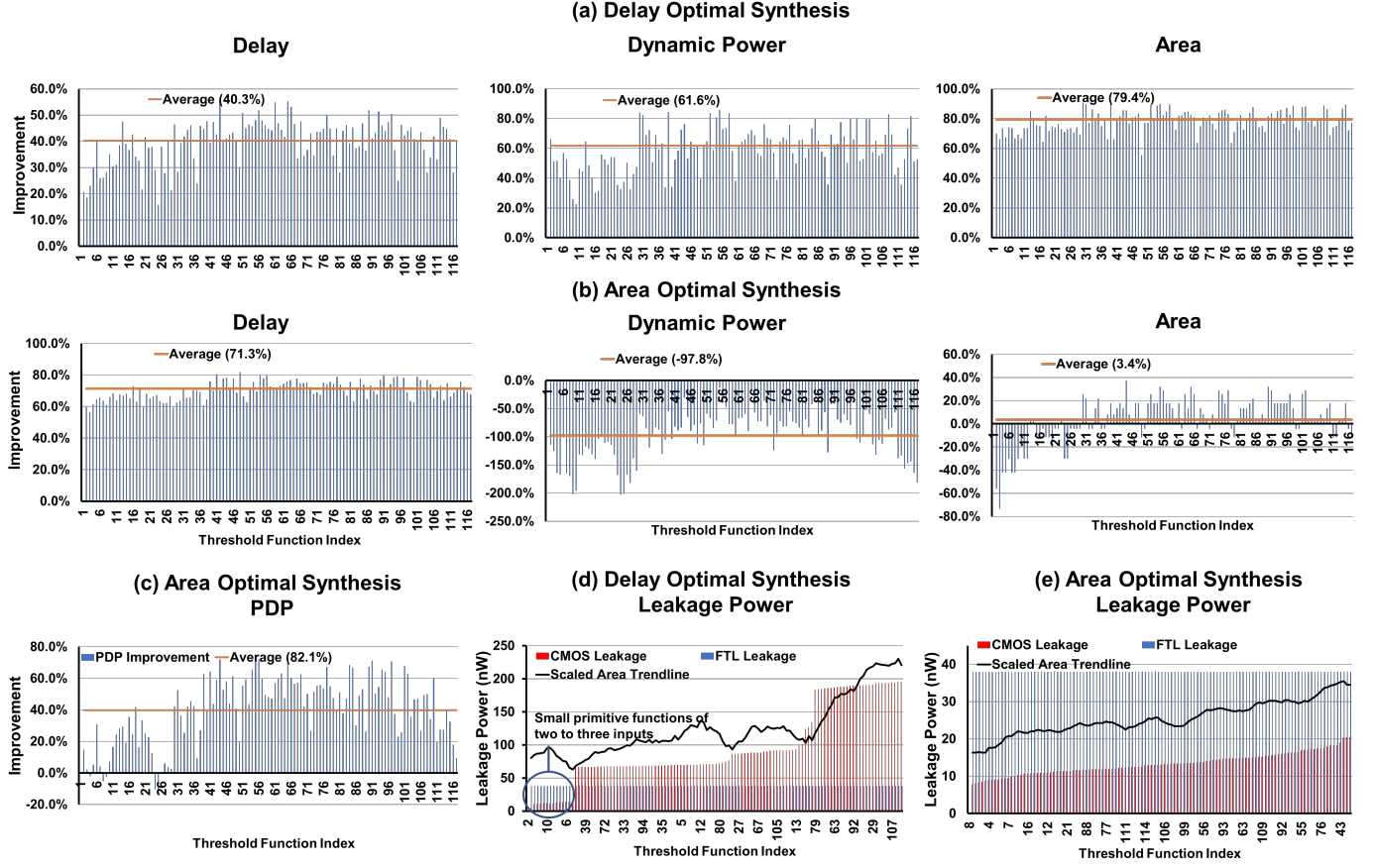


Fig. 9: PPA improvements of FTL over CMOS implementations. Simulations done at 25°C assuming a 20% input switching activity.

between the two is insignificant. Note that the FTL instances with large delays can be *re-programmed* to reduce the delay further. This capability is not possible for the CMOS versions.

by 2.8X), power varies by 5.9X (CMOS equivalent by 1.9X), and the PDP (energy) varies by 2.3X (CMOS equivalent by 1.43X), as the supply voltage varies over [0.8V, 1.1V].

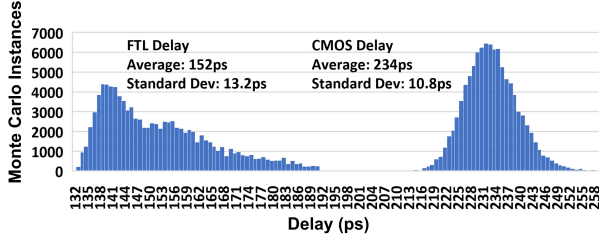


Fig. 10: Delay histogram of FTL_{35} and $CMOS_{35}$ with 100K Monte Carlo simulations. $PVT = [TT, 0.9V, 25^\circ C]$.

