

# RAPID: AppRoximAte Pipelined Soft Multipliers and Dividers for High-Throughput and Energy-Efficiency

Zahra Ebrahimi, Muhammad Zaid, Mark Wijtvliet, Akash Kumar, *Senior Member, IEEE*

**Abstract**—The rapid updates in error-resilient applications along with their quest for high throughput has motivated designing fast approximate functional units for Field-Programmable Gate Arrays (FPGAs). Studies have proposed various imprecise functional techniques, albeit posed with three shortcomings: first, most existing inexact multipliers and dividers are specialized for Application-Specific Integrated Circuit (ASIC) platforms. Therefore, due to the architectural differences of underlying building blocks in FPGA and ASIC, ASIC-customized designs have not yielded comparable improvements when directly synthesized and ported to FPGAs. Second, state-of-the-art (SoA) approximate units are substituted, mostly in a single kernel of a multi-kernel application. Moreover, the *end-to-end* assessment is adopted on the Quality of Results (QoR), but not on the *overall* gained performance. Finally, existing imprecise components are not designed to support a pipelined approach, which could boost the operating frequency/throughput of, e.g., division-included applications. In this paper, we propose RAPID, the first pipelined approximate multiplier and divider architectures, customized for FPGAs. The proposed units efficiently utilize 6-input Look-up Tables (6-LUTs) and fast carry chains to implement Mitchell’s approximate algorithms. Our novel error-refinement scheme not only has negligible overhead over the baseline Mitchell’s approach, but also boosts its accuracy to 99.4% for arbitrary size of multiplication and division.

Experimental results obtained with Xilinx Vivado demonstrate the efficiency of the proposed pipelined and non-pipelined RAPID multipliers and dividers over accurate counterparts. In particular, 4-stage pipelined architecture of 32-bit RAPID multiplier (divider) enables  $3.3\times$  ( $5.1\times$ ) higher throughput,  $2.3\times$  ( $6.8\times$ ) higher throughput/Watt, and 52% (31%) savings of LUTs, over their 4-stage pipelined, accurate IP counterparts. Moreover, the end-to-end evaluations of non-pipelined RAPID, deployed in three multi-kernel applications in the domains of bio-signal processing, image processing, and moving object tracking for Unmanned Air Vehicles (UAV) indicate up to 35%, 33%, and 45% improvements in area, latency, and Area-Delay-Product (ADP), respectively, over accurate kernels, with negligible loss in QoR. To springboard future research in reconfigurable and approximate computing communities, our implementations will be available and open-sourced at <https://cfaed.tu-dresden.de/pd-downloads>.

**Index Terms**—Field-Programmable Gate Arrays, Approximate Computing, Pipeline, Multiplier, Divider, Mitchell’s Algorithm, Bio-signal Processing, Unmanned Air Vehicles, High-Throughput, Energy-Efficiency.

## I. INTRODUCTION

The ever-growing demand for edge computing has become pronounced in the Internet of Things (IoT) era for a wide domain of applications, from bio-signal to various cutting-edge image processing. For example, wearable 24/7 health monitoring gadgets are becoming ubiquitous, especially noting that 47% of cardiac diseases – the main cause of death, worldwide – occur outside of hospitals [1, 2]. Unmanned Aerial Vehicles (UAVs) such as drones are also proliferating for e.g., object/self tracking, search and surveillance, agricultural operations, and entertainment. Although Application-Specific Integrated Circuits (ASICs) are highly power-efficient for implementing the above-mentioned programs, off-the-shelf Field-Programmable Gate Arrays (FPGAs) have shown to serve as commercially-viable options owing to their rapid prototyping and post-fabrication datapath versatility which can address the outpacing speed of algorithmic updates over the updates in hardware [3, 4, 5]. For example, the hardware of health gadgets should be able to adapt with various patients’ physiological traits and the changes in heart’s activity. High-throughput and/or energy-efficiency are also of high desire for the acceleration of such parallelizable applications that are repeatedly fed with a bulk of data.

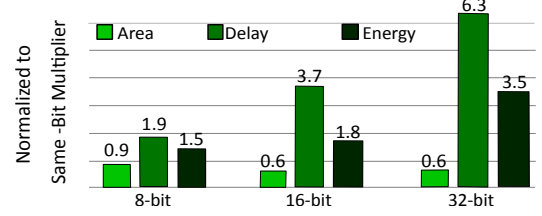


Fig. 1: Comparing area, delay, and energy of 8-, 16-, and 32-bit multiplication and division functions, implemented in Virtex-7 FPGA through Look-up Table (LUT) based accurate IPs.

FPGAs are equipped with hard-wired DSP blocks to speed-up multiplication, being the commonly-used function in bio-signal or image processing workloads. However, exploiting DSP blocks may not fulfill design requirements due to three reasons: first, their limited ratio versus Look-up Tables (LUTs), i.e.,  $<0.001$  could be insufficient for concurrent or multiplication-intensive applications. Second, their fixed locations in FPGAs impose routing overhead and may result in reduced performance for some applications [6, 7, 8]. Finally, DSPs are unable to be efficiently-utilized for multiplication with precision, smaller than  $18\times 18$  bit [9, 10]. Therefore, designers are compelled to also employ soft Intellectual Properties (IPs) for e.g., multiplication and division functions, proposed by major FPGA vendors such as Xilinx and Intel [11, 12]. In fact, exploitation of soft IPs instead of DSPs, for low bit-width operations, has also been suggested by academia and industry [13, 14]. Nonetheless, the long latency and high resource footprint of LUT-based IPs still needs to be decreased to facilitate the deployment of off-the-shelf FPGAs in aerial platforms and wearable gadgets. On the other hand, to boost the throughput in despite of stagnant clock speed, *pipelined* IPs have been proposed as a promising solution, albeit leaving the the quest for reducing the resource-gap unaddressed.

To minimize the resource footprint in aforementioned error-resilient programs, many approximation techniques have been emerged, however, three main challenges are attributed to them. First, approximation approaches customized in ASIC platforms have not yielded comparable performance gains when directly synthesized and ported to FPGAs, due to their different architectural specifications [7]. Second, approximation is often applied on a single kernel of a multi-kernel application, e.g., replacing multiplication (or division) with imprecise versions in the DCT (quantization) stage of JPEG compression. Their *performance gain* is also often evaluated and reported for that single kernel and not in the end-to-end implementation of the application. Third, while most efforts are concentrated on multiplication, studies like [15, 16] and our analysis shown in Fig. 1, consider the longer latency and higher energy of division compared to the multiplication. In fact, division is often the speed-bottleneck arithmetic operation in soft processors and can constrain the speed of application. Therefore, approximation of division has also become pronounced as, although less frequent, this operation is unavoidable in image processing and vision applications (Harris corner detection and K-means in unsupervised clustering) as well as bio-signal processing (heartbeat detection). In this context, a pipelined/low-latency divider is of high desire. Nevertheless, among

the imprecise dividers in the literature only one is targeted for LUT-based architectures [15] and none is specialized toward a pipelined design. In fact, pipelining has been mostly applied, in coarse-grain granularity, in a processor datapath while the fine-grained approach has been overlooked, especially in approximate designs.

