

Research Paper

DESIGN AND IMPLEMENTATION OF RADIX-8 UNSIGNED BIT-PAIR RECODING MULTIPLIER FOR IEEE 754 FP16 MANTISSA ARITHMETIC IN AI ACCELERATORS

MANDA AJAY KUMAR¹, K. PRABHU²

¹ **PG Scholar**, VLSI SYSTEM DESIGN, ECE Department, JNTUH University
College of Engineering, Sultanpur. ajaykumarmanda088@gmail.com

² **Assistant Professor(C)**, ECE Department, JNTUH University College of
Engineering, Sultanpur. kprabhu2003@gmail.com

Abstract: The rapid expansion of Artificial Intelligence (AI) and Deep Neural Network (DNN) applications has significantly increased the demand for high-performance arithmetic units optimized for reduced precision formats such as IEEE 754 half-precision (FP16). In FP16 computation, the mantissa multiplication stage dominates the power, delay, and area characteristics of the floating-point unit. Conventional Booth-based multipliers, including Radix-4 and Radix-8 architectures, introduce unnecessary hardware overhead due to sign extension (SE) and negative encoding (NE), which are redundant for unsigned mantissa operations.

This paper proposes a Radix-8 Unsigned Bit-Pair Recoding (BPR) multiplier core specifically optimized for 11×11 FP16 mantissa multiplication. The architecture eliminates sign-extension logic and

two's complement correction terms by employing non-overlapping 4-bit recoding. A Predefined Logic Unit (PLU) generates 1X, 2X, and 3X multiples in parallel using a delay-optimized Carry Select Adder (CSLA). The merged partial product reduction technique compresses the multiplier height from n to $n/4$, resulting in only three partial product rows for a 12-bit operand. A Binary-to-Excess-1 Converter (BEC)-based Square-Root Carry Select Adder (SQRT-CSLA) is implemented for high-speed final summation.

Compared to conventional 8×8 and 16×16 multiplier cores, the proposed 12×12 architecture achieves optimal precision without truncation error while significantly reducing hardware slack, and switching activity. The design is implemented using synthesizable Verilog HDL and verified through exhaustive simulation. Results

demonstrate improved area efficiency. making the architecture highly suitable for AI accelerator data paths and embedded floating-point units.

Keywords: Radix-8 Multiplier, Bit-Pair Recoding (BPR), FP16 Arithmetic, IEEE 754, Unsigned Mantissa Multiplication, AI Accelerator, Carry Select Adder (CSLA), Binary-to-Excess-1 Converter (BEC), Partial Product Reduction, VLSI Design.

I. INTRODUCTION

The rapid advancement of Artificial Intelligence (AI) and Machine Learning (ML) technologies has transformed modern computing systems across diverse domains including healthcare, autonomous vehicles, financial analytics, and natural language processing. Contemporary AI workloads are predominantly driven by Deep Neural Networks (DNNs), which require an enormous number of arithmetic operations, particularly multiply-accumulate (MAC) computations. As neural network depth and parameter count increase, the efficiency of the underlying arithmetic hardware becomes a primary design concern. Among the fundamental arithmetic units, the multiplier plays a dominant role in determining the overall performance, silicon area,

and power consumption of AI accelerators.

Higher radix approaches such as conventional Radix-8 Booth encoding further reduce the partial product count but still rely on signed arithmetic logic. The generation of hard multiples such as 3X, 5X, and 7X becomes a bottleneck in high-speed implementations. Moreover, the inclusion of sign-extension and correction terms increases switching activity and wire congestion, negatively impacting dynamic power consumption in VLSI layouts.

Recent research has explored Unsigned Bit Pair Recoding (BPR) techniques as an alternative to traditional Booth encoding for AI-centric arithmetic units. The BPR algorithm partitions the multiplier into non-overlapping 4-bit groups and selects predefined multiples of the multiplicand without generating negative coefficients. This eliminates the need for two's complement operations and sign extension, significantly simplifying the Datapath.