E. Dynamic Voltage Scaling

Voltage scaling is a common mechanism to trade off performance against power. Table I shows the results of training FTL_{35} at 0.9V. The FTL cell was programmed with the determined set of flash threshold voltages, and then operated over the voltage range [0.8V, 1.1V]. To ensure proper operation across all voltages, the gate voltages of the flash transistors were scaled accordingly. This result demonstrates how a single $VT^+(f)$ assignment can be used for dynamic voltage scaling. The delay of the FTL_{35} varies by 2.5X (its CMOS equivalent

Supply Voltage (V)	Flash Gate Voltage (V)	Power (uW)	Delay (ps)	PDP
0.8	0.8	14.3	198.1	2837.1
0.85	0.825	20.5	157.6	3228.7
0.9	0.85	26.1	130.2	3396.9
0.95	0.875	40.3	111.2	4482.7
1	0.9	53.1	97.0	5148.6
1.05	0.925	76.0	86.4	6562.9
1.1	0.95	85.0	78.2	6644.0

TABLE I: Delay, total power and power-delay-product (PDP) of FTL_{35} , trained at $V_{DD} = 0.9V$, and $C_0 = C_1 = 0.1fF$.

F. Number of programming pulses

Figure 11 shows the number of high voltage pulses needed to program the 117 threshold functions. The number of high voltage pulses is estimated, assuming that each high voltage pulse would increment the threshold voltage of a flash transistor by 20mV. This assumption will vary across flash transistors. As shown in Figure 11, the number of high voltage pulses needed to program a given FTL cell increases with an increase in the number of variables of the threshold function being implemented.

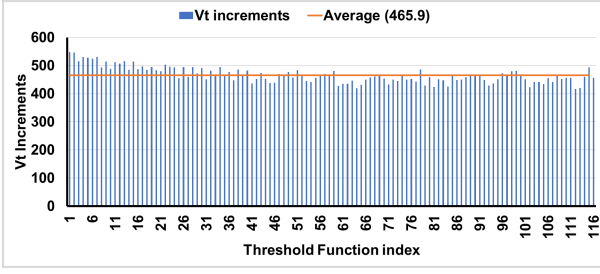


Fig. 11: Number of high voltage (HiV) pulses needed to program the FTL cells with all 117 threshold functions of up to 5 inputs

G. Experiments on Training for Robustness

In this section, we present the results of Algorithms mPLA⁺, and mPLA⁺⁺ for training FTL cells taking into account parasitics and manufacturing variations. The test function $f_{35}(a, b, c, d, e) : (\mathbf{W}, T) = [3, 3, 2, 1, 1; 8] = ab(c + de)$ was chosen for this evaluation as this function generated the most number of error-types ($M_f=61$) out of all the 117 threshold functions when Monte Carlo simulations were run on 20K training samples.

The first experiment consisted of training FTL_{35} using mPLA⁺ for various values of the capacitances C_1 and C_0 , and for each solution, extracting the delay values. The results for this experiment are as shown in Table II. There are two important observations to be made here. First, even though the weights of the inputs d and e are equal, the corresponding flash transistors (V_4 and V_5) may be assigned different threshold voltages. This is because mPLA⁺ compensated for the irregular layout parasitics of both the flash transistors using threshold voltages to realize equal weights. Second, the delay improves with increasing robustness, as discussed earlier in Section V-C. This is because the separation between the lines $G_R = G_L - \Delta_L$ and $G_R = G_L + \Delta_R$ increases with increase in C_1 and C_0 . This increased separation results in a higher voltage difference at the inputs of the sense amplifier, which leads to a faster evaluation of the FTL cell.

C_1, C_0	Average Vt Values (V) ($V_1, V_2, V_3, V_4, V_5; V_{I0}, V_{r0}$)	Delay (ps)
0	0.64, 0.64, 0.66, 0.70, 0.72; 1.00, 0.58	224
0.01	0.60, 0.60, 0.64, 0.68, 0.70; 1.00, 0.50	178
0.02	0.60, 0.60, 0.64, 0.68, 0.70; 1.00, 0.50	178
0.05	0.60, 0.60, 0.64, 0.70, 0.70; 1.00, 0.50	172
0.1	0.56, 0.56, 0.60, 0.66, 0.66; 1.00, 0.42	163
0.15	0.52, 0.54, 0.58, 0.64, 0.64; 1.00, 0.34	154

TABLE II: Delay values of $FTL_{35} = [3, 3, 2, 1, 1; 8]$, trained for robustness using various capacitor values (fF).

The second experiment was aimed at validating mPLA⁺⁺. We used f_{35} as a test function. The first step is to create the database $\{\mathbf{VT}^{++}(f_{35})\}$. Algorithm mPLA⁺⁺ was given f_{35} and $N_{MC} = 20K$ as inputs. The erroneous instances were grouped into $M_{f_{35}} = 61$ error-types. Algorithm mPLA⁺⁺ generated $\{\mathbf{VT}^{++}(f_{35})\} = \{\mathbf{VT}^+(f_{35,1}^e), \dots, \mathbf{VT}^+(f_{35,61}^e)\}$.

