

# Area-Efficient LUT-Based Multipliers for AMD Versal FPGAs

Zetao Miao, Xander Pottier, Jonas Bertels, Wouter Legiest, Ingrid Verbauwhede

*COSIC, KU Leuven*

Leuven, Belgium

{zetao.miao, xander.pottier, jonas.bertels, wouter.legiest, ingrid.verbauwhede}@kuleuven.be

**Abstract**—AMD Versal FPGAs introduce a new CLB micro-architecture featuring the LOOKAHEAD8 carry structure in place of the legacy CARRY4/8 chains, on which existing area-efficient LUT-based multiplier designs map inefficiently. This paper proposes a LUT-based integer multiplier architecture tailored to the Versal fabric. By jointly exploiting radix-4 modified Booth recoding and the new Versal LUT micro-architecture, only  $\sim n^2/4$  LUTs are required to generate the partial-product bit heap for an  $n$ -bit multiplication. A new heuristic for compressor-tree synthesis further improves the area–delay product by 8–20% over state-of-the-art Versal heuristics. Overall, the proposed multipliers achieve up to 40% LUT reduction relative to AMD LogiCORE IP multipliers at comparable critical-path delay. An open-source Python RTL generator with configurable operand widths and pipeline depths is provided for scalable deployment.

**Index Terms**—LUT-based multiplier, radix-4 Booth recoding, compressor tree synthesis, FPGA, Versal

## I. INTRODUCTION

Multiplication is a fundamental and performance-critical operation across many field-programmable gate array (FPGA) applications. Different applications require multipliers of different operand precisions and performance characteristics. For instance, modern neural networks increasingly adopt low-precision integer arithmetic, typically INT8 or even INT4. DSP algorithms such as filtering and spectral analysis commonly require medium-precision fixed-point or floating-point multiplications. In contrast, cryptographic algorithms rely on large integer multiplications, often up to several hundred bits.

To accelerate multiplication, modern FPGAs integrate fixed-size DSP blocks that provide efficient embedded multipliers. However, their fixed operand formats, rigid placement, and limited quantity make them an inflexible, scarce resource. In contrast, look-up-table (LUT)-based multipliers support

This work was supported in part by the Horizon 2020 ERC Advanced Grant (101020005 Belfort) and the CyberSecurity Research Flanders with reference number VOWEICS02. Xander Pottier is funded by FWO (Research Foundation – Flanders) as Strategic Basic (SB) PhD fellow (project number 1S93126N). Wouter Legiest is funded by FWO (Research Foundation – Flanders) as Strategic Basic (SB) PhD fellow (project number 1S57125N).



European Research Council  
Established by the European Commission

arbitrary operand sizes and flexible placement, built from the most abundant computational resource in FPGA fabrics. They can also be combined with DSP blocks to construct wide-precision multipliers. Consequently, efficient LUT-based multiplier design remains critical for precision-diverse FPGA workloads.

Most prior work has focused on optimizing LUT-based multipliers for AMD Xilinx 7-series and UltraScale/UltraScale+ architectures, exploiting the configurable logic block (CLB) micro-architecture of 6-input LUTs coupled with CARRY4/8 primitives to generate and accumulate partial products [1]–[6] or to compress partial-product bit heaps [6]–[8]. The CLB micro-architecture of AMD Versal devices differs substantially, most notably by replacing CARRY4/8 with the LOOKAHEAD8 carry structure [9]. Consequently, many prior designs no longer map efficiently to Versal, motivating new multiplier architectures tailored to Versal CLBs.

This paper presents a LUT-based multiplier architecture designed for Versal CLBs. Partial products are generated efficiently by exploiting the new cascade multiplexers in Versal LUTs, and accumulated through compressor trees synthesized by a new heuristic under the LOOKAHEAD8 carry constraints. We further provide an open-source Python-based RTL generator<sup>1</sup> that produces either standalone compressor trees or complete multipliers, with user-configurable operand widths and pipeline depths. Evaluation shows up to 40% LUT reduction over AMD LogiCORE IP multipliers [10] at comparable critical-path delay, and an 8–20% area–delay product improvement over state-of-the-art Versal compressor-tree heuristics.

## II. BACKGROUND

### A. Versal CLB Adaptations

While Versal CLBs retain the overall organization of LUTs, carry logic, and flip-flops, the internal architectures of both the LUT and the carry chain differ significantly from previous FPGA generations.