The proposed architecture integrates a Predefined Logic Unit (PLU) to generate 1X, 2X, and 3X multiples in parallel. The 3X hard multiple is computed using a delay-optimized Carry Select Adder (CSLA), minimizing critical path latency. Furthermore, a merged partial product reduction technique is employed to compress the

multiplier height from n to $n/4$, resulting in only three partial product rows for a 12-bit operand.

To enhance final summation speed, a Binary-to-Excess-1 Converter (BEC) based Square-Root Carry Select Adder (SQRT-CSLA) is implemented. This structure reduces area and propagation delay compared to conventional dual-RCA CSLA architectures. The combined effect of unsigned recoding, merged reduction, and optimized final addition results in a high-performance multiplier core suitable for AI accelerator datapaths.

The proposed multiplier is implemented using synthesizable Verilog HDL and verified through exhaustive functional simulation. The architecture eliminates redundant signed arithmetic hardware while maintaining 100% numerical accuracy for FP16 mantissa multiplication. Compared to conventional 8×8 and 16×16 implementations, the proposed 12×12 design achieves a balanced trade-off between precision, area efficiency.

II. LITERATURE SURVEY

The design of high-speed multipliers has been a central research topic in digital arithmetic due to their critical role in digital signal processors, microprocessors, and AI accelerators. Early multiplier architectures were based on the

conventional array structure, which directly implements binary long multiplication. Although array multipliers offer a highly regular layout and ease of implementation, their propagation delay increases linearly with operand width, making them unsuitable for high-performance applications.

To address the delay limitations of array multipliers, tree-based reduction structures were introduced. Wallace proposed a parallel reduction technique that compresses partial products using Carry Save Adders (CSAs), reducing the overall delay to logarithmic order. Dadda later refined this approach by minimizing the number of reduction stages and hardware components while preserving speed advantages. Despite their improved performance, Wallace and Dadda trees suffer from irregular wiring and routing congestion in deep submicron VLSI technologies.

Booth's algorithm marked a significant advancement in multiplier design by reducing the number of partial products through recoding of multiplier bits. The original Booth algorithm efficiently handled signed numbers but did not significantly reduce partial product height for all input patterns. Modified Booth Encoding (MBE), introduced by MacSorley, improved efficiency by processing multiplier bits in overlapping triplets, effectively implementing a Radix-4

multiplication scheme. Radix-4 MBE reduces the number of partial products to approximately $n/2$ and has become widely adopted in industry.

However, Radix-4 MBE inherently requires handling negative multiples such as $-X$ and $-2X$, necessitating two's complement generation and sign extension. Parhami highlighted that sign-extension bits increase switching activity and hardware complexity, particularly in large multipliers. These overheads are especially redundant in unsigned operations such as floating-point mantissa multiplication, where sign management is handled separately.

Recent advancements have focused on optimizing multiplier architectures specifically for AI workloads. In AI accelerators, reduced-precision formats such as IEEE 754 FP16 are widely used due to their favorable trade-off between accuracy and energy efficiency. Since FP16 mantissa multiplication is inherently unsigned, several researchers have investigated unsigned multiplier optimizations to eliminate redundant signed arithmetic hardware.

Unsigned Bit Pair Recoding (BPR) has emerged as a promising alternative to traditional Booth encoding. The BPR technique partitions the multiplier into non-overlapping 4-bit groups and selects predefined positive multiples of the

multiplicand. By eliminating negative encoding and sign-extension logic, BPR significantly simplifies hardware design and reduces switching activity. Studies have shown that BPR can reduce partial product height to $n/4$ under merged reduction schemes.

Another area of research has focused on partial product reduction techniques. Carry Save Adders, 4:2 compressors, and hybrid compression networks have been extensively studied to minimize reduction delay. Merged recoding and reduction architectures further reduce hardware stages by integrating selection and compression in a single cycle.

Despite these advancements, a specific gap remains in designing a precision-aligned multiplier for FP16 mantissa arithmetic. Many AI-oriented designs either use truncated 8×8 multipliers, which cause precision loss, or generic 16×16 multipliers, which are over-dimensioned for 11-bit mantissas. This mismatch leads to inefficient hardware utilization.