Next, 100K new MC instances were generated and programmed first with $\mathbf{VT}^+(f_{35})$. Among the erroneous instances, 99.96% of them were one of 61 error-types that

were previously found. When each FTL cell in group j , ($1 \leq j \leq 61$) was programmed with the threshold voltage set $\mathbf{VT}^+(f_{35,j})$, all the erroneous instances correctly computed f_{35} . The remaining .04% of the 100K were correctly programmed by executing mPLA⁰ directly to the chip, starting with $\mathbf{VT}^+(f_{35})$. This required fewer than five iterations on the average for the instances. Since f_{35} had the most number of failure types, all of the other 117 functions, which exhibit fewer failure types, would be equally easy to program correctly in the presence of variations. Thus, all errors caused by process variations were corrected, with the vast majority requiring a single, precomputed VT set and a small fraction requiring on-chip programming.

Stage	Procedure	Yield (%)
Training (20K instances) $M_f=61$	mPLA ⁺⁺	100%
Testing (100K instances)	mPLA ⁺⁺	99.96%
	mPLA ⁰ (On-chip)	100%

TABLE III: Yield when mPLA⁺⁺ and mPLA⁰ (on-chip) are used for programming instances of $FTL_{35} = [3, 3, 2, 1, 1; 8]$.

H. Robustness Against PVT Variations

Figure 12 shows the delay variations in delay of five sample threshold functions [38] w.r.t process, temperature and temperature variations. As expected, FTL cells are slowest in the SS corner and fastest in the FF corner. Furthermore, as the process moves from the SS corner to the FF corner, the delay improves, as expected. When the voltage increases from 0.81 V to 0.99 V, the delay improves. The FTL cells were also tested for reliability for the consumer temperature range of 0°C, 25°C, and 55°C. This result demonstrates that a $\mathbf{VT}^+(f)$ solution, generated using TT 0.9V 25°C can reliably work with PVT variations.

I. Robustness Against VT Drift

Over the lifetime of an FTL cell, the charge stored in the gate of flash transistors eventually leaks into the channel due to the deterioration of thin oxide layer [40], signal disturbances [41], etc. This leakage effectively changes the VT of the flash transistors. By extension, it also changes the weights programmed on the FTL cell. Table IV shows the effect of decreasing VT on the threshold functions programmed on the FTL cells. All 117 FTL cells operated correctly with a VT drift of up to 5mV. Beyond 5mv, some cells failed. However, after testing, their VT s can be reprogrammed to compensate for this drift. Furthermore, all the FTL cells that were selected by Genus when synthesizing ASIC designs (See Section VI-K) operated correctly with 20mV drift in VT .

Vt Drift (mV)	% FTL cells operated correctly
1	100
2	100
5	100
10	96.55

TABLE IV: Robustness against V_T drift for FTL cells programmed with all 117 threshold functions of up to 5 inputs.

Index	Threshold Function	Sum of Products Form
2	[1,1,0,0,1]	$a+b$
7	[2,1,1,0,2]	$a+bc$
18	[3,2,1,1,0,5]	$ab+acd$
37	[2,2,2,1,1,6]	$abc+abde+acde$
60	[4,3,2,2,1,8]	$abc+abd+acd+abe+bcde$

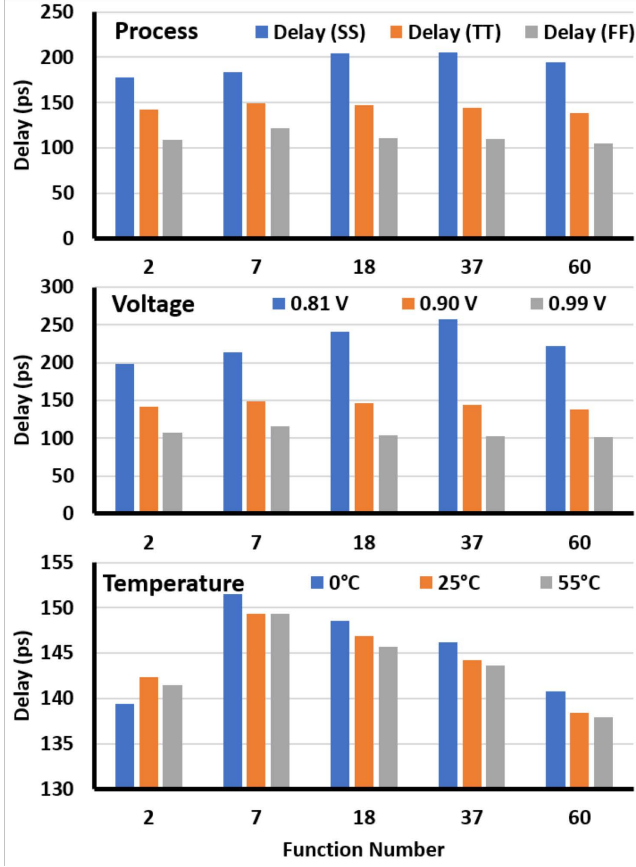


Fig. 12: Delay of an FTL cell for threshold functions, with process (SS, TT, FF), voltage (0.81 V, 0.9 V, 0.99 V), and temperature (0°C, 25°C, 55°C) variations.

J. Post-fabrication Timing Correction

The experiments described in Sections VI-G, VI-D and VI-E demonstrate the flexibility of FTL due to the possibility of configuring its function after fabrication. This characteristic can also be used to correct timing errors.