Fig. 1 compares the Versal LUT micro-architecture with those of UltraScale/UltraScale+ devices. One major change is the introduction of two independently configurable cascade multiplexers, which control two A5 multiplexers, allowing the A5, A6, or CASC pins to serve as the fifth input of a

<sup>1</sup><https://github.com/KULeuven-COSIC/crypt-arith>

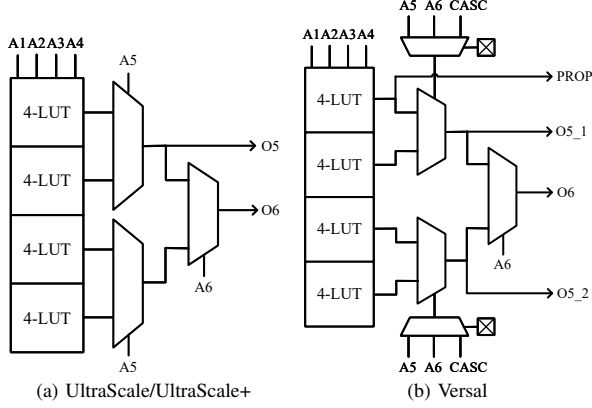


Fig. 1: LUT micro-architectures based on [9].

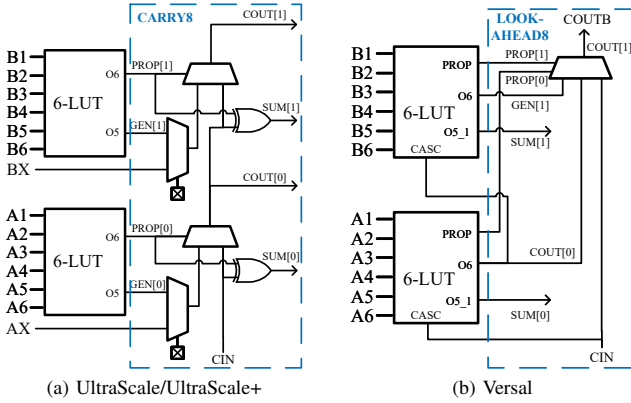


Fig. 2: First two-bit section of carry hardware with LUTs.

Versal LUT. For example, in dual 5-LUT mode, the two 5-input sub-LUTs of an UltraScale/UltraScale+ LUT must share all five inputs (A1–A5), whereas Versal allows them to take independent fifth inputs from A5 and A6. The output of the second 5-input sub-LUT now appears on pin O5\_2 instead of O6. Versal LUTs also expose an extra PROP output, an internal signal dedicated to the LOOKAHEAD8 carry structure.

Versal also introduces a dedicated cascade path between each pair of adjacent LUTs, as illustrated in Fig. 2b: a direct wire connects the O6 output of the lower LUT to the cascade input (CASC) of the upper LUT, enabling fast signal propagation without consuming general-purpose routing resources.

Fig. 2 illustrates the first two-bit section of the UltraScale/UltraScale+ CARRY8 block and the corresponding Versal LOOKAHEAD8 block. Compared to CARRY8, LOOKAHEAD8 removes the XOR gates and multiplexers that previously allowed signals outside the CLB to directly access the carry chain. Propagation signals are now generated by a dedicated 4-input sub-LUT inside each Versal LUT and exposed through the PROP port instead of O6. The carry multiplexers are organized in two-bit sections, with carry propagation within each two-LUT pair routed through the

LUT cascade paths. Each LOOKAHEAD8 block contains four such sections in a carry-lookahead arrangement, detailed in Section IV.

### B. Partial Product Generation

Given two signed integers  $A$  and  $B$  in two's complement form with  $n$  and  $m$  bits respectively, partial-product bits can be generated in several ways. The most straightforward is radix-2 generation, in which each bit is the two-input AND (or NAND) of one bit of  $A$  and one bit of  $B$ , producing approximately  $n \times m$  partial-product bits, with the Baugh–Wooley technique [11] applied to simplify sign extension.

To reduce the number of partial products, radix-4 modified Booth recoding [12] is widely used. It recodes one operand (typically  $B$ ) into radix- $2^2$  digits  $b'_i = -2b_{2i+1} + b_{2i} + b_{2i-1}$ , where  $b_{-1} = 0$ ,  $i \in \{0, 1, \dots, m/2\}$ , and  $B$  is sign-extended to an even number of bits if necessary. The number of partial products is thereby roughly halved, and since each digit satisfies  $b'_i \in \{-2, -1, 0, 1, 2\}$ , the corresponding partial products are obtained by shift-and-add/subtract operations on the multiplicand  $A$ .