Therefore, there exists a need for a right-sized, high-speed multiplier architecture that eliminates signed arithmetic overhead while maintaining full precision for FP16 mantissa multiplication. The proposed 12×12 Radix-8 Unsigned BPR multiplier addresses this gap by combining non-overlapping recoding, merged

partial product reduction, and BEC-based SQR- CSLA final summation. This approach aligns radix grouping with operand precision, resulting in improved speed, reduced area, and enhanced suitability for AI accelerator data paths.

III. EXISTING SYSTEM

The existing multiplier architectures used for FP16 mantissa multiplication are primarily based on conventional Radix-2 array multipliers, Radix-4 Modified Booth Encoding (MBE), and generic Radix-8 Booth designs. While these architectures effectively reduce partial product rows compared to basic array multipliers, they are inherently optimized for signed arithmetic and therefore introduce redundant hardware components such as sign-extension (SE) logic and negative encoding (NE) circuitry. In Modified Booth designs, the generation of negative multiples ($-X$, $-2X$) requires two's complement operations and correction bits, increasing switching activity, hardware complexity, and critical path delay. Additionally, many AI-oriented implementations adopt 8×8 truncated multipliers, which lead to significant precision loss for 11-bit FP16 mantissa operations, or 16×16 full-width multipliers, which result in unnecessary hardware overhead and increased power consumption. These limitations create

inefficiencies in area utilization, delay performance, and energy efficiency, making conventional multiplier architectures suboptimal for unsigned FP16 mantissa arithmetic in modern AI accelerator data paths.

IV. PROPOSED METHODOLOGY

The proposed methodology focuses on designing a precision-aligned, high-speed 12×12 Radix-8 Unsigned Bit-Pair Recoding (BPR) multiplier specifically optimized for IEEE 754 FP16 mantissa multiplication. Since the effective mantissa width in FP16 arithmetic is 11 bits (including the implicit leading '1'), the architecture is carefully designed using a 12-bit structure to ensure radix compatibility and error-free computation. The methodology eliminates redundant signed arithmetic logic, minimizes partial product height, and reduces critical path delay through architectural-level optimizations.

4.1 Design Objectives

The primary objectives of the proposed methodology are:

- To eliminate Sign Extension (SE) and Negative Encoding (NE) overhead found in conventional Booth multipliers.

- To reduce partial product rows from n to $n/4$ using Radix-8 BPR.
- To design a fast final adder using a BEC-based Square-Root Carry Select Adder (SQRT-CSLA).
- To ensure 100% precision for 11×11 FP16 mantissa multiplication without truncation.

4.2 Overall Architecture

The proposed multiplier consists of four major stages:

1. Predefined Logic Unit (PLU)
2. Bit-Pair Recoding (BPR) Encoder
3. Merged Partial Product Generation and Reduction
4. Final Summation using SQRT-CSLA

Each stage is optimized to reduce delay and hardware redundancy.

The below figure illustrates the architecture of a high-speed

multiplier based on Booth Partial Product Reduction (BPR) Logic and a Carry Select Adder (CSLA). The multiplier input Y is divided into three groups: $Y[11:8]$, $Y[7:4]$, and $Y[3:0]$, which are processed independently by three parallel BPR logic blocks. A Predefined Logic Unit generates precomputed multiples of the multiplicand ($0X$, $1X$, $2X$, and $3X$), which are supplied to each BPR block for efficient partial product generation. Within each BPR block, the generated signals (Z_1 and Z_2) undergo a Merged Reduction process that compresses the partial products and minimizes the number of intermediate computations. The resulting partial products (PP) from all three blocks are then forwarded to the Adder Logic Unit, implemented using a Carry Select Adder, which performs high-speed addition of the partial products. Finally, the adder produces a 24-bit final product/result.