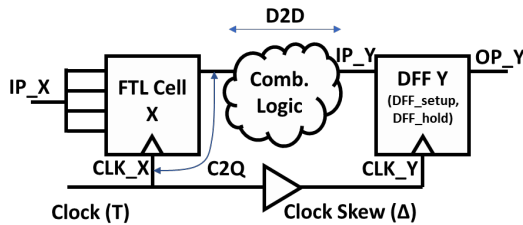


Fig. 13: Datapath to Demonstrate Post-Fabrication Timing Corrections.

Figure 13 shows a small datapath that was constructed to demonstrate how to correct setup-time and hold-time violations after fabrication in an FTL design. The datapath con-

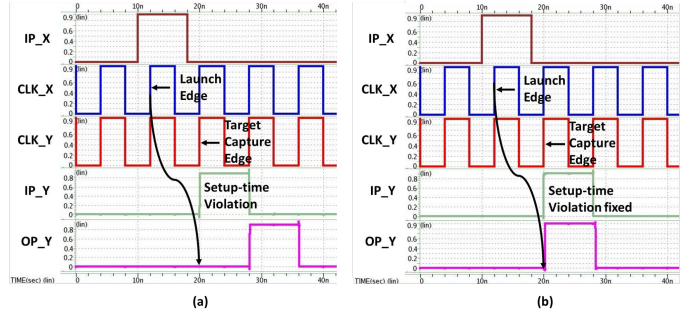


Fig. 14: Post-fabrication setup-time correction using an FTL cell.

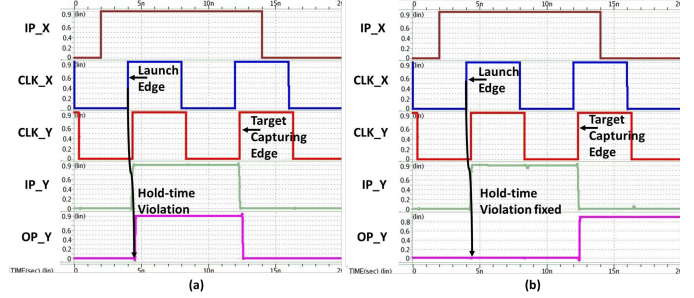


Fig. 15: Post-fabrication hold-time correction using an FTL cell.

sists of clock-to-Q (C2Q) delay, combinational delay (D2D) and DFF specifications for setup (DFF_{setup}) and hold (DFF_{hold}) times. The clock is skewed by an appropriate amount Δ , to generate either a setup-time or a hold-time violation. The violations are corrected by reprogramming the FTL cell to produce different C2Q values.

Figure 14(a) shows how the data launched from FTL X misses the target clock edge at DFF Y, thereby violating setup-time. Figure 14(b) then shows that decreasing the C2Q of FTL X fixes the setup-time violation. Similarly, Figure 15(a) shows how the data launched from FTL X is captured by the target clock edge at DFF Y one cycle early, thereby violating hold-time. Figure 15(b) then shows that increasing the C2Q of FTL X fixes the setup-time violation. Note that we can extend post-fabrication V_T adjustment to also mitigate delay increases due to aging.

K. Delay Optimal Synthesis of ASICs with FTLs

In this section, we show how commercial design tools can accommodate FTL cells in synthesis, and placement and routing. Five circuit blocks were synthesized using the 40nm TSMC standard cell library, which was augmented with FTLs to realize 117 positive forms of all 5-input threshold functions. This was done by creating one cell and making 117 copies and then determining the V_T s of the flash transistors and signal assignments to realize each threshold function. Then each FTL standard cell was characterized in the conventional way. Only the positive forms of the threshold functions were included in the library to keep the increase in the library size to a minimum (about 7%) and to exploit the capability of Genus to recognize NPN equivalents of the cells (see below).

The ASIC benchmarks are: 1) 32-bit Wallace multiplier (Mul), 2) 28-bit FIR filter (FIR), 3) 64-bit floating-point unit multiplier, 4) 16-bit Fast Fourier Transform (FFT), and 5)

Design	Freq. (MHz)	Conventional					FTL-integrated					Improvements			Prog. Time (μsec)
		Std. Cells	DFF	Area (μm^2)	Power (mW)	Wire-length (μm)	Std. Cells	DFF/FTL	Area (μm^2)	Power (mW)	Wire-length (μm)	Area	Power	Wire-length	
Mul	417	19536	343	51855	6.00	118906	11493	272/71	32339	4.64	88592	37.6%	22.7%	25.5%	204.3
Filter	406	53588	529	157482	36.26	436000	41711	281/248	107420	28.41	322400	31.8%	21.6%	26.1%	585.8
FPU	392	48992	1734	132655	27.82	484096	40937	1693/41	98879	24.14	406091	25.5%	13.2%	16.1%	113.2
FFT	667	156242	9614	443356	100.14	1405565	140650	9286/328	368509	86.62	1199160	16.9%	13.5%	14.7%	807.1
SHA	308	33204	2161	109170	15.95	396267	26511	2147/14	66001	13.23	343852	39.5%	17.1%	13.2%	47.9
Avg.				139290	25	425689			96468	21	343504	30.7%	17.7%	19.3%	351.7

TABLE V: Improvement in area, power, and wirelength improvement in ASICs with FTL integrated, over conventional ASICs, without trading off performance. Average improvements are calculated using the geometric mean.