### C. Generalized Parallel Counters

Once partial products are generated, they are reduced to fewer operands before the final addition stage. This reduction is typically performed by parallel counters that compress multiple input bits of the same or different weights into fewer output bits while preserving the overall numerical value. Generalized parallel counters (GPCs) provide an abstract representation of such circuits, enabling systematic analysis and construction of compressor trees [13]. A GPC is denoted as  $(p_{m-1}, p_{m-2}, \dots, p_0 : q_{n-1}, q_{n-2}, \dots, q_0)$ , where  $p_i$  and  $q_i$  denote the numbers of input and output bits of weight  $2^i$ , respectively, and the weighted sum of outputs equals that of the inputs. When each output column contains exactly one bit, we abbreviate the notation as  $(p_{m-1}, p_{m-2}, \dots, p_0 : n)$ , where  $n$  is the total number of output bits. For example,  $(1, 5 : 1, 1, 1)$  is written  $(1, 5 : 3)$ .

Two metrics are commonly used to compare GPCs: *efficiency*  $E$  and *strength*  $S$ . Efficiency  $E$ , defined as the ratio of bits reduced to LUTs used, is implementation-dependent and measures how effectively LUTs reduce bits [13]. Strength  $S$ , defined as the ratio of input to output bits, is implementation-independent and captures the inherent bit-reduction capability [14].

$$E = \frac{\sum_{i=0}^{m-1} p_i - \sum_{i=0}^{n-1} q_i}{\#\text{LUTs}}, \quad S = \frac{\sum_{i=0}^{m-1} p_i}{\sum_{i=0}^{n-1} q_i}. \quad (1)$$

Compressor-tree synthesis on FPGAs iteratively applies GPCs across bit positions to reduce bit-heap height until a final adder can resolve the remaining bits. Existing methods fall into two categories: heuristic approaches select GPCs by metrics such as *efficiency* and *strength*, achieving good area-delay trade-offs while being computationally inexpensive [15], [16], [18], whereas optimization-based approaches formulate

TABLE I: Basic Versal GPCs used for Compressor Tree Synthesis in [18]

Counter	#LUTs	E	S	LO.AH. <sup>1</sup>	Row/Column <sup>2</sup>
(3 : 2]	1	1.0	1.5	✓	Both
(6 : 3]	3	1.0	2.0	×	Neither
(10 : 4, 2)	3	1.33	1.67	×	Neither
(2, 5 : 1, 2, 1)	2	1.5	1.75	×	Column
(1, 5 : 3]	2	1.5	2.0	✓	Row
(2, 2, 3 : 4]	2	1.5	1.75	×	Row

<sup>1</sup>if the GPC is compatible with carry-lookahead structure

<sup>2</sup>if the GPC is used to construct row/column counters

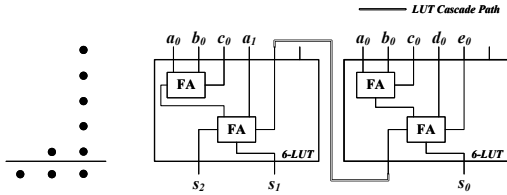


Fig. 3: (1, 5 : 3] GPC in [18].

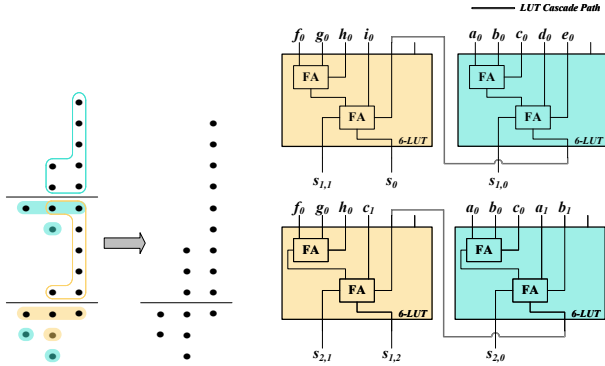


Fig. 4: (3, 9 : 2, 3, 1) GPC in [18], composed of two (1, 5 : 3] GPCs.

the problem as integer linear programming (ILP), yielding optimal or near-optimal solutions but at significant computational expense [16], [17]. Practical FPGA designs therefore mostly rely on heuristic synthesis.