To cope with the first challenge, we have designed *LeAp* [17] as the first logarithmic multiplier specialized for FPGAs. The motivation behind devising *LeAp* is based on translation of multiplication to addition in the logarithmic presentation, through Mitchell’s algorithm [18] ( $P = A \times B \xrightarrow[\text{Log}]{\text{Approx.}} \widetilde{\text{Log}}_P = \widetilde{\text{Log}}_A + \widetilde{\text{Log}}_B \xrightarrow[\text{Anti-Log}]{\text{Approx.}} \widetilde{P} = 2^{\widetilde{\text{Log}}_P}$ ). Transforming the 2D array structure of multiplication to 1D addition significantly reduces the design complexity. Specifically, it reduces the long latency and energy of an accurate divider, close to those of an accurate multiplier of the same-size. This translation is also suited for FPGAs and enables substantial gains as they already encompass fast carry chains to expedite addition. Implementation of Mitchell’s algorithm consists of three steps: approximate the log of inputs (by finding the position of the leading one), addition of two logs, and finally anti-log (a shift operation). *LeAp* has taken major leaps toward speeding-up Leading One Detection (LOD) and error-reduction steps: in the former, the LOD is accelerated by probing the 4-bit segment in parallel, and then finding the leading one in the most significant segment. For the latter, we have augmented the original Mitchell’s algorithm with three novel error-reduction schemes (independent from multiplier-size). The proposed error-reduction schemes allow the *addition* of the error coefficient *concurrently* with the fractional addition step and within the same resource footprint, used for implementing the baseline Mitchell’s designs [19]. Furthermore, *LeAp* [17] is customized in a way that effectively utilizes FPGA primitives, based upon the fact that LUTs and carry chain (having fixed latency), can also be configured to perform a ternary addition ( $\text{frac}_1 + \text{frac}_2 + \text{error coefficient}$ ). This is in contrast to previous approaches (MBM [20] and INZeD [16]), for which an additional circuitry was necessary for the addition of error-reduction terms to the original Mitchell’s circuits.

To surmount the rest of foregoing challenges in FPGAs, we further expand *LeAp* multiplier [17] and propose *RAPID*, which sets out to enable the first architectures for approximate *fine-grained pipelined* multiplier and divider. Our novel contributions in *RAPID* architectures are highlighted from three perspectives: first, following the light-weight error-reduction scheme of *LeAp*, we propose a logarithmic divider. Second, we further boost the accuracy of our designs by proposing various configurations for both multiplication and division that also surpass the SoA counterparts in terms of accuracy and/or performance metrics. Finally, we tailor our non-pipelined architectures for fine-grain pipelining and propose various versions of approximate multiplier and divider, having different number of stages. Such a spectrum of configurations also enables diverse power-throughput trade-offs, suited for different operating frequency levels within FPGAs. It is worth underlining that featuring fine-grained pipelining at an intra-unit granularity could enable operating at a higher frequency-level versus coarser-granularity (i.e., inter-unit), especially in applications including division operation. Note that although our circuits are specialized for FPGAs, our fine-grain pipelining approach is also applicable to ASIC-based architectures. In short, we make the following key technical research contributions in this journal version:

- **Near-zero biased multiplier & divider with extendable accuracy.** We further increase accuracy of *LeAp* multiplier and propose three versions of divider. The accuracy of error-reduction schemes are further increased so that the minimally-biased designs confine the average of absolute relative error (ARE) from  $>3.8\%$  to  $<0.6\%$ , through a limited number of coefficients.

- **First approximate pipelined multiplier and divider.** We design 2-, 3-, and 4-stage *pipeline* architectures that achieve up to  $8.1\times$  higher throughput and  $6.8\times$  higher throughput/Watt over accurate IP-based counterparts from Xilinx Vivado.
- **End-to-end performance-QoR evaluation on three application domains.** The efficacy of the proposed multiplier and divider is evaluated in the end-to-end implementation of three applications: Pan-Tompkins heartbeat QRS detection, JPEG compression, and Harris Corner Detection (HCD).
- **Open-source model.** The implementations of *RAPID* multipliers, dividers, and FPGA-customized programs will be open-sourced at <https://cfaed.tu-dresden.de/pd-downloads>, to fuel further research in reconfigurable and approximate computing communities.

The rest of this article is organized as follows: Section II presents a brief survey on the imprecise multipliers and dividers and distinguish the contribution of this work. Section III summarizes a background on Mitchell’s Mul/Div algorithms We elaborate upon the proposed architectures and pipelining approach in Section IV, respectively. Experimental setup, circuit- and application-level results (on three multi-kernel applications) are detailed in Section V. Finally, Section VI draws the conclusion with an outlook to interesting future tracks.

## II. RELATED WORK

Although substantial effort has been dedicated to ASIC-based inexact multipliers and dividers (a quantitative evaluation of them can be found in [21]), FPGA-specialized counterparts have also recently gained traction. The approaches in both landscapes are discussed herein and their compendium is highlighted in Table I.

**Partial product (PP) approximation:** a wide class of works targeting imprecise multipliers have mostly applied approximation on: 1) PP generation by simplification of truth table in e.g.,  $2\times 2$  and  $4\times 4$  multipliers and use them in a modular design [22]. 2) Reduction/accumulation of PP rows into two, using 3:2 and 4:2 imprecise compressors [23, 24, 25, 26, 14, 27, 28] in Dadda- or array-based multipliers. or 3) Adding the reduced PP rows by e.g., splitting the carry propagation path [7, 29]. The main drawback of these approaches is weak-scalability to larger input-width, as error may drastically increase, when accumulated in a hierarchical design approach. Furthermore, compressor-based designs are posed with high average relative error (ARE), or render moderate performance gain when applied on few least significant PP columns.

**Resizing to narrower-width multiplier/divider:** another category of works utilize a smaller instance of the arithmetic unit, based on a dynamic selection of most significant bits, starting from the position of the leading one ([35, 34, 36, 47] for multiplier and [37, 38, 48] for divider). Although offering high resource improvements, these truncation techniques suffer from error cases up to 100%. In addition, the latency of an accurate smaller-divider could be still multiple times of a same-sized multiplier.

**Division with inexact subtractor:** this branch of studies have truncated or replaced accurate subtractors with imprecise counterparts, in  $2N$ -by- $N$  non-restoring and restoring array dividers [49, 50, 51, 52]. The main difference between the two is in the remainder output, non-restoring version does not correct the remainder when subtraction has a negative result. Thus, a correction circuit is added to its hardware implementation. It has been shown that an approximate restoring version dissipates less power than its non-restoring counterpart, while introducing slightly larger accuracy degradation [52]. Such cells are also exploited in combinational implementation of higher radix SRT dividers. Generally approximate array dividers offer high accuracy, however, the resource savings achieved by them are small due to their array structure [37, 38]. Furthermore, their latency remains multiple times of a same-sized multiplier.