512-bit Secure Hash Algorithm (SHA). Designs were synthesized using Cadence Genus and then placed and routed using Cadence Innovus. Standard cell libraries for FTL cells were characterized using Synopsys HSPICE and generated in Liberty format. Timing checks were performed using cross-corner analysis at {SS, 125C, 0.81V}, {TT, 25C, 0.9V} and {FF, 0C, 0.99V} corners. After placement and routing, the *select* cells and the FTL programming logic cells (see Figure 5 are *paired*. Then engineering change order (ECO) commands stitch the programming scan chain. Since the latter uses high voltage nets, shielding nets are added to protect neighboring nets from high voltage signals. Both versions of each ASIC were verified using Cadence Conformal.

The results of synthesis and P&R, summarized in Table V, demonstrate significant improvements in the area (30.7%), power (17.7%), and wirelength (19.3%) averaged over the designs. These improvements include the *overhead* of the programming infrastructure described in Section IV, which was less than 5% in the worst case. Note that these *across the board* improvements were obtained under *delay-optimal synthesis*. This would not be the case for area-optimal synthesis.

Wherever it was beneficial to improve timing, Genus found and replaced *threshold logic cones* (not necessarily maximal fanout-free cones) driving DFFs with the appropriate FTL cell. This led to a reduction in the number of standard cells. It ranged from 10% to 42%. There are two causes for this reduction. First is the absorption of part of the fanin cone that is a threshold function driving the DFF into the FTL. This eliminates all those cells. A second source is the reduction of the subcircuit (e.g. *C* in Figure 2(b)) that *feeds* the fanin cone. The significant speed advantage of the FTL cell creates large positive slack at the outputs of the *feeder* subcircuit. Consequently, to meet timing, Genus re-synthesizes the feeder with slower logic. Standard logic primitives such as inverters, 2-input gates, 3-input gates, inverters, and even AOI/OAI gates are reduced and the number of complex cells increased, reducing the total cell count.

The last column of Table V shows estimates of the time (i.e., number of pulses) required to program the FTL cells, which increases linearly with the number of FTLs. Although the actual programming time will depend on the technology, it is expected to be on the order of microseconds [23].

Table VI shows the run-time of Genus during synthesis, for all the ASIC designs. While the inclusion of all the 117 FTL cells increases the library size slightly (about 7%), FTL cells allow *faster timing closure* by generating positive

slack. Table VI also shows the peak memory usage of Genus during synthesis, for all the ASIC designs. The peak memory requirements are almost identical even after adding the 117 threshold functions in the library.

	Runtime(sec)		Peak memory (MB)	
	Conv.	FTL-integ.	Conv.	FTL-integ.
Multiplier	1451	636	1269	1292
Filter	2596	2893	1401	1439
FPU	2947	2724	1273	1262
FFT	3102	2653	1421	1416
SHA	1838	1790	1297	1292

TABLE VI: Runtime and Peak memory usage for the synthesis of ASIC designs.

To demonstrate that Genus can recognize NPN equivalences of positive-form threshold functions, we selected a number of threshold functions and negated and permuted their inputs and negated their output. Table VII shows the result of one of the more complex functions. The interpretation of Table VII is as follows. Consider the threshold function $ab + ace + ade + bcd + acd$. The weight-threshold description is $[4, 3, 2, 2, 1; 7] = 4a + 3b + 2c + 2c + d \geq 7$, which is an FTL_{93} . When Genus found a sub-circuit with input negation, $\bar{a}b + \bar{a}\bar{c}e + \bar{a}de + bcd + a\bar{c}d$, it replaced it with an FTL_{93} with \bar{a} and \bar{c} driving inputs *a* and *c*. The last row shows that Genus can detect output negation and maps it to a different cell FTL_{94} whose positive form is $[4, 3, 2, 2, 1; 6] = 4a + 3b + 2c + 2c + d \geq 6$. In each case, the synthesis tool detected the threshold functions and their NPN equivalents, and added inverters as necessary, without using any additional standard logic gates such as AND, OR, etc.

The last experiment conducted was aimed at discovering what threshold functions would be detected if there were no area or delay constraints. Figure 16 shows all possible threshold functions that could be detected in the 32-bit Wallace tree multiplier. The multiplier has 343 DFFs. Excluding the 64 input DFFs, all 279 remaining DFFs and cones of logic feeding them were replaced by FTL cells, showing that complex multi-level logic circuits that are threshold functions frequently occur in logic circuits and synthesis tools can recognize them.

Threshold function	Verilog description of NPN equivalent	Synthesis result
ab+ace+ade+ bcd+acd	$y \leq ((4*a) + (3*b) + (2*c) + (2*d) + (1*e)) \geq 7 ? 1:0;$	FTL_93 (4,3,2,2,1;7)
[4,3,2,2,1;7]	$y \leq ((4*(!a)) + (3*b) + (2*(!c)) + (2*d) + (1*e)) \geq 7 ? 1:0;$	FTL_93 (4,3,2,2,1;7) and two inverters for "a" and "c"
	$y \leq !(4*(!b)) + (3*c) + (2*(!d)) + (2*a) + (1*e) \geq 7 ? 1:0;$	FTL_94 (4,3,2,2,1;6) and three inverters for "a", "c" and "e"

TABLE VII: Detection of NPN equivalents of threshold functions using a library of 117 5-input FTL cells.