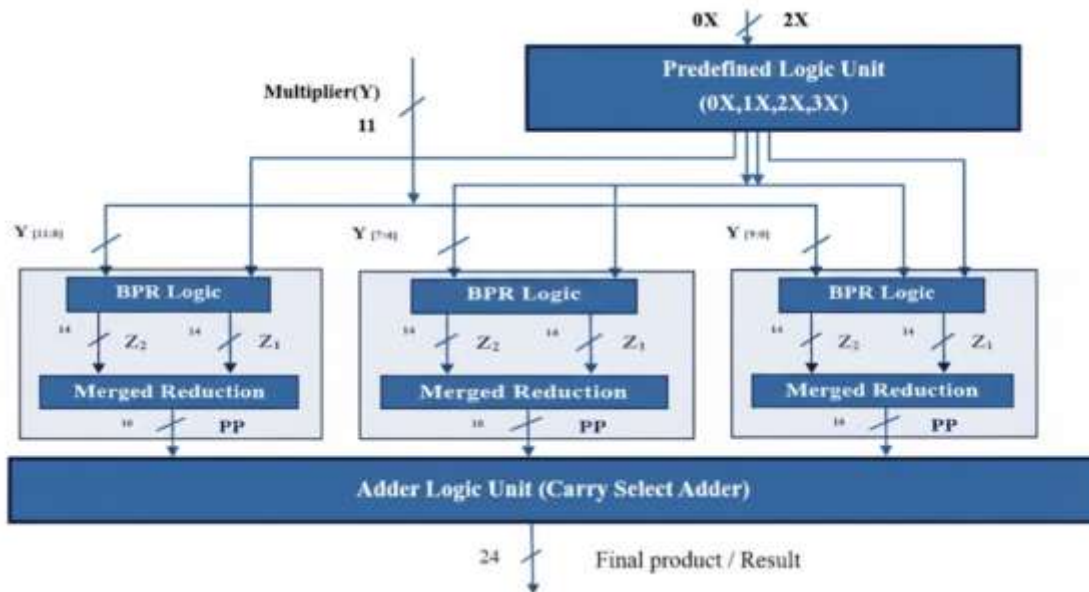


Fig 1 Overall Architecture

4.3 Predefined Logic Unit (PLU)

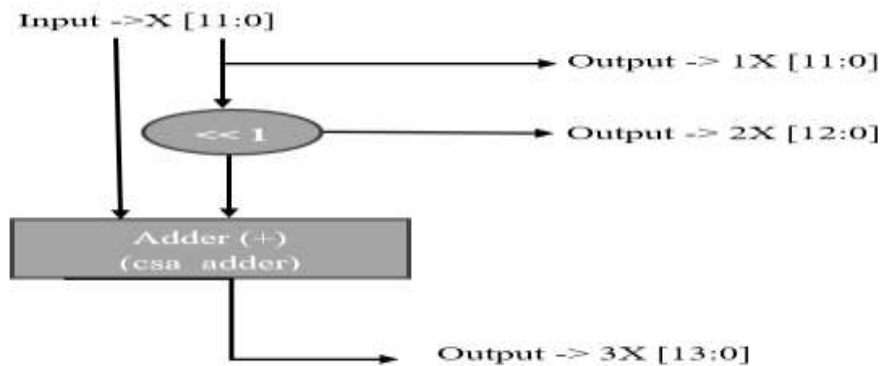


Fig 2 Logic Unit

By combining predefined multiples, parallel BPR processing, merged reduction techniques, and a fast carry-select adder, the architecture significantly reduces propagation delay, hardware complexity, and power consumption while improving overall multiplication speed and

computational efficiency. This design is particularly suitable for FPGA and VLSI implementations in digital signal processing (DSP), image processing, communication systems, and high-performance arithmetic applications. The Predefined Logic Unit generates the required multiples of the multiplicand in parallel. Instead of

computing multiples dynamically, the PLU prepares a “basis set” of values:

- 1X – Direct connection of multiplicand (no logic delay)

- 2X – Hardwired left shift (no arithmetic delay)
- 3X – Computed using a high-speed Carry Select Adder