TABLE I: Summary of state-of-the-art approximation approaches for ASIC- and FPGA-based multipliers and dividers [30]

Approach	Mul/Div	Pipelined	ARE <sup>1</sup> up to (%)	Description	Platform	Reported Mul/Div Gain <sup>2</sup>	End-to-end <sup>4</sup>		
Partial Product Generation/ Addition/ Accumulation	✓/✗	✗	7.6	Inexact 4:2 compressor in Dadda Mul [24, 25, 26, 28, 27]	ASIC	{A, D, P} +	✗		
			1.7	Asymmetrically utilize inexact compressor in 3/4 of LSB columns [23]		{A, P} ++, D +			
			8.4	Simplified Karnaugh map of 4:2 compressor in Booth multiplier [31]		A +, E ++			
			Config.	Library of larger multiplier and adders using 2x2 instances [32, 22]	FPGA	{A, D, P} ++			
			0.3	Cutting the carry propagation path in 4-, 8-bit array multiplier [29, 7]		{A, D, P} +			
			8.5	Truth table simplification 3:2/4:2 compressor for Dadda multiplication [14]		{A, P} +			
Truncated Mul/Div	✓/✗	✗	10.9	Leading one based: with error compensation [34], with rounding [35, 36]	ASIC	{A, P} ++	✗		
	✗/✓		6.7	Leading-one position based 2k+2/k+1 Div plus error reduction circuit [37, 38]		{A, D} +, E ++			
	✓/✗		1.2-4.7	Variable-precision multiplier based on 8-bit truncated instances [37, 38]		{A, E} +			
Multiplicative Dividers	✗/✓	✗	2.9	Piecewise linear approximation and rounding of reciprocal of divisor [39]	ASIC	{D, E} ++	✗		
			6.4	Approximating reciprocal by bit manipulation [40], with truncation [40]		{A, D, E} +++			
			16.3	Approximating reciprocal using a table indexed by upper bits of divisor [41]		{A, D} +, E ++			
		✓	4.9	Incremental approximation of reciprocal of divisor using Taylor series [42]		{D, E} +++		✗	
Mitchell's Multiplication and Division Algorithms [18]	✓/✗	✗	2.9	Enhance Log accuracy: round rather truncation in piecewise approximation [43]	ASIC	{A, P} ++, D +	✗		
			> 3.9	Use different approximate adders in Mitchell's multiplier [44]		{A, P} +++			
			2.7	Improving accuracy of Mitchell's Mul with adding one error-correction [20]		{A, P} ++			
			2.7	Adding up to 256 error-coefficient to Mitchell's multiplier [45]		{A, P} +			
	✓/✗	3.0	Add one error-correction (with a similar approach to [20]) [16]	A +, {D, E} ++					
	✓/✗	✗	1-1.6	Adding one to five error-coefficient to Mitchell's multiplier [17]		FPGA		{A, E} ++, T +	✗
	✓/✓	✗	0.8	Adding 64 error-coefficients to Mitchell's multiplier and divider [15, 30, 46]		FPGA/ASIC		{A, E} ++, T +++	In [30, 46]
✓/✓	✓	0.6-1	Different pipelined designs for multiplier and divider with tunable accuracy	FPGA	{A, E} +++, T +++	✓			

<sup>1</sup> Average of Absolute Relative Error (a.k.a MRED), <sup>2</sup> Area/Delay/Energy/Power/Throughput <sup>4</sup> End-to-end performance gain

**Multiplicative dividers:** in another class of dividers, the reciprocal of the divisor is calculated first and subsequently multiplied with the dividend. Approximation is applied on the reciprocal of the divisor (using Taylor series [53] or linear piecewise approximation [39]). Encoding the reciprocal to lie in the range of (0.5, 1] is performed usually using a table indexed by the upper bits of the divisor (similar to [41]), or a series of bit manipulations, e.g., inverting all bits and appending '1' at MSB (like [40]). In some studies operands are also truncated/rounded to lie in some specific bit-width range. Overall, these approaches impose significant resource cost in FPGAs, due to requiring both reciprocal and multiplier IPs, separately. Moreover, when truncating divisor goes beyond relatively few bits, the accuracy of the reciprocated divisor degrades significantly [39].

**FPGA-customized multipliers:** it has been shown that ASIC-based approximation approaches have not yielded comparable performance gains when directly synthesized and ported to FPGAs [7, 33] (primarily, due to the underlying architectural differences). To cope with this challenge, recent works narrowed their focus toward specializing such techniques taking into account the LUT-based structure of FPGAs. For example, modification of LUTs which calculate LSBs in  $4 \times 4$  multiplications have been customized separately, for array based- [29, 7, 33], and Booth-based [54] architectures. Furthermore, an approximate compressor has been proposed in [14], geared toward an LUT-oriented implementation. Despite the specialized customization in these schemes for reconfigurable platforms, their resource savings have not been significant compared to logarithmic designs (discussed in the following). Moreover, Booth-based designs suffer from an increased critical path length [54]. This has highlighted the need for exploring other avenues for FPGAs which can achieve higher performance gain with an acceptable accuracy.

**Logarithmic multiplier and dividers:** In a more recent trend, some works have adopted Mitchell's approximate algorithms in lieu of array-based designs. Mitchell's algorithms translate multiplication (division) into log and addition (subtraction) in logarithmic representation. Such logarithmic transformations generally render higher resource gains compared to other approximation approaches. Interestingly, the latency of a division is also reduced, comparable to a same-sized multiplier, via this algorithm. However, these im-

provements come with the cost of relatively high error (ARE of 3.8% for multiplication and 4.1% for division). Therefore, various schemes have been presented to reduce the error in the original Mitchell's algorithms. The authors of MBM [20] have proposed a single error-reduction term for multiplication, and they employed a similar scheme for division in INZeD [16]. However, a single error-reduction term weakly fits all input combinations and eventuates in many output overflow cases, when adding the error-reduction term. In addition, their approaches are optimized for ASIC platforms. To alleviate error and also targeting FPGAs, we have recently proposed two FPGA-specialized designs, LeAp [17] and SIMDive [15] (a similar error-reduction approach to SIMDive has also been adopted in the REALM multiplier [45]). The approaches in these works considers  $F$  MSBs of fractional parts for each operand. Therefore, the possible combinations for the possible pairs makes a squarish region which is then partitioned to  $2^F \times 2^F$  sub-regions, each assigned with a unique error-coefficient. Although the error-reduction approaches of these works are similar in spirit, each has been designed for either ASIC or FPGA and targeted for different design goal. While REALM [45] is an ASIC multiplier, SIMDive [15] implements an FPGA-specific hybrid Mul/Div, aimed at applications with data parallelism opportunities that can benefit from Single Instruction, Multiple Data (SIMD) architectures. Moreover, as will be shown later, RAPID enables higher accuracy than SIMDive, even with smaller number of coefficients/LUTs.

**Pipelined architectures:** *Fine-grain* pipelining has recently gained attention. In this context, authors in [55, 56] have applied pipelining on PP generation and accumulation in the array-based *accurate* multipliers. Compared to accurate-, limited attention has been dedicated to approximate-designs. SoA divider SAADI-EC [42] and a modified version of the DRUM multiplier [57] have been the only works that adopted intra-unit pipelining. Despite the possibility of increasing throughput, these work lack a well-established pipelining strategy as no analysis has been performed to uniformly divide the critical path over pipeline stages, resulting in poor performance. To the best of our knowledge, no other work has proposed approximate pipelined multiplier or divider architectures, which is highly desired for processing of cutting-edge applications. In fact, the achieved

substantial boost in throughput is pivotal in e.g., image processing and health monitoring applications.