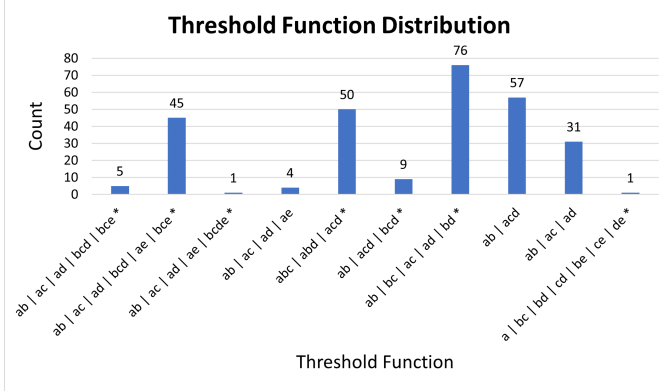


Fig. 16: Distribution of threshold functions in 32-bit multiplier when synthesized using FTL cells with zero-delay zero-power.

VII. CONCLUSION

In this paper we demonstrated that there could be substantial value in going beyond the traditional use of flash technology as memory and using it in CMOS logic. Unlike the many emerging memory technologies, flash technology is mature and compatible with CMOS fabrication. Using flash transistors in conjunction with CMOS transistors, we developed a design of a binary neuron, referred to as FTL, that can realize a large number threshold functions in a single standard cell. We demonstrated several novel features of an FTL cell: (1) it is a *configurable standard cell*, whose function can be configured after fabrication; (2) the configuration is achieved by conventional techniques of tuning the threshold voltages of flash transistors with high fidelity; (3) its design could be optimized to make it very robust in the presence of circuit parasitics and improving robustness also improves its performance; (4) the ability to tune its performance after fabrication provides a novel way to improve the yield in the presence of process variations and correct timing errors; (5) it was designed so that it can automatically be embedded within ASICs using commercial CAD tools, and resulting in significantly improved area and power while still operating at the maximum possible frequency.

REFERENCES

- [1] S. Muroga. *Threshold Logic and its Applications*. John Wiley & Sons, 1971.
- [2] K. Siu, V. Roychowdhury, and T. Kailath. *Discrete Neural Computation: A Theoretical Foundation*. Prentice-Hall, Inc., 1995.
- [3] Y. Cai, E. F. Haratsch, O. Mutlu, and K. Mai. Threshold voltage distribution in MLC NAND flash memory: Characterization, analysis, and modeling. In *2013 Design, Automation Test in Europe Conference Exhibition (DATE)*, pages 1285–1290, March 2013.

- [4] R. Perricone, I. Ahmed, Z. Liang, M. G. Mankalale, X. S. Hu, C. H. Kim, M. Niemier, S. S. Sapatnekar, and J. Wang. Advanced spintronic memory and logic for non-volatile processors. In *Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2017, pages 972–977, March 2017.
- [5] J. Yang, N. Kulkarni, S. Yu, and S. Vrudhula. Integration of threshold logic gates with RRAM devices for energy efficient and robust operation. In *2014 IEEE/ACM International Symposium on Nanoscale Architectures (NANOARCH)*, pages 39–44. IEEE, July 2014.
- [6] P. Gupta and N.K. Jha. An Algorithm for Nanopipelining of RTD-Based Circuits and Architectures. *IEEE Transactions On Nanotechnology*, 4(2):159–167, Mar 2005.
- [7] K.S. Berezowski and S. Vrudhula. Automatic design of binary and multiple-valued logic gates on RTD series. In *8th Euromicro Conference on Digital System Design (DSD'05)*, page 139–142, Aug 2005.
- [8] N. Kulkarni, J. Yang, J. Seo, and S. Vrudhula. Reducing Power, Leakage, and Area of Standard-Cell ASICs Using Threshold Logic Flip-Flops. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 24(9):2873–2886, September 2016.
- [9] F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 1958.
- [10] W. S. McCulloch and W. Pitts. A logical calculus of the ideas immanent in nervous activity. In James A. Anderson and Edward Rosenfeld, editors, *Neurocomputing: Foundations of Research*. MIT Press, 1988.
- [11] A. A. Mullin. Threshold Logic: A Synthesis Approach. *SIAM Review*, 8(3):405–406, Jul 1966.
- [12] V. Beiu, J.M. Quintana, and M.J. Avedillo. VLSI implementations of threshold logic- a comprehensive survey. *IEEE Transactions on Neural Networks*, 14(5):1217–1243, Sep 2003.
- [13] V. Beiu. A survey of perceptron circuit complexity results. In *Proceedings of the International Joint Conference on Neural Networks*, 2003., volume 2, pages 989–994. IEEE, 2003.
- [14] P. Celinski, S. D. Cotozana, J. Lopez, S. F. Al-Sarawi, and D. Abbott. State of the art in CMOS threshold logic VLSI gate implementations and applications. In *VLSI Circuits and Systems*, volume 5117, pages 53–64. SPIE, 2003.
- [15] C. Lageweg, S. Cotozana, and S. Vassiliadis. A full adder implementation using SET based linear threshold gates. In *9th International Conference on Electronics, Circuits and Systems*, volume 2, pages 665–668. IEEE, 2002.
- [16] S.N. Mozaffari, S. Tragoudas, and T. Haniotakis. A Generalized Approach to Implement Efficient CMOS-Based Threshold Logic Functions. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 65(3):946–959, March 2018.
- [17] R. Zhang, P. Gupta, L. Zhong, and N.K. Jha. Threshold Network Synthesis and Optimization and Its Application to Nanotechnologies. *IEEE Transaction On Computer-Aided Design Of Integrated Circuits And Systems*, 24:107–118, 2005.
- [18] V. Annampedu and M.D. Wagh. Decomposition of threshold functions into bounded fan-in threshold functions. *Information and Computation*, 227:84–101, Jun 2013.
- [19] S.N. Mozaffari and S. Tragoudas. Maximizing the Number of Threshold Logic Functions Using Resistive Memory. *IEEE Transactions on Nanotechnology*, 17(5):897–905, September 2018.
- [20] S. Kaya, H.F.A. Hamed, D.T. Ting, and G. Creech. Reconfigurable threshold logic gates with nanoscale DG-MOSFETs. *Solid-State Electronics*, 51(10):1301–1307, October 2007.
- [21] J. Yang, J. Davis, N. Kulkarni, J. Seo, and S. Vrudhula. Dynamic and leakage power reduction of ASICs using configurable threshold logic gates. In *2015 IEEE Custom Integrated Circuits Conference (CICC)*, pages 1–4. IEEE, September 2015.
- [22] R. Fowler and L. Nordheim. Electron Emission in Intense Electric Fields. *Proc. Royal Soc. of London. Series A*, 119(781), May 1928.