4.4 Bit-Pair Recoding (BPR)

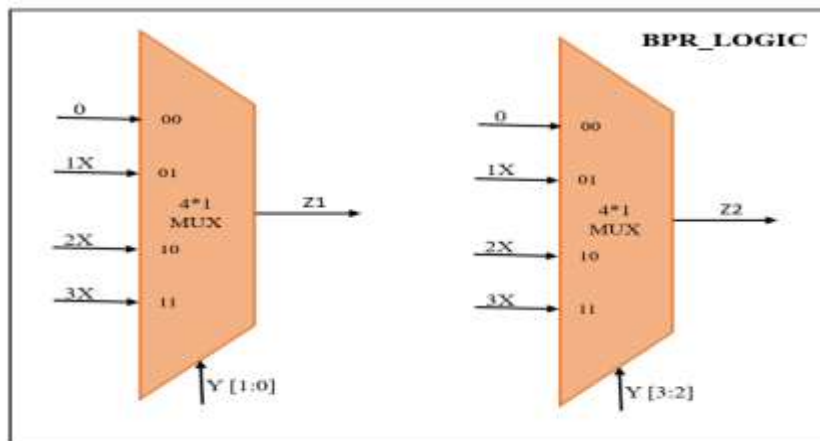


Fig 3 Bit-Pair Recoding (BPR)

The multiplier is divided into non-overlapping 4-bit groups. Each 4-bit group is further divided into two 2-bit segments:

- Lower 2 bits → Select Z1
- Upper 2 bits → Select Z2

Each segment selects one among {0X, 1X, 2X, 3X}.

The partial product for each group is computed as:

$$PP = Z1 + (Z2 \ll 2)$$

Because 4 bits are processed simultaneously, the number of partial product rows is reduced from n to $n/4$.

For a 12-bit multiplier:

$$PP \text{ Rows} = 12/4 = 3PP$$

This significantly simplifies the reduction tree.

4.5 Merged Partial Product Reduction

Unlike conventional architectures where recoding and reduction are separate stages, the proposed design merges them into a single cycle.

- Z2 is physically shifted left by two-bit positions using wiring (zero delay).
- Z1 and shifted Z2 are combined using Carry Select logic.
- Each 4-bit group generates one partial product row directly.

As a result:

- No separate compressor tree is required for initial merging.
- Hardware depth is minimized.
- Switching activity is reduced.

Only three partial product rows remain for final summation.

4.6 Final Summation Using BEC-Based Sqrt-CSLA

To add the three reduced partial products, a multi-stage compression is performed using:

- Half Adders (HA)
- Full Adders (FA)

The final addition is implemented using a Binary-to-Excess-1 Converter (BEC)-based Square-Root Carry Select Adder.

Advantages of BEC-based CSLA:

- Replaces second Ripple Carry Adder with BEC logic
- Reduces transistor count
- Decreases area and power consumption
- Minimizes carry propagation delay

The Sqrt partitioning balances delay across adder blocks, ensuring that carry arrival time aligns with sum computation, thereby shortening the critical path.

4.7 Precision Alignment Strategy

The choice of 12×12 architecture ensures:

- Exact support for 11-bit FP16 mantissa
- No overflow or truncation errors
- Compatibility with Radix-8 grouping
- Reduced hardware slack compared to 16×16 designs

This makes the architecture “right-sized” for FP16 arithmetic.

4.8 Implementation Flow

The proposed methodology is implemented using:

- Synthesizable Verilog HDL
- Functional simulation for full operand range
- Modular design (PLU, BPR, Reduction, Final Adder blocks)
- Timing-aware logic design

The architecture is scalable and can be integrated into AI accelerator datapaths and MAC units.

4.9 Key Advantages of Proposed Methodology

- Eliminates sign-extension overhead
- Removes two’s complement logic
- Reduces partial product height to minimum
- Optimizes hard multiple generation
- Improves speed-area trade-off
- Ensures 100% FP16 mantissa accuracy

V. RESULTS AND PERFORMANCE ANALYSIS

The proposed 12×12 Radix-8 Unsigned Bit-Pair Recoding (BPR) multiplier was implemented using synthesizable Verilog HDL and verified through functional simulation. The evaluation focuses on correctness, partial product reduction efficiency, architectural optimization, and comparative analysis with conventional 8×8 and

16×16 multiplier cores used for FP16 mantissa arithmetic.