### III. RELATED WORK

#### A. Versal GPCs and Compressor Tree Synthesis

Höbfeld *et al.* [18] were the first to design GPCs for Versal, revisiting UltraScale/UltraScale+ counters and proposing new structures tailored to the Versal architecture. Table I summarizes the basic GPCs in their compressor-synthesis framework. A defining feature of these GPCs is their use of the new LUT cascade path: depending on the construction strategy, the cascade propagates either carry signals, as in the (1, 5 : 3] counter shown in Fig. 3, or sum signals, as in the (3, 9 : 2, 3, 1) counter in Fig. 4.

Beyond basic GPCs, Höbfeld *et al.* [18] compose two larger structures: row counters and column counters. Row counters extend horizontally, with eligible GPCs forwarding their carry

outputs across bit positions, while column counters grow vertically by cascading sum signals to target tall bit heaps. Column counters such as  $(2n+1 : n, 1)$  and  $(n+1, 4n+1 : n, n+1, 1)$  are formed by cascading (3 : 2] and (2, 5 : 1, 2, 1) GPCs, respectively, with rippled sum. For example, the (3, 9 : 2, 3, 1) counter in Fig. 4 (called the dual-rail ripple-sum counter in [18]) is built by cascading two (2, 5 : 1, 2, 1) counters ( $n = 2$ ). Row counters, on the other hand, cascade (1, 5 : 3], (3 : 2], and (2, 2, 3 : 4] GPCs through carry propagation. Column counters are incompatible with the LOOKAHEAD8 carry-lookahead structure, so their GPCs must cascade through general-purpose routing. Among row counters, only (1, 5 : 3] and (3 : 2] are LOOKAHEAD8-compatible; the (2, 2, 3 : 4] GPC is therefore limited, as its cascading likewise relies on general-purpose routing.

To construct compressor trees, Höbfeld *et al.* [18] proposed a stage-wise heuristic that iteratively reduces the height of the bit heap through dynamically selected GPCs. At each stage, candidate counters, including row counters and column counters, are evaluated and selected by prioritizing their efficiency or strength. Experimental results demonstrate significant LUT savings (around 45%) compared to Vivado adder-tree implementations.

#### B. Versal LUT-Based Multipliers Using Gate Absorption

For LUT-based tree multipliers, gate absorption is a compressor-tree optimization that merges radix-2 partial-product logic into the compressors instead of implementing it as a separate pre-processing stage. By exploiting unused LUT inputs in feasible GPCs, it eliminates the additional logic layer in the first compression stage. Höbfeld *et al.* [18] report that absorbing two-input gates into the compressor tree reduces LUT utilization by at least 21%.

## IV. PROPOSED MULTIPLIERS

We propose a tree multiplier with two stages: partial-product bit heap generation followed by GPC-based compression. Both stages are optimized for the Versal fabric.

#### A. Partial-Product Generation

Partial-product generation on FPGAs depends on the chosen scheme. For radix-2 partial products, two AND (or NAND) gates fit in a single LUT5:2 on both UltraScale/UltraScale+ and Versal devices, requiring approximately  $n^2/2$  LUTs for an  $n$ -bit multiplication. Radix-4 modified Booth recoding halves the partial-product count but makes each bit a five-input Boolean function. On UltraScale/UltraScale+ FPGAs, this typically yields a similar LUT cost as radix-2, since one LUT produces only one such bit.

On Versal, however, radix-4 modified Booth recoding can be implemented much more efficiently thanks to the new LUT micro-architecture. The proposed multiplier employs this scheme as summarized in Table II, where  $P'_{i,j}$  denotes a partial-product bit with weight  $2^{2i+j}$  and  $c_i$  a carry bit with weight  $2^{2i}$ . Each signed partial product  $P'_i$  consists of bits