- [23] D. Richter. *Fundamentals of Non-Volatile Memories*, page 5–110. Springer Series in Advanced Microelectronics. Springer Netherlands, 2014.
- [24] K. Nii, Y. Taniguchi, and K. Okuyama. A Cost-Effective Embedded Nonvolatile Memory with Scalable LEE Flash®-G2 SONOS for Secure IoT and Computing-in-Memory (CiM) Applications. In *2020 International Symposium on VLSI Design, Automation and Test (VLSI-DAT)*, pages 1–4, Aug 2020.
- [25] S. Tsuda, Y. Kawashima, K. Sonoda, A. Yoshitomi, T. Mihara, S. Narumi, M. Inoue, S. Muranaka, T. Maruyama, T. Yamashita, and et al. First demonstration of FinFET split-gate MONOS for high-speed and highly-reliable embedded flash in 16/14nm-node and beyond. In *2016 IEEE International Electron Devices Meeting (IEDM)*, pages 11.1.1–11.1.4. IEEE, Dec 2016.
- [26] F. Khan, D. Moy, D. Anand, E.H. Schroeder, R. Katz, L. Jiang, E. Banghart, N. Robson, and T. Kirihata. Turning Logic Transistors into Secure, Multi-Time Programmable, Embedded Non-Volatile Memory Elements for 14 nm FinFET Technologies and Beyond. In *2019 Symposium on VLSI Technology*, pages T116–T117. IEEE, Jun 2019.
- [27] F. Khan, E. Cartier, J.C.S. Woo, and S.S. Iyer. Charge Trap Transistor (CTT): An Embedded Fully Logic-Compatible Multiple-Time Programmable Non-Volatile Memory Element for High- k Metal-Gate CMOS Technologies. *IEEE Electron Device Letters*, 38(1):44–47, Jan 2017.
- [28] V. Agrawal, V. Prabhakar, K. Ramkumar, L. Hinh, S. Saha, S. Samanta, and R. Kapre. In-Memory Computing array using 40nm multibit SONOS achieving 100 TOPS/W energy efficiency for Deep Neural Network Edge Inference Accelerators. In *2020 IEEE International Memory Workshop (IMW)*, page 1–4, May 2020.
- [29] M. Abusultan and S.P. Khatri. A flash-based digital circuit design flow. In *2016 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, pages 1–6, November 2016.
- [30] M. Abusultan and S.P. Khatri. Implementing low power digital circuits using flash devices. In *2016 IEEE 34th International Conference on Computer Design (ICCD)*, pages 109–116, October 2016.
- [31] M. Abusultan and S.P. Khatri. A Ternary-Valued, Floating Gate Transistor-Based Circuit Design Approach. In *2016 IEEE Computer Society Annual Symposium on VLSI (ISVLSI)*, pages 719–724, July 2016.
- [32] M. Abusultan and S.P. Khatri. Design of a Flash-based Circuit for Multi-valued Logic. In *Proceedings of the on Great Lakes Symposium on VLSI 2017*, page 41–46. ACM, May 2017.
- [33] V. Bohossian, P. Hasler, and J. Bruck. Programmable neural logic. *IEEE Transactions on Components, Packaging, and Manufacturing Technology: Part B*, 21(4):346–351, November 1998.
- [34] R. Bez, E. Camerlenghi, A. Modelli, and A. Visconti. Introduction to flash memory. *Proceedings of the IEEE*, 91(4):489–502, Apr 2003.
- [35] Y. Cai, Y. Luo, S. Ghose, and O. Mutlu. Read Disturb Errors in MLC NAND Flash Memory: Characterization, Mitigation, and Recovery. In *2015 45th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*, page 438–449. IEEE, Jun 2015.
- [36] S. Boboila and P. Desnoyers. Write endurance in flash drives: measurements and analysis. In *Proceedings of the 8th USENIX conference on File and storage technologies, FAST’10*, page 9. USENIX Association, Feb 2010.
- [37] M. Mehri and N. Masoumi. A thorough investigation into active and passive shielding methods for nano-VLSI interconnects against EMI and crosstalk. *AEU - International Journal of Electronics and Communications*, 69(9):1199–1207, Sep 2015.
- [38] <https://sites.google.com/view/threshold-functions/home/>.
- [39] Cadence. <http://www.cadence.com>.
- [40] R. Degraeve, B. Kaczer, and G. Groeseneken. Degradation and breakdown in thin oxide layers: mechanisms, models and reliability prediction. *Microelectronics Reliability*, 39(10):1445–1460, Oct 1999.
- [41] G. Bersuker, Y. Jeon, and H.R. Huff. Degradation of thin oxides during electrical stress. *Microelectronics Reliability*, 41(12):1923s–1931, Dec 2001.
- [42] S. Dechu, M.K. Goparaju, and S. Tragoudas. A Metric of Tolerance for the Manufacturing Defects of Threshold Logic Gates. In *2006 21st IEEE International Symposium on Defect and Fault Tolerance in VLSI Systems*, page 318–326, Oct 2006.
- [43] A. Neutzling, J.M. Matos, A. Mishchenko, A. Reis, and R.P. Ribas. Effective Logic Synthesis for Threshold Logic Circuit Design. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 38(5):926–937, May 2019.