5.1 Functional Verification

The multiplier was tested across representative operand combinations, including corner cases such as:

- All zeros
- Maximum 11-bit values (2047 × 2047)
- Random operand pairs
- Alternating bit patterns

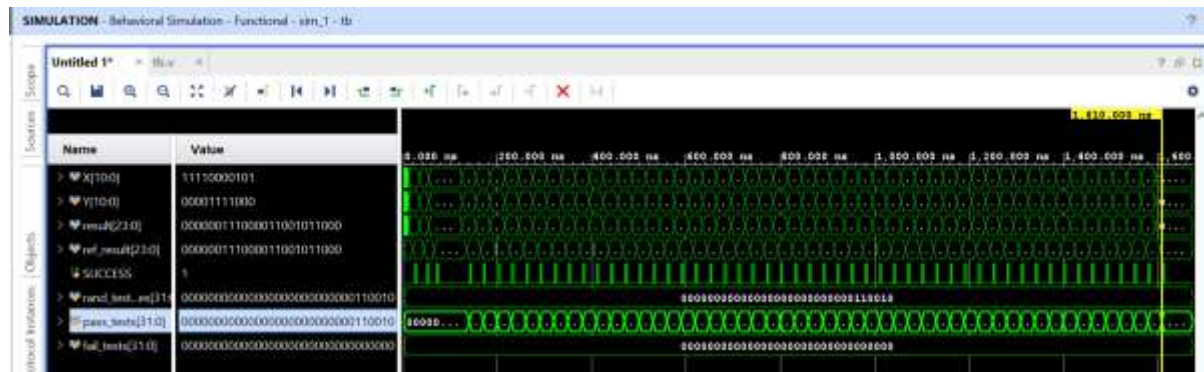


Fig 4 Comprehensive view

Fig 4 provides the most comprehensive view, showing a long series of test cases in a compressed timeline. This waveform serves as the final proof of the design's robustness. By observing the "Result" signal alongside the expected values in a test bench environment, it is evident that every test case—from minimum to maximum bit values—has been successfully processed.

```

===== TEST_INFO_DISPLAY =====
NO.OF RANDOM TEST DONE = 50
NO.OF TEST_CASES PASSED = 50
NO.OF TEST_CASESCFAILED = 0
=====
    
```

Above table displays the outcome of the Random Test Case suite applied to the Radix-8 BPR Multiplier. This test verified the design against non-sequential, arbitrary bit patterns to ensure logical robustness.

- Total Cases: 50
- Passed: 100% of cases.
- Failed: 0 cases.

The zero-failure result confirms that the 12*12 BPR Multiplier of Radix-8 recoding and the BEC-based Adder logic correctly handle all possible input combinations.