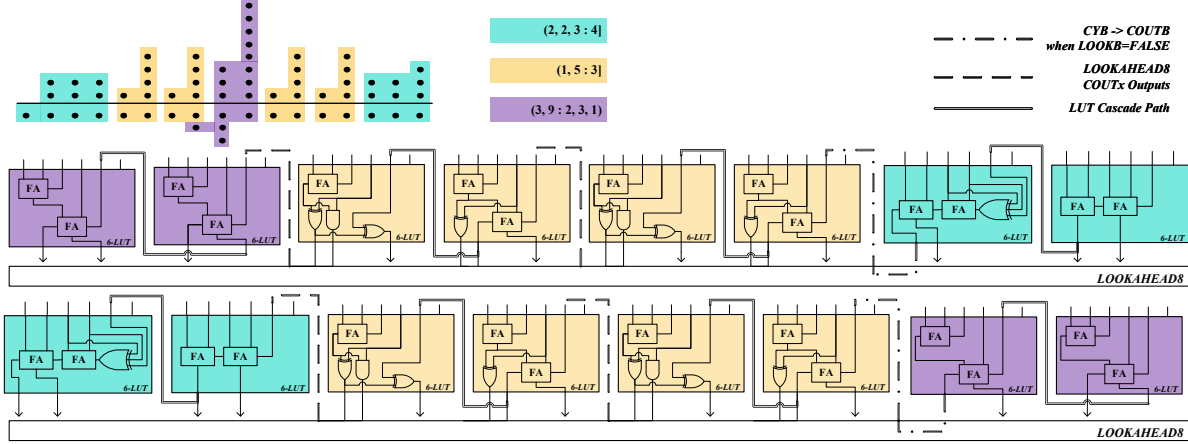


Fig. 8: Row counter benefits from dual-rail property of (3, 9 : 2, 3, 1) GPC.

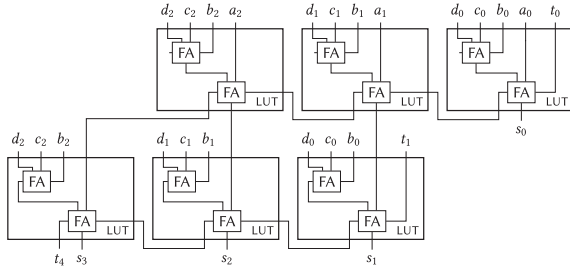


Fig. 9: Quaternary adder implementation in [18].

TABLE V: Counter Necessity Conditions

Counter	Necessity Conditions
(5, 17 : 4, 5, 1)	Always Necessary
(4, 13 : 3, 4, 1)	$H_c \geq 16$
(3, 9 : 2, 3, 1)	$H_c \geq 12$
(9 : 4, 1)	$H_c \geq 12, 5 H_c > 17 H_{c+1}$
(6 : 3]	$H_c = 9, H_{c+1} \leq 3, H_{c+2} \leq 3$
(2, 2, 3 : 4]	$5 \leq H_c \leq 6, 4 \leq H_{c+1} \leq 5, 4 \leq H_{c+2} \leq 5$
(3 : 2]	$5 \leq H_c \leq 6$
(1, 5 : 3]	Always Necessary

the heuristic prioritizes starting or extending a row counter (i.e., cascading eligible GPCs) under the LOOKAHEAD8 constraints of the previous subsection. After placement, if the counter is a row-counter candidate, the scan advances to the column of its MSB output to continue the cascade; otherwise, it moves to the next column. The process repeats until no further counters can be scheduled in the current stage.

Stages are generated iteratively until the reduced bit heap can be finalized by the terminal quaternary adder. A consolidation step is then performed: if counters in the last GPC-compression stage can be under-utilized (i.e., with some inputs supplied by leftover bits from the preceding stage) without preventing terminal addition, the last two GPC-compression stages are merged. This reduces critical-path delay without increasing LUT cost.

### F. RTL Generator for Proposed Multipliers

We developed a Python-based RTL generator for the proposed multipliers. The generator performs compressor-tree synthesis with the proposed heuristic and produces synthesizable SystemVerilog via primitive instantiation and explicit wiring. From user-defined parameters, it generates the partial-product logic and compressor-tree structures, along with XDC constraints, testbenches, and random test vectors; standalone compressor generation is also supported. Pipeline insertion is automatic: for a given pipeline depth, registers are placed across the partial-product generation, GPC-based compression, and terminal-addition logic to balance stage delays. All experimental results in Section V are obtained with designs from this generator.

## V. EXPERIMENTAL RESULTS

### A. Evaluation Method

Designs are synthesized with default settings in Vivado, targeting xcvc1902-vsva2197-2MP-e-S device. They are embedded within a register sandwich to ensure consistent timing evaluation. The reported delay is the critical-path delay at the tightest successful timing closure, found by iteratively tightening the timing constraint in 0.1 ns steps. LUT utilization reported by Vivado is used as the area metric. For fair comparison, Vivado 2023.1 is used when comparing compressors against the designs of Hoßfeld *et al.* [18]; Vivado 2025.2 is used for all other designs.