### III. PRELIMINARIES AND BACKGROUND

*Mitchell's Multiplication and Division Algorithms*: as shown in Eq. 1, Mitchell's algorithms perform imprecise multiplication and division in the logarithmic representation of numbers. Consider the binary representation for  $N$ -bit unsigned input  $A$ , which can be written as Eq. 2, where  $k$  indicate the position of the leading one. The rest of the bits (starting from position  $k-1$  to 0) are considered as the fractional part and fall in the range of  $0 \leq x < 1$ .

$$\begin{cases} P = A \times B \xrightarrow[\text{Log}]{\text{Approx.}} \widetilde{\text{Log}}_P = \widetilde{\text{Log}}_A + \widetilde{\text{Log}}_B \xrightarrow[\text{Anti-Log}]{\text{Approx.}} \tilde{P} = 2^{\widetilde{\text{Log}}_P} \\ D = A \div B \xrightarrow[\text{Log}]{\text{Approx.}} \widetilde{\text{Log}}_D = \widetilde{\text{Log}}_A - \widetilde{\text{Log}}_B \xrightarrow[\text{Anti-Log}]{\text{Approx.}} \tilde{D} = 2^{\widetilde{\text{Log}}_D} \end{cases} \quad (1)$$

$$A = 2^k + \sum_{i=0}^{k-1} 2^i b_i = 2^k(1+x) \xrightarrow{\text{e.g.}} 58 = 2^5(1+0.11010)_2, 18 = 2^4(1+0.001)_2 \quad (2)$$

In linear mathematics,  $\log_2(1+x)$  is approximated to  $x$  for this range of  $0 \leq x < 1$ ; thus, the approximate log value of input  $A$  is:

$$\text{Log}_2(A) \simeq k+x \rightarrow \text{Log}_2(58) \simeq (101.11010)_2, \text{Log}_2(18) \simeq (100.001)_2 \quad (3)$$

After applying the same step on the second input to get its approximate log, the summation (subtraction) of two parts is obtained in Eq. 4 (Eq. 5).

$$\widetilde{\text{Log}}_2(\tilde{P}) = (k_1 + k_2) + (x_1 + x_2) \rightarrow K_s = (1001)_2, X_s = (0.1111)_2 \quad (4)$$

$$\widetilde{\text{Log}}_2(\tilde{D}) = (k_1 - k_2) + (x_1 - x_2) \rightarrow K_s = (1)_2, X_s = (0.1011)_2 \quad (5)$$

Finally, by applying the anti-log (which mathematically is a shift operation), binary representation of the approximate product (quotient) are derived by Eq. 6 (Eq. 7):

$$\tilde{P} = \begin{cases} 2^{k_1+k_2}(1+x_1+x_2), & x_1+x_2 < 1 \\ 2^{k_1+k_2+1}(x_1+x_2), & x_1+x_2 \geq 1 \end{cases} \rightarrow \tilde{P} = 992, P_{acc} = 1044 \quad (6)$$

$$\tilde{D} = \begin{cases} 2^{k_1-k_2-1}(2+x_1-x_2), & x_1-x_2 < 0 \\ 2^{k_1-k_2}(1+x_1-x_2), & x_1-x_2 \geq 0 \end{cases} \rightarrow \tilde{D} = (11)_2 = 3, D_{acc} = 3 \quad (7)$$

### IV. PROPOSED PIPELINED AND NON-PIPELINED MULTIPLIER AND DIVIDER ARCHITECTURES

In this section we initially present the proposed error-reduction schemes, and afterwards elaborate on the structure of RAPID multiplier and divider architectures (non-pipelined and pipelined).

#### A. Proposed light-weight & minimally-biased error-reduction schemes

Mitchell's error for 8-bit multiplication and division are formulated in the Equation 8 and 9. Inspecting the behavior of the error has provided the following insights, some of which are also noted and endorsed by SoAs [45, 15, 20, 16]:

$$E_P = P - \tilde{P} = \begin{cases} 2^{k_1+k_2}(x_1x_2), & x_1+x_2 < 1 \\ 2^{k_1+k_2}(1-x_1-x_2+x_1x_2), & x_1+x_2 \geq 1 \end{cases} \quad (8)$$

$$E_D = D - \tilde{D} = \begin{cases} 2^{k_1-k_2} \frac{(x_1(x_2-1)+x_2-(x_2)^2)}{2(1+x_2)}, & x_1-x_2 < 0 \\ 2^{k_1-k_2} \frac{(x_1x_2-(x_2)^2)}{1+x_2}, & x_1-x_2 \geq 0 \end{cases} \quad (9)$$

- The equations demonstrate the different error magnitude in each power-of-two interval. Therefore, merely adding a single correction term to the output (as proposed in INZeD [16] and MBM [20]) is not efficient and results in many output overflow cases [45, 15].

- The equations also reflect the proportional replication of error in each power-of-two. This behavior is repeated for each size of multiplication or division (irrespective of  $k_1$  and  $k_2$ ). This repetition allows applying the same error-reduction approach for all multiplier or divider sizes. In fact, a set of coefficients can be added to fractional parts, before that their summation become scaled [20, 16].

- As can be seen in Fig. 2 (a), and (e), some of nearby sub-regions have very small error (less than 2%), while the changes of error occurs is more pronounced in other sub-regions. This motivates grouping the regions having similar error as well as assigning more coefficients to the sub-regions having higher error (to efficiently reduce the error to a certain bound).

It is also worth noting that an efficient partitioning should result in a small number of coefficients to minimize the overhead of error-reduction. Partitioning approaches such as REALM [45] and SIMDive [15] results in superfluous number of error-coefficients when targeting high accuracy. Note, considering even four MSBs of fractional parts results in  $2^4 \times 2^4 = 256$  coefficients and limits the scalability of such schemes (especially for FPGAs, as will be discussed later). It fact, the resource-efficiency of the error-reduction strategy proposed in SIMDive depends on the number of LUT inputs and is not easily extendable when targeting higher accuracy. Going into the details, to generate 256 coefficient through SIMDive approach, four MSBs from each fractional part should be considered, which requires the direct utilization of 8-LUTs for each bit of error-reduction term. Otherwise the usage of a 256-to-1 Mux (each input is 8- or 16-bit) is inevitable. Implementing this mux or 8-LUT out of 6-LUTs increases the error-reduction overhead of SIMDive to nearly  $4 \times$  of its original idea (well-suited for considering 3 or less MSBs of fractional part). This overhead is not negligible, especially when targeting a small-sized multiplier or divider. Building upon this discussion, we propose RAPID to enable reaching higher accuracy with even smaller overhead.

**Proposed light-weight error-reduction scheme**: to cope with the overflow problem in INZeD and MBM, and the over-provisioned number of coefficients (i.e., 256) in REALM, the insights obtained from aforementioned observations incentivize reducing the error-reduction terms by efficiently clustering them in fewer groups. In RAPID, we divide the squarish region between each power-of-two pair to *different, and fewer*, regions than REALM/SIMDive. In our partitioning we consider the following factors:

- 1) Consider four rather than three MSBs of fractional parts, to increase the accuracy.
- 2) Minimize the number of sub-regions while considering four MSBs, to constrain resource cost for selection of the coefficient.
- 3) Minimize *error distribution*  $\times$  *error magnitude* in each region (can be estimated to the integral of error-magnitude for that region).