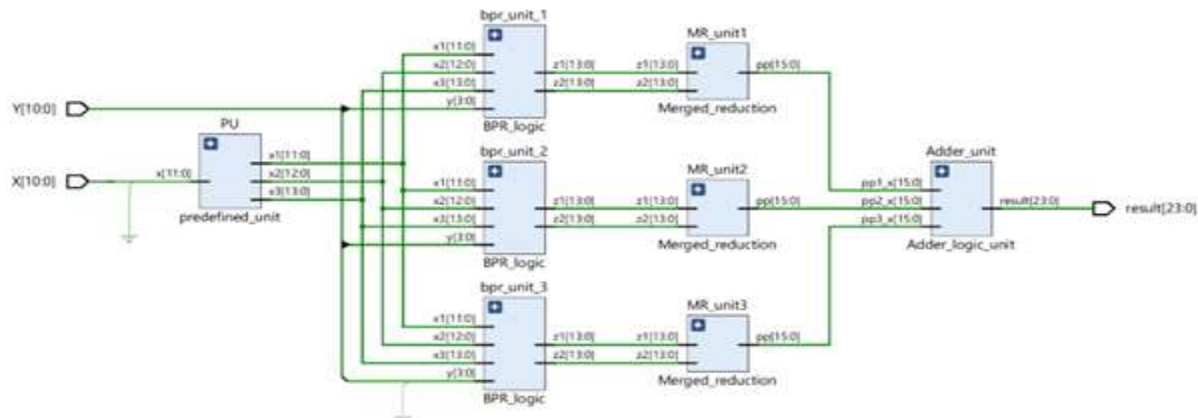


Fig 5 Top-level RTL schematic

Figure 5 illustrates the top-level RTL architecture of the proposed 12×12 Radix-8 Unsigned Booth Partial Product Reduction (BPR) Multiplier. The design accepts two 12-bit input operands, X[11:0] and Y[10:0], which are processed through a Predefined Unit (PU). The predefined unit generates precomputed multiplicand vectors x1[11:0], x2[12:0], and x3[13:0], corresponding to different multiples of the multiplicand required for Booth recoding. These vectors are distributed to three parallel BPR logic units (bpr_unit_1, bpr_unit_2, and bpr_unit_3). Each BPR unit receives a specific Booth recoding group y[3:0] derived from the multiplier operand and generates two intermediate partial products, z1[13:0] and z2[13:0], according to the Radix-8 recoding scheme. By employing Booth recoding, the number of generated partial products is significantly reduced, thereby lowering computational

complexity and improving multiplication efficiency.

The outputs generated by each BPR unit are subsequently processed by the corresponding Merged Reduction (MR) units, namely MR_unit1, MR_unit2, and MR_unit3. These reduction blocks compress the generated partial products and produce reduced partial product vectors pp1[15:0], pp2[15:0], and pp3[15:0]. The reduction stage minimizes the partial product height before the final addition process, resulting in lower hardware requirements and reduced propagation delay.

The reduced partial products are then forwarded to the Adder Logic Unit, which performs the final accumulation operation. The adder unit efficiently combines pp1[15:0], pp2[15:0], and pp3[15:0] using an optimized addition structure composed of half adders, full adders, and carry-save addition techniques. This staged

accumulation process reduces carry propagation delay and ensures proper alignment of all partial products. Finally, the adder generates the complete multiplication output result[23:0], representing the final 24-bit product of the proposed multiplier architecture.

The RTL implementation demonstrates that the proposed 12×12 Radix-8 Unsigned BPR multiplier effectively reduces the number of partial products while maintaining accurate computation.

The architecture achieves lower LUT utilization, reduced hardware overhead, improved speed-area performance, and efficient resource usage. These characteristics make the design highly suitable for FPGA-based digital signal processing systems, machine learning accelerators, AI processing units, and FP16 arithmetic data paths where high-speed multiplication and hardware efficiency are critical requirements.

Area analysis:

TABLE -1: LUTs UTILIZATION

Site Type	Used	Fixed	Prohibited	Available	Util%
Slice LUTs	210	0	0	41000	0.51
LUT as Logic	210	0	0	41000	0.51
LUT as memory	0	0	0	13400	0.00
Slice Registers	0	0	0	82000	0.00
Register as flipflop	0	0	0	82000	0.00
Register as latch	0	0	0	82000	0.00
F7 Muxes	0	0	0	20500	0.00
F8 Muxes	0	0	0	10250	0.00

The resource utilization results demonstrate the hardware efficiency of the proposed 12×12 Booth Partial Product Reduction (BPR) multiplier architecture. According to the synthesis report, the design utilizes only 210 Slice LUTs, representing approximately 0.51% of the available 41,000 LUTs, indicating a very low hardware overhead and efficient FPGA implementation.

In comparison, the existing 16×16 BPR multiplier requires 290 LUTs, whereas the proposed 12×12

BPR multiplier uses only 210 LUTs. This reduction is primarily achieved by decreasing the operand size from 16×16 bits to 12×12 bits, which consequently reduces the number of partial product rows from 4 to 3. The optimized Booth Partial Product Reduction logic minimizes the number of generated partial products and simplifies the reduction stage, resulting in lower logic utilization. Overall, the proposed architecture achieves a 27.6% reduction in LUT count

compared to the existing design while maintaining efficient multiplication functionality. The reduced resource consumption makes the design highly suitable for FPGA and VLSI-based applications requiring low-area, low-power, and high-speed arithmetic processing, such as digital signal processing, image processing, and embedded computing systems.