### B. Evaluation Results of Compressor Trees

Compressor trees are generated using the proposed generator with the same input shapes as those used by Hoßfeld *et al.* [18]; results are shown in Fig. 11.

The proposed heuristic achieves the best area efficiency across all evaluation cases. For single-column bit heaps, it produces critical-path delay comparable to the efficiency-first heuristic; for two-column and multi-column bit heaps, it achieves delay close to the strength-first heuristic. Overall, it

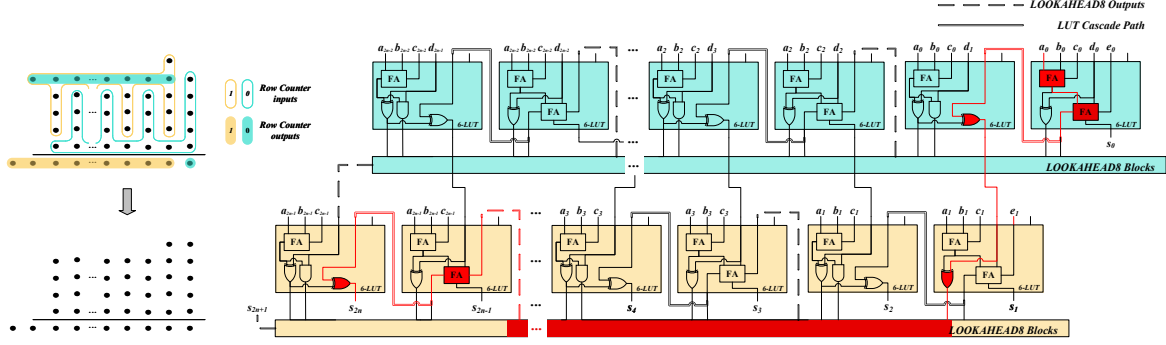


Fig. 10: Proposed implementation of quaternary adder as two layers of row counters.

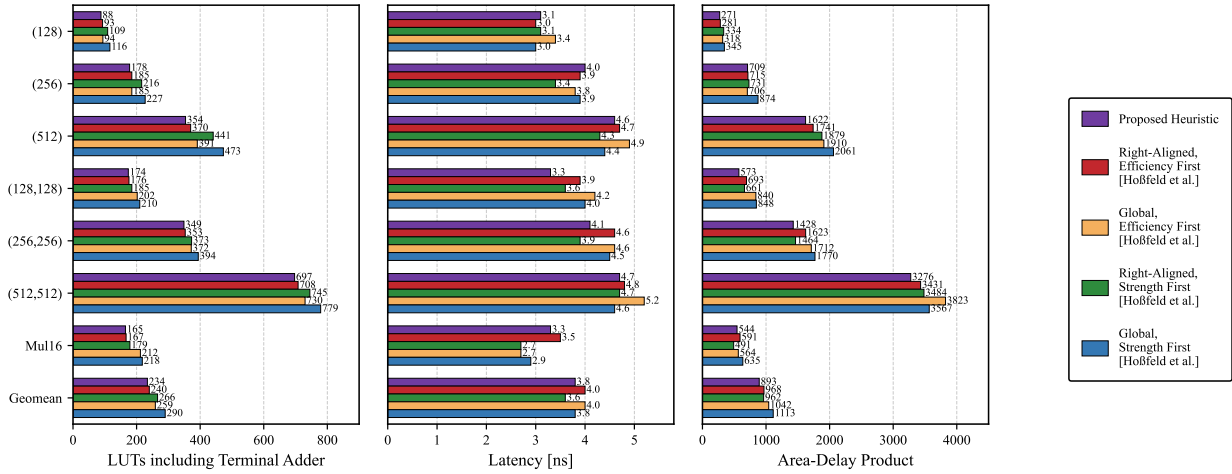


Fig. 11: Comparison of proposed compression heuristic and heuristics used by Hoßfeld *et al.* [18]. (128) denotes a bit heap consisting of a single column with 128 bits, (128,128) denotes a bit heap with two columns containing 128 bits each, and so on. Mul16 denotes the bit heap of 16-bit radix-2 multiplication. Geomean represents the geometric mean of all seven cases. The same Vivado version and Versal device as in [18] are used for fair comparison.