We have proposed three error-reduction schemes in Fig. 2, based on the in-depth analysis of error and the resource-usage reports from Vivado. To minimize the average of absolute relative error (ARE), various partitionings have been investigated. The goal of partitioning is to keep the error of grouped sub-regions nearly uniform and below pre-defined thresholds, e.g.,  $\sim 4\%$ ,  $3\%$ , and  $2.5\%$  for 3-, 5-, and 10-coefficient multiplier schemes, respectively. Further, the error-reduction coefficient for each group of sub-intervals, are derived by following the mathematical approach detailed in [45]. The binary representation of the proposed coefficients for multiplier and divider are also shown in Table II. From the implementation point of view, partitioning is implemented via small-sized multiplexers (coded via casex statement in HDL). In such implementation, the LUT usage depends on the *number of error-coefficients as the Mux-inputs* and the *complexity of conditional statements in the casex statement*. In order to minimize the former (number of Mux-inputs) we have proposed three schemes, having most 10 coefficients. In order to minimize the latter (complexity of conditional statements for selecting the coefficients), we only consider comparing 4 MSBs of fractional parts for the partitioning. Moreover, in this partitioning, neighbouring sub-regions having the same coefficient are packed to reduce the complexity for each conditional statement. From the hardware-usage perspective, each 6-LUT functions as a 4-MUX (in which 4 inputs and the 2 select lines are the LUT inputs). Further, as also denoted by [58], a 16:1 multiplexer requires a single FPGA-slice having four 6-LUTs. Based upon this, we have also chosen a squarish partitioning through a MUX-based approach, to minimize the overhead for the proposed error-reduction strategy. This is while, checking the rhombus borders (see the original error shape in Fig. 2) is very costly in hardware and could nullify the LUT-saving, gained from approximating approach. It should be noted that the proposed partitioning is also scalable, as the resource cost for choosing one of few coefficients does not grow exponentially, contrary to the REALM/SIMDive approaches. In fact, the proposed partitioning surpasses the SoAs in terms of resource-error trade-off: not only with 10 error-coefficients and considering 4 MSBs we achieve Average Relative Error (ARE) of 0.6%, which is better than SIMDive/REALM (0.8% ARE, considering 3 MSBs), but also the resource footprint of this scheme is 193 LUTs, still lower than SIMDive/REALM counterparts (see Table III).

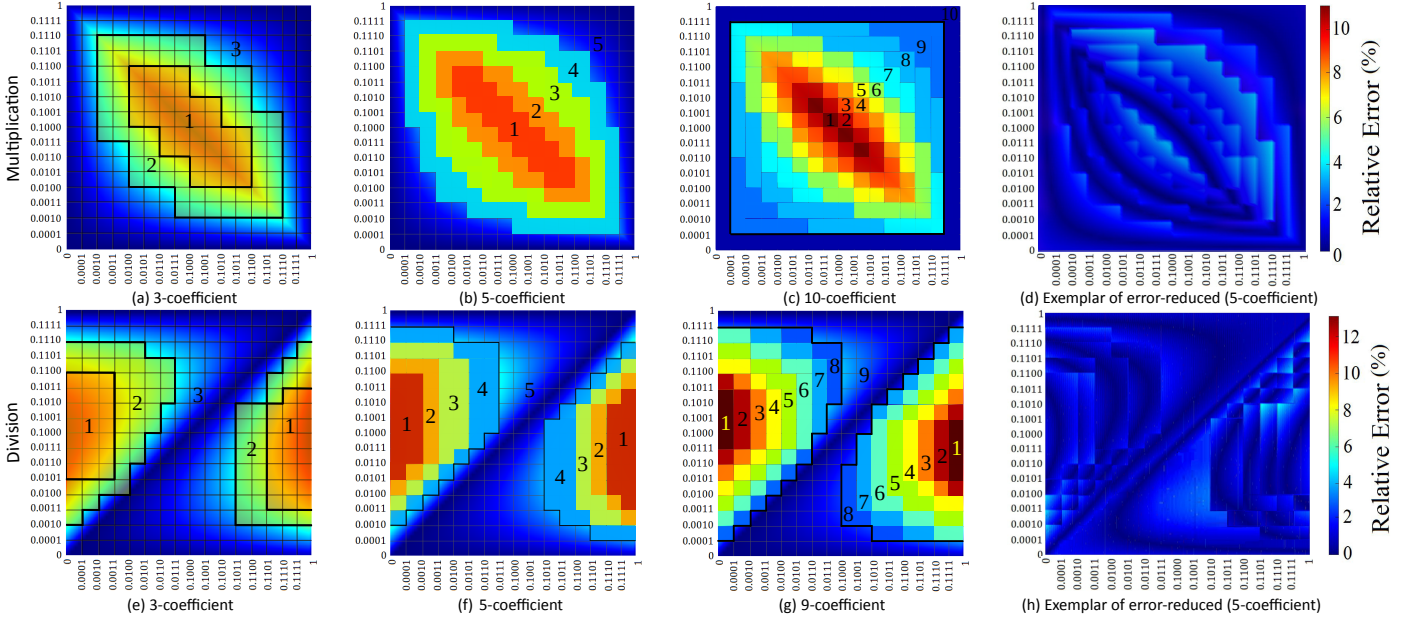


Fig. 2: Proposed error reduction schemes of RAPID for multiplication and division, based on 4 MSBs of fractional parts.

TABLE II: Binary representation of error-reduction coefficients in 16-bit multiplier & divider (3/4 MSBs are zero for Mul/Div and excluded)

Multiplier				Divider			
3-coefficient	5-coefficient	10-coefficient		3-coefficient	5-coefficient	9-coefficient	
1) 10000100111	1) 1001111111111	1) 1001111000110	6) 0101110011111	1) 1000011111111	1) 1001111000100	1) 1001110001111	6) 0110010100101
2) 010011101100	2) 1000011011101	2) 1000110110001	7) 0100101000011	2) 0100010111111	2) 1000010001111	2) 1000110111100	7) 0101000101011
3) 000100101001	3) 0110010001010	3) 0111111000100	8) 0100001011101	3) 0001011111111	3) 0110110001101	3) 100000010100	8) 0100111101000
	4) 0011110010111	4) 0111000110101	9) 0011110000011		4) 0101010100111	4) 0111001100010	9) 0100001101100
	5) 0000111110000	5) 0110010100011	10) 0010101111111		5) 0011011100100	5) 0110100001101	

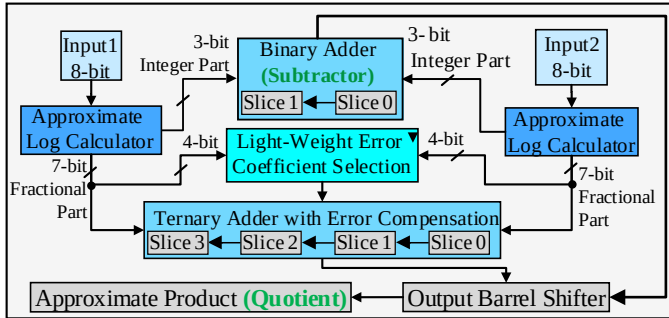


Fig. 3: Overall structure of proposed RAPID multiplier and divider

### B. Structure of RAPID multiplier and divider (Non-Pipelined)

Fig. 3 illustrates the structure of the proposed multiplier and divider. The overall design is designed based on  $Res = A \times B \xrightarrow{\text{Approx. Log}} \widetilde{LogRes} = \widetilde{LogA} \pm \widetilde{LogB} \xrightarrow{\text{Approx. Anti-Log}} \widetilde{Res} = 2^{\widetilde{LogRes}}$ .

**Leading-one detection:** To accelerate leading one detection in our FPGA-customized method, this process is calculated based upon 4-bit LODs (implemented by directly configuring LUTs). In the first step, we probe the presence of bit value ‘1’, simultaneously, in each 4-bit segment of the operands. To this end, one LUT is configured as a logical OR function, applied on 4-bits of each group to reveal whether the segment contains a bit with value ‘1’ (acts as a zero-detection flag). In parallel, another 6-LUT is configured to two 5-LUTs in such a way that it determines the position of leading-one in the 4-bit segment (LOD4-LUT). Finally, based on the resulting bits from the output of these 6-LUTs, we determine the position of leading one in the most significant group, through the priority logic.