Arizona State University, Tempe, AZ, USA since 2016. His current research interests include new circuit architectures and design algorithms using threshold logic gates, and their applications to the design of energy efficient digital application-specified integrated circuit, field-programmable gate array, and neural network accelerators.



comm Inc., San Jose, CA, USA as a Hardware Engineering intern. His current research interest includes the design of threshold logic gates, In-memory computing, near memory processing enabling high throughput and energy-efficient systems for data-intensive applications.



publications. Among these papers, five received a best paper award, while six others received best paper nominations. He has coauthored nine research monographs and one edited research monograph, three book chapters, and 13 invited conference papers or workshop papers. He holds six U.S. patents. He was invited to serve as a Panelist at a conference seven times and have presented two conference tutorials. His current research interests include VLSI IC/system-on-a-chip design [including energy efficient design of custom ICs and field-programmable gate arrays (FPGAs), radiation and variation tolerant design, clocking], algorithm acceleration using hardware (FPGA as well as custom IC based) and software (uniprocessor and GPU based), and interdisciplinary extensions of these topics to other areas.



low power circuit design, and energy management of circuits and systems. Specific topics include: energy optimization of battery powered computing systems, including smartphones, wireless sensor networks and IoT systems that rely on energy harvesting; system level dynamic power and thermal management of multicore processors and system-on-chip (SoC); statistical methods for the analysis of process variations; statistical optimization of performance, power and leakage; new circuit architectures of threshold logic circuits for the design of ASICs and FPGAs.

Ankit Wagle (M’17) received the B.S. degree in Electronics and Telecommunication from the University of Pune, Maharashtra, India, in 2013, and the M.S. degree in VLSI Design from Vellore Institute of Technology, Vellore, TN, India, in 2015. He spent his graduate research internships at Intel, Bangalore, KA, India in 2015 and Maxlinear, Carlsbad, CA, USA in 2017. He also worked with Open-Silicon, Bangalore, KA, India from 2015 to 2016. He is currently pursuing the Ph.D. degree with the School of Computing and Augmented Intelligence (SCAI),

Gian Singh (M’19) received his B.Tech degree in Electronics and Communication Engineering from the National Institute of Technology (NIT-H), Hamirpur, India, in 2017. He worked as Project Associate at NIT-H under SMDP-C2SD project sponsored by the Govt. of India from 2017 to 2018. He started his Ph.D. degree at the School of Computing and Augmented Intelligence (SCAI), Arizona State University, Tempe, AZ, the USA in Fall 2018. He spent Fall’19 as an SoC Tech intern at Maxlinear Inc., Carlsbad, CA, USA and Summer’20 at Qualcomm Inc., San Jose, CA, USA as a Hardware Engineering intern. His current research interest includes the design of threshold logic gates, In-memory computing, near memory processing enabling high throughput and energy-efficient systems for data-intensive applications.

Sunil Khatri received the B.Tech. degree in electrical engineering from IIT Kanpur, Kanpur, India, the M.S. degree in electronics and communication engineering from The University of Texas at Austin, Austin, TX, USA, and the Ph.D. degree in electrical engineering and computer sciences from the University of California at Berkeley, Berkeley, CA, USA. He is currently a Professor of Electronics and Communication Engineering at Texas A&M University, College Station, TX, USA. He has authored or coauthored more than 250 peer-reviewed

Sarma Vrudhula (M’85-SM’02-F’16) is a Professor of Computer Science and Engineering with Arizona State University, and the Director of the NSF I/UCRC Center for Embedded Systems. He received the B.Math. degree from the University of Waterloo, Waterloo, ON, Canada, and the M.S.E.E. and Ph.D. degrees in electrical and computer engineering from the University of Southern California, Los Angeles, CA, USA. His work spans several areas in design automation and computer aided design for digital integrated circuit and systems, focusing on