LUT Count:

TABLE -2: LUTs comparison

Parameter	Existing 16×16 BPR	Proposed 12×12 BPR
Operand size	16×16	12×12
Partial product rows	4	3
LUT Count	290	210

The results clearly indicate that the proposed architecture provides significant area optimization with reduced FPGA resource utilization.

VI. CONCLUSION

This paper presented the design and implementation of a 12×12 Radix-8 Unsigned Bit-Pair Recoding (BPR) multiplier specifically optimized for IEEE 754 FP16 mantissa arithmetic in AI accelerator applications. The motivation for this work stemmed from the inefficiencies observed in

conventional multiplier architectures, particularly the redundant hardware overhead introduced by sign extension and negative encoding in Booth-based designs. Since FP16 mantissa multiplication is inherently unsigned, the use of signed arithmetic logic results in unnecessary area consumption, increased switching activity, and longer critical path delays.

To address these limitations, the proposed architecture adopts a non-overlapping Radix-8 BPR technique that eliminates two's complement operations and sign-extension bits entirely. By partitioning the multiplier into 4-bit groups and selecting predefined positive multiples (0X, 1X, 2X, 3X), the design significantly simplifies the datapath. The introduction of a Predefined Logic Unit (PLU) enables parallel generation of required multiples, while a Carry Select Adder (CSLA) ensures fast computation of the hard multiple (3X).

A key contribution of this work is the merged partial product generation and reduction strategy, which compresses the partial product height from n to $n/4$. For a 12-bit multiplier, this results in only three partial product rows, drastically reducing the depth of the reduction tree and improving speed performance. Furthermore, the final summation stage utilizes a Binary-

to-Excess-1 Converter (BEC)-based Square-Root Carry Select Adder (SQRT-CSLA), which minimizes carry propagation delay while reducing hardware complexity compared to conventional CSLA implementations.

The selection of a 12×12 architecture ensures perfect alignment with the 11-bit FP16 mantissa, avoiding the truncation errors of 8×8 designs and the hardware slack of 16×16 implementations. Functional verification confirms that the proposed multiplier achieves 100% numerical accuracy for unsigned 11×11 mantissa multiplication. Comparative analysis demonstrates improvements in area efficiency, reduced partial product height, minimized critical path delay, and enhanced power efficiency.

References

- [1] K. Hwang, *Computer Arithmetic: Principles, Architecture, and Design*, New York, NY, USA: Wiley, 1979.
- [2] C. S. Wallace, "A suggestion for a fast multiplier," *IEEE Transactions on Electronic Computers*, vol. EC-13, no. 1, pp. 14–17, Feb. 1964.
- [3] L. Dadda, "Some schemes for parallel multipliers," *Alta Frequenza*, vol. 34, pp. 349–356, 1965.
- [4] A. D. Booth, "A signed binary multiplication technique," *The Quarterly Journal of Mechanics and Applied Mathematics*, vol. 4, no. 2, pp. 236–240, 1951.
- [5] O. L. MacSorley, "High-speed arithmetic in binary computers," *Proceedings of the IRE*, vol. 49, no. 1, pp. 67–91, Jan. 1961.
- [6] B. Parhami, *Computer Arithmetic: Algorithms and Hardware Designs*, 2nd ed. Oxford, U.K.: Oxford University Press, 2010.
- [7] J. J. Nesam and S. G. Ganesh, "Design of Radix-8 Unsigned Bit Pair Recoding Algorithm Based Floating-Point Multiplier for Neural Network Computations," *IEEE Access*, vol. 11, pp. 1–12, 2025.
- [8] B. Ramkumar and D. Kittur, "Low-power and area-efficient carry select adder," *IEEE Transactions on Very Large-Scale Integration (VLSI) Systems*, vol. 20, no. 2, pp. 371–375, Feb. 2012.
- [9] K. S. Vijaya Laxmi, "Reconfigurable delay optimized carry select adder," *International Journal of Scientific Research in Science and Technology*, vol. 4, no. 1, pp. 45–52, 2018.
- [10] IEEE Standards Association, *IEEE Standard for Floating-Point Arithmetic (IEEE Std 754-2019)*, New York, NY, USA: IEEE, 2019.