For example, the position of leading one in the 8-bit LOD equals to the concatenation of  $\{Location\ index\ of\ most\ significant\ segment, Leading\ one\ position\ in\ that\ segment\}$ , e.g., the leading one position in “01010101” is  $\{\{1\}, \{10\}\} \Rightarrow \{110\}$  in binary. The first part ( $\{1\}$ ) is the output of Flag-LUT that has been employed for upper 4-bit segment. The second part ( $\{10\}$ ) is the result of LOD4-LUTs on the upper segment. A similar method has been also exploited for 16- and 32-bit LODs. For example, in a 16-LOD, if the upper half of the operand is zero, the 16-bit LOD is equal to the lower 8-bit LOD. Else, the position of leading-one is  $8 + \text{leading-one position in the upper-half 8-LOD}$ . This can be obtained by applying a logical-OR function on the outputs of Flag-LUTs on the 4-bit segments of the upper-half. In LeAp [17], LOD step was orchestrated through an FSM and performed in, at most, five clock-cycles. In order to achieve efficient pipelining and minimize the number of registers, the LOD is implemented as combinational logic. Subsequently, the critical path is analyzed in order to achieve balanced partitioning for pipelining.

**Addition of integer parts:** as shown in Fig. 3, each 4-bit addition is fulfilled by one Virtex-7 slice which includes four 6-LUTs and its associated fast carry chains. Together, these implement a Carry Look-Ahead Adder (CLA). Extending the 4-bit addition to 8-bits is easily achievable by connecting the  $C_{out}$  from a previous slice to the  $C_{in}$  of the next slice. Recalling Eq. 4 and Eq. 5, division is performed by changing additions to subtractions, through 2’s complement modules.

**LUT-optimised ternary addition:** recalling that in LeAp [17], we have proposed error reduction coefficients such that it only depends on the fractional bits (and not the intermediate result of Mitchell Mul/Div, contrary to MBM/INZeD [20, 16]). On the other hand, LUTs and their associated fast carry chain in Xilinx UNISIM library [59] can be configured to implement a ternary adder. To

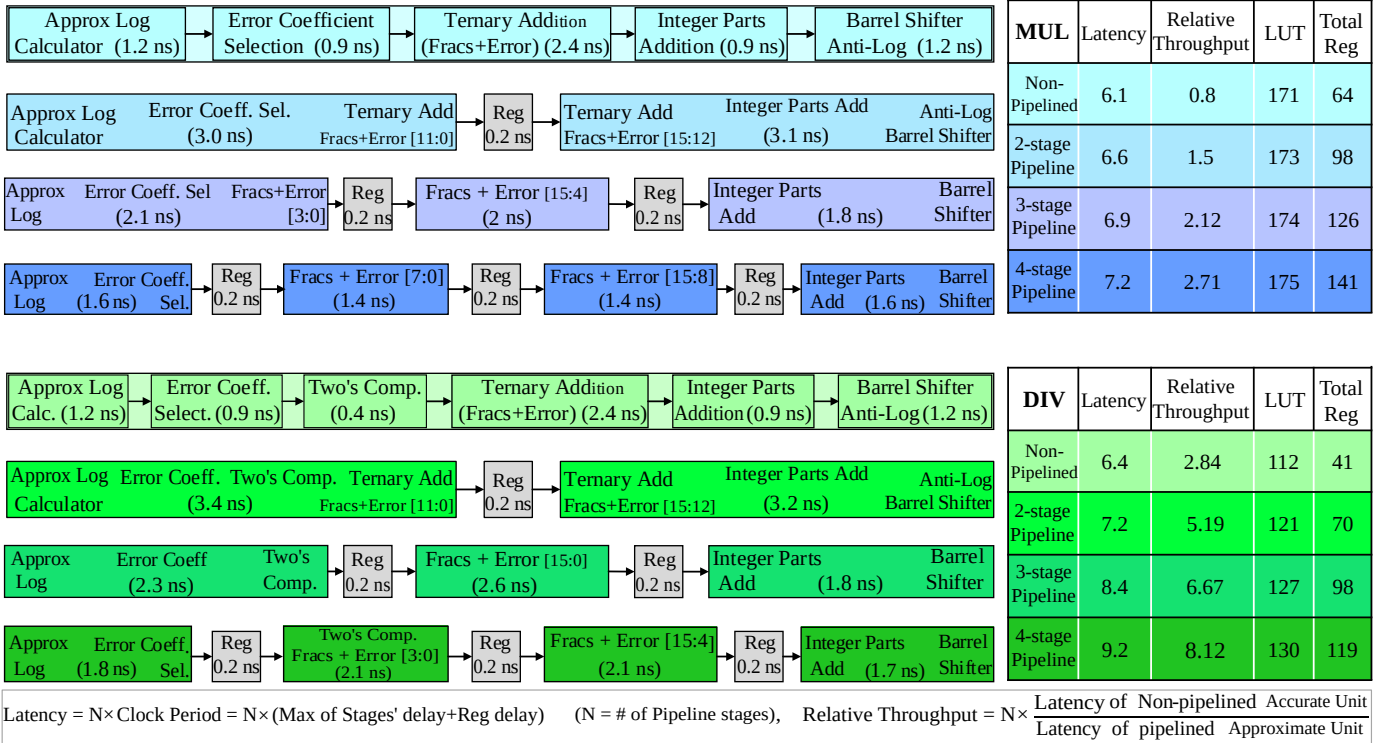


Fig. 4: Proposed 2-, 3-, and 4-stage pipelining for 16x16 RAPID-5 multiplier (top) and RAPID-9 16/8 divider (bottom).

this end, we have manually configured the LUTs and carry chain primitives of the FPGA in such a way that they implement a ternary adder. This highly suits our error-reduction approach as *enables combining the process of adding error-reduction coefficient with fractional parts with the same resource footprint and in a single step*. In fact, regardless of the ternary adder size and compared to the binary version, only one more bit at MSB position is needed, as  $\text{frac}_1 + \text{frac}_2 + \text{error\_coefficient}_i + C_{in}$  ( $C_{out}$  from prior bit) may result in 3 bits, requiring another LUT at the end of the chain [19]. Moreover, the delay of FPGA primitives is fixed. Therefore, adding error-reduction term at the same time as fractional parts does not impose additional overhead to the design. On the contrary, in REALM [45], MBM [20], and INZeD [16] an additional circuit is needed to add the error-reduction term (or half of it) to the original Mitchell's circuits, based on the intermediate summation/subtraction of fractional parts.

It should be noted that to prevent overflow in  $2N$ -by- $N$  bit standard division, the condition for  $\text{dividend} < 2^N \times \text{divisor}$  needs to be satisfied [60]. This means that scaling of the divider lies in the range of  $2^0$  to  $2^{N-1}$  (similar to [16]). Therefore, only  $N-1$  bits are used from the subtractors for the output and  $N$  LSBs from  $\log_{\text{dividend}}$  is neglected. This also reduces the resource footprint for subtractor and barrel shifter and does not affect the accuracy.