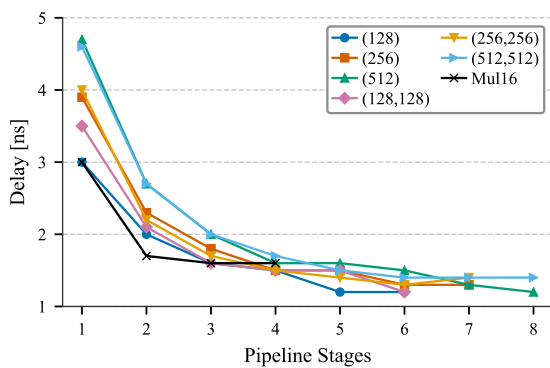


Fig. 12: Critical-path delay of the proposed compressors versus pipeline depth. Results obtained by using Vivado 2025.2.

reduces the area-delay product by 8–20% compared to state-of-the-art heuristics. Fig. 12 reports the critical-path delay

of the proposed compressors under different pipeline depths, demonstrating the effectiveness of the generator’s pipeline insertion.

### C. Evaluation Results of Proposed Multipliers

We compare the proposed multipliers against AMD LogiCORE IP multipliers [10], configured for speed optimization since the area-optimized configuration consumes significantly more LUTs in our evaluation. Results are shown in Fig. 13. The proposed multipliers reduce LUT utilization by up to 40% compared to the speed-optimized LogiCORE designs while maintaining comparable critical-path delay. For operand widths between 28 and 32 bits, the proposed multipliers exhibit slightly longer critical-path delay due to the earlier introduction of an additional compression stage in the compressor tree.

Compared with the 16-bit multiplier with gate absorption from [18], which requires 245 LUTs, the proposed 16-bit multiplier reduces LUT usage by more than 25%, and even the proposed 18-bit multiplier requires over 10% fewer LUTs.

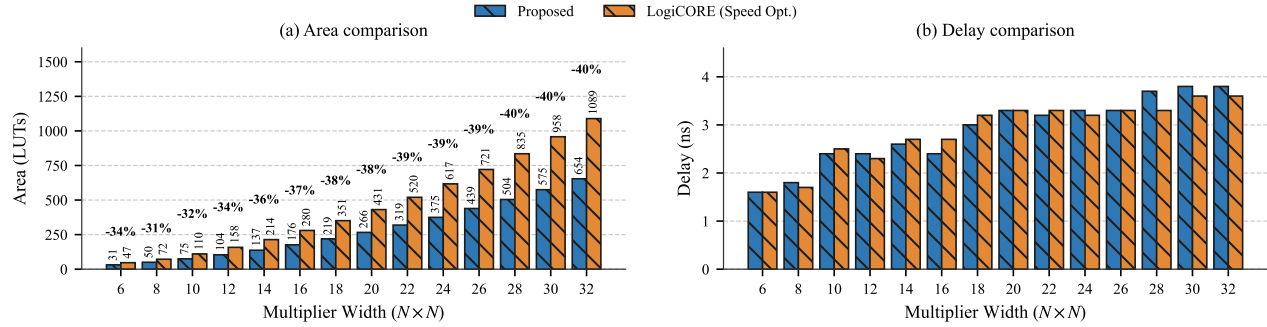


Fig. 13: Comparison of proposed multipliers and AMD LogiCORE IP multipliers (speed optimized).

## VI. CONCLUSION

This paper presents an area-efficient LUT-based integer multiplier architecture tailored to AMD Versal FPGAs. By jointly exploiting radix-4 modified Booth recoding, the Versal LUT micro-architecture, and the LOOKAHEAD8 carry structure, the proposed design enables efficient partial-product generation and compressor-tree construction. We further introduce an architecture-aware heuristic for GPC-based compressor-tree synthesis that improves delay while preserving high area efficiency, and an automated Python-based RTL generator supporting compressor-tree synthesis and pipeline insertion. Experimental results show up to 40% LUT reduction over AMD LogiCORE IP multipliers at comparable delay, and an 8–20% area–delay product improvement over existing Versal compressor-tree heuristics.

Future work will extend the proposed architecture-aware methodology to other arithmetic operators on Versal, further leveraging the new LUT structures and LOOKAHEAD8 carry architecture.