### C. Proposed Pipelined Architectures for Multiplier and Divider

In order to minimize the latency overhead, the combinational datapath of multiplier and divider should be partitioned for near uniform latency over the pipelined stages. To achieve this, we have adopted the following steps: first, each stage of the multiplication or division is synthesized in isolation to get an *estimation* of the delay for each stage. Fig. 4 shows the latency of individual stages in non-pipelined and various pipelined configurations for a 5-coefficient multiplier and divider (the results for all designs are shown in Table III). Note that after integrating the sub-modules and synthesizing the entire component, the end-to-end latency may change due to

the default-optimizations and structural-modification applied by the tool. Based on this analysis, the pipelining registers were inserted in the proper locations of the non-pipelined design. Afterwards, re-synthesis has been performed to assess whether a marginal fine-tuning for the adopted partitioning can result in better end-to-end latency. It can be inferred from Fig. 4, that through the adopted pipelining approach, the latency of stages become similar. Moreover, the delay overhead of a pipeline register is small, particularly since the divider is in the critical path in most applications. Moreover, each of the proposed soft-core designs has a different operating frequency and can be utilized w.r.t the different frequency levels offered by a single or various FPGA families. Finally, different pipeline versions enable a spectrum of latency-throughput trade-offs. Therefore, the proper configurations can be loaded w.r.t. the application's requirement.

## V. RESULTS AND DISCUSSION

In this section, we first present the implementation results of the proposed RAPID multipliers and dividers (pipelined and non-pipelined configurations). We have also implemented the compared the proposed RAPID against the following designs: DSP- and pipelined/non-pipelined IP-based accurate counterparts, dynamically truncated DRUM [47] and AAXD [37], hierarchical-based AFM multiplier [29], SAADI-EC pipelined divider [42], and Mitchell-based counterparts (SIMDive multiplier and divider [15], MBM multiplier [20], and INZeD divider [16]). We compare the designs w.r.t different performance metrics including *Performance (herein throughput) per Watt* which is considered by industry to be the very new Moore's Law [61]. Afterwards, the end-to-end performance gain on three multi-kernel applications is assessed.

### A. RAPID versus SoA multipliers & dividers (circuit-level evaluation)

**Experimental Setup:** to evaluate performance metrics, all the designs are developed in Verilog HDL, synthesized, and implemented in Xilinx Vivado 2019.2 for the Virtex-7 FPGA. To ensure scalability

TABLE III: Accuracy-resource trade-off of accurate/approximate multipliers &amp; dividers (pipelined/non-pipelined) in FPGA implementation

Multiplier													Divider										
8×8 Mul	LUT	FF	E2E <sup>1</sup> Latency (ns)	Rel. <sup>2</sup> Tput	Circuit Power (mW)	Clk Power (mW)	Rel. <sup>3</sup> Energy per inst	Rel. Tput/ Watt	ARE (%)	PRE (%)	Error Bias (%)	8/4 Div	LUT	FF	E2E Latency (ns)	Rel. Tput	Circuit Power (mW)	Clk Power (mW)	Rel. Energy per inst	Rel. Tput/ Watt	ARE (%)	PRE (%)	Error Bias (%)
DSP-based	7 + 1 DSP	31	3.72	1.97	17.52	10.3	1.77	1.13	-	-	-	DSP-based	11+ 2 DSPs	63	7.48	1.44	19.80	11.6	2.25	0.44	-	-	-
Acc IP <sub>NP</sub> <sup>4</sup>	60	48	3.67	1.0	15.91	-	1.0	1.0	-	-	-	Acc IP <sub>NP</sub>	51	42	10.74	1.0	9.62	-	1.0	1.0	-	-	-
Acc IP <sub>P2</sub>	60	88	4.30	1.71	23.19	37.5	2.23	0.45	-	-	-	Acc IP <sub>P2</sub>	55	59	13	1.67	21.09	18.4	2.43	0.41	-	-	-
Acc IP <sub>P3</sub>	60	112	5.10	2.16	47.5	43	2.64	0.38	-	-	-	Acc IP <sub>P4</sub>	69	89	23.7	1.81	26.27	29.1	3.14	0.32	-	-	-
RAPID-3 <sub>NP</sub>	57	39	4.97	0.74	12.06	-	1.03	0.97	1.02	6.1	0.06	RAPID-3 <sub>NP</sub>	41	20	5.20	2.06	9.72	-	0.48	2.06	0.99	5.74	0.02
RAPID-5 <sub>P2</sub> <sup>5</sup>	62	56	5.45	1.35	18.17	22.9	1.91	0.52	0.91	4.45	0.05	RAPID-5 <sub>P2</sub>	44	20	5.18	4.15	14.74	20.8	0.88	1.13	0.79	4.34	0.01
RAPID-10 <sub>P3</sub>	71	69	6.88	1.60	24.31	34.0	2.29	0.44	0.64	3.69	0.05	RAPID-9 <sub>P3</sub>	51	20	5.34	8.05	22.66	30.6	0.91	1.47	0.58	3.48	0.01
AFM1	69	58	5.50	0.67	18.66	-	1.76	0.57	0.23	16.52	0.23	SIMDiv-DIV	44	20	5.23	2.09	9.23	-	0.45	2.20	0.77	5.20	0.01
SIMDiv-MUL	61	32	5.13	0.72	16.34	-	1.44	0.7	0.82	4.76	0.05	INZed	47	21	6.12	1.75	11.14	-	0.65	1.53	2.93	9.53	0.02
MBM	64	33	5.16	0.71	17.48	-	1.54	0.65	2.60	8.59	0.09	Mitchell	36	20	5.10	2.11	8.53	-	0.42	2.40	3.90	13.00	3.90
Mitchell	51	32	4.82	0.76	13.39	-	1.11	0.9	3.77	11.11	3.77	SAADI-EC (16)	103	42	13	0.84	24.42	-	2.99	0.33	2.37	8.82	1.92
DRUM-4	53	32	5.08	0.72	14.33	-	1.25	0.81	5.82	25.35	1.84	AAXD (6/3)	38	20	6.06	1.77	9.01	-	0.52	1.91	2.08	100.00	1.49
16×16 Mul	LUT	FF	E2E Latency (ns)	Rel. Tput	Circuit Power (mW)	Clk Power (mW)	Rel. Energy per inst	Rel. Tput/ Watt	ARE (%)	PRE (%)	Error Bias (%)	16/8 Div	LUT	FF	E2E Latency (ns)	Rel. Tput	Circuit Power (mW)	Clk Power (mW)	Rel. Energy per inst	Rel. Tput/ Watt	ARE (%)	PRE (%)	Error Bias (%)
DSP-based	8 + 1 DSP	32	4.11	1.19	17.48	9.8	0.48	2.08	-	-	-	DSP-based	197+ 7DSP	131	10.36	1.76	51.62	25.1	2.34	0.43	-	-	-
Acc IP <sub>NP</sub>	287	64	4.88	1.0	47.81	-	1.0	1.0	-	-	-	Acc IP <sub>NP</sub>	169	76	18.23	1.0	17.97	-	1.0	1.0	-	-	-
Acc IP <sub>P2</sub>	249	176	6.14	1.60	64.86	65.3	1.71	0.59	-	-	-	Acc IP <sub>P2</sub>	175	104	19.59	1.86	38.82	12.3	1.47	0.68	-	-	-
Acc IP <sub>P3</sub>	249	245	8.88	1.65	94.49	75.0	2.15	0.46	-	-	-	Acc IP <sub>P4</sub>	181	168	20.09	3.63	56.21	24.5	1.68	0.6	-	-	-
Acc IP <sub>P4</sub>	249	343	9.60	2.03	150.73	130.0	2.89	0.35	-	-	-	RAPID-3 <sub>NP</sub>	112	41	6.38	2.98	18.67	-	0.34	2.98	1.02	5.74	0.02
RAPID-3 <sub>NP</sub>	168	64	5.90	0.83	31.43	-	0.86	1.17	1.03	6.1	0.06	RAPID-5 <sub>P2</sub>	121	70	7.07	5.16	27.77	19.5	0.49	2.04	0.79	4.34	0.01
RAPID-3 <sub>P2</sub>	169	98	6.11	1.52	47.35	40.7	1.21	0.85	1.03	6.1	0.06	RAPID-9 <sub>P3</sub>	127	98	8.35	6.62	30.29	24.3	0.44	2.28	0.58	3.48	0.01
RAPID-5 <sub>P3</sub>	177	126	6.87	2.25	75.35	58.2	1.27	0.81	0.93	4.45	0.05	RAPID-9 <sub>P4</sub>	130	119	9.2	8.01	34.68	27.7	0.42	2.40	0.58	3.48	0.01
RAPID-10 <sub>P4</sub>	193	141	7.25	2.52	84.75	87.4	1.46	0.7	0.56	3.69	0.23	SIMDiv-DIV	143	64	5.68	3.28	23.84	-	0.39	2.57	0.78	5.20	0.01
AFM1	261	66	7.32	0.67	44.78	-	1.41	0.71	1.34	17.80	1.34	INZed	165	41	6.28	2.90	27.50	-	0.51	1.97	2.93	9.54	0.02
SIMDiv-MUL	216	64	5.95	0.82	37.06	-	0.95	1.06	0.82	4.90	0.05	Mitchell	106	64	5.56	3.39	17.34	-	0.32	3.11	4.11	13.00	4.11
MBM	204	65	6.59	0.74	35.34	-	1.0	1.0	2.63	8.83	0.09	SAADI-EC (16)	342	126	25.70	0.71	57.01	-	4.31	2.23	2.14	8.82	1.76
Mitchell	167	64	5.51	0.99	31.46	-	0.64	1.56	3.85	11.11	3.85	AAXD (8/4)	151	155	12.51	1.46	25.17	-	0.93	1.08	2.99	100	0.90
DRUM-6	233	64	5.34	0.91	38.43	-	0.88	1.14	1.47	6.31	0.04	AAXD (12/6)	207	233	21.26	0.86	34.51	-	2.16	0.46	0.74	100	0.30
32×32 Mul	LUT	FF	E2E Latency (ns)	Rel. Tput	Circuit Power (mW)	Clk Power (mW)	Rel. Energy per inst	Rel. Tput/ Watt	ARE (%)	PRE (%)	Error Bias (%)	32/16 Div	LUT	FF	E2E Latency (ns)	Rel. Tput	Circuit Power (mW)	Clk Power (mW)	Rel. Energy per inst	Rel. Tput/ Watt	ARE (%)	PRE (%)	Error Bias (%)
DSP-based	18 + 4 DSPs	75	6.93	1.04	35.48	14.8	0.47	2.28	-	-	-	DSP-based	359+ 9DSP	228	11.65	3.62	92.91	28.7	0.98	1.02	-	-	-
Acc IP <sub>NP</sub>	1012	128	6.69	1.0	110.56	-	1.0	1.0	-	-	-	Acc IP <sub>NP</sub>	597	139	42.24	1.0	34.16	-	1.0	1.0	-	-	-
Acc IP <sub>P2</sub>	1012	291	7.82	1.71	154.26	92.8	1.31	0.76	-	-	-	Acc IP <sub>P2</sub>	607	208	47.48	1.78	74.55	10.4	1.4	0.74	-	-	-
Acc IP <sub>P3</sub>	1012	835	9.74	2.06	218.32	146.6	1.6	0.63	-	-	-	Acc IP <sub>P4</sub>	607	339	56.95	2.97	111.04	21.9	1.75	0.76	-	-	-
Acc IP <sub>P4</sub>	1014	1119	11.73	2.28	330.64	250.8	2.31	0.43	-	-	-	RAPID-3 <sub>NP</sub>	378	85	6.70	6.30	42.70	-	0.19	5.04	1.04	5.74	0.02
RAPID-3 <sub>NP</sub>	434	128	6.30	1.06	90.08	-	0.77	1.3	1.05	6.1	0.07	RAPID-5 <sub>P2</sub>	389	128	8.43	10.02	49.07	16.5	0.2	5.22	0.79	4.34	0.01
RAPID-3 <sub>P2</sub>	450	184	6.74	1.98	131.22	62.7	0.88	1.13	1.05	6.1	0.07	RAPID-9 <sub>P3</sub>	399	173	9.86	12	59.20	20.2	0.2	5.24	0.61	3.48	0.01
RAPID-5 <sub>P3</sub>	462	242	7.38	2.72	181.53	95.1	0.92	1.09	0.95	4.45	0.05	RAPID-9 <sub>P4</sub>	417	213	11.10	15.22	67.59	26.7	0.2	5.53	0.61	3.48	0.01
RAPID-10 <sub>P4</sub>	490	276	8.04	3.33	253.12	124.2	1.03	0.97	0.58	3.64	0.06	SIMDiv-DIV	381	80	6.84	6.18	43.04	-	0.2	4.90	0.81	5.16	0.02
AFM1	995	317	10.76	0.62	119.63	-	1.74	0.57	2.88	22.40	2.88	INZed	422	81	8.15	5.18	47.67	-	0.27	3.71	2.96	9.47	0.03
SIMDiv-MUL	521	128	6.88	0.97	83.50	-	0.78	1.29	0.91	4.72	0.05	Mitchell	349	80	6.21	6.91	39.42	-	0.17	5.99	4.19	13.00	4.19
MBM	533	129	7.51	0.89	89.94	-	0.91	1.09	2.69	8.74	0.10	SAADI-EC (16)	822	228	51.60	0.82	92.86	-	3.32	0.30	2.33	9.04	1.85
Mitchell	428	128	6.23	1.12	60.25	-	0.55	1.83	3.91	11.11	3.91	AAXD (8/4)	361	278	24.66	1.71	40.78	-	0.70	1.43	3.04	100	1.10
DRUM-6	616	128	6.35	1.05	88.85	-	0.76	1.31	1.53	5.88	0.05	AAXD (12/6)	513	505	37.20	1.14	57.95	-	1.49	0.67	0.79	100	0.35

<sup>1</sup>End-to-end <sup>2</sup>Throughput <sup>3</sup>Energy per instruction = total dynamic power  $\times$  clock period <sup>4</sup>Non-Pipelined <sup>5</sup>5-Coeff. 2-stage pipelined

of multipliers and dividers, they are compared for precisions of 8-, 16-, and 32-bit. Area and latency are collected from Vivado reports. Power and energy dissipations<sup>1</sup> are obtained through simulations in Xilinx Power Estimator (XPE) over 100 million inputs, uniformly distributed in a random order over the whole input interval. For a precise measurement, we used the clock gating command to prevent superfluous switching activity in unused resources. To assess accuracy metrics for different partitioning, the behavioral structure of multipliers and dividers are developed in C++. To calculate the average of absolute relative error (ARE), peak absolute relative error (PRE), and error bias in 8- and 16-bit designs, exhaustive testing