## REFERENCES

- [1] E. G. Walters, “Partial-product generation and addition for multiplication in FPGAs with 6-input LUTs,” in *Proc. 48th Asilomar Conf. on Signals, Systems and Computers*, pp. 1247–1251, 2014.
- [2] M. Kumm, S. Abbas, and P. Zipf, “An efficient software multiplier architecture for Xilinx FPGAs,” in *Proc. IEEE 22nd Symp. on Computer Arithmetic*, pp. 18–25, 2015.
- [3] M. Shu and Q. Liu, “LHAM: Low-cost and high-accuracy approximate multiplier for FPGA-based computing,” *ACM Trans. Reconfigurable Technol. Syst.*, vol. 18, no. 4, Art. no. 48, pp. 1–25, Dec. 2025.
- [4] S. Ullah, S. Rehman, B. S. Prabhakaran, F. Kriebel, M. A. Hanif, M. Shafique, and A. Kumar, “Area-optimized low-latency approximate multipliers for FPGA-based hardware accelerators,” in *Proc. 55th Annu. Design Automation Conf. (DAC)*, pp. 1–6, 2018.
- [5] R. Chen, Y. Lyu, H. Bao, and B. da Silva, “A power-efficient hardware implementation of L-Mul,” arXiv preprint arXiv:2412.18948, 2024.
- [6] A. Böttcher and M. Kumm, “Multiplier design addressing area-delay trade-offs by using DSP and logic resources on FPGAs,” in *Proc. IEEE 35th Int. Conf. on Application-Specific Systems, Architectures and Processors (ASAP)*, pp. 217–225, 2024.
- [7] A. Böttcher and M. Kumm, “Small logic-based multipliers with incomplete sub-multipliers for FPGAs,” in *Proc. IEEE 31st Symp. on Computer Arithmetic (ARITH)*, pp. 124–131, 2024.
- [8] A. Böttcher and M. Kumm, “Towards globally optimal design of multipliers for FPGAs,” *IEEE Trans. Comput.*, vol. 72, no. 5, pp. 1261–1273, 2023.
- [9] AMD, “Versal Adaptive SoC Configurable Logic Block Architecture Manual, AM005,” AMD, Release 1.4, May 14, 2025. [Online]. Available: <https://docs.amd.com/r/en-US/am005-versal-clb/Overview>
- [10] AMD, “Multiplier v12.0 Product Guide (PG108),” AMD, Nov. 18, 2015. [Online]. Available: <https://docs.amd.com/v/u/en-US/pg108-mult-gen>
- [11] C. R. Baugh and B. A. Wooley, “A two’s complement parallel array multiplication algorithm,” *IEEE Trans. Comput.*, vol. C-22, no. 12, pp. 1045–1047, Dec. 1973.
- [12] O. L. MacSorley, “High-speed arithmetic in binary computers,” *Proc. IRE*, vol. 49, no. 1, pp. 67–91, 1961.
- [13] H. Parandeh-Afshar, A. Neogy, P. Brisk, and P. Ienne, “Compressor tree synthesis on commercial high-performance FPGAs,” *ACM Trans. Reconfigurable Technol. Syst.*, vol. 4, no. 4, pp. 1–19, 2011.
- [14] T. B. Preußer, “Generic and universal parallel matrix summation with a flexible compression goal for Xilinx FPGAs,” in *Proc. 27th Int. Conf. on Field Programmable Logic and Applications (FPL)*, pp. 1–7, 2017.
- [15] H. Parandeh-Afshar, P. Brisk, and P. Ienne, “Efficient synthesis of compressor trees on FPGAs,” in *Proc. Asia and South Pacific Design Automation Conf. (ASP-DAC)*, pp. 138–143, 2008.
- [16] M. Kumm and J. Kappauf, “Advanced compressor tree synthesis for FPGAs,” *IEEE Trans. Comput.*, vol. 67, no. 8, pp. 1078–1091, 2018.
- [17] H. Parandeh-Afshar, P. Brisk, and P. Ienne, “Improving synthesis of compressor trees on FPGAs via integer linear programming,” in *Proc. Design, Automation and Test in Europe Conf. (DATE)*, pp. 1256–1261, 2008.
- [18] K. Hoffeld, H. J. Damsgaard, J. Nurmi, M. Blott, and T. B. Preußer, “High-efficiency compressor trees for latest AMD FPGAs,” *ACM Trans. Reconfigurable Technol. Syst.*, vol. 17, no. 2, pp. 1–32, 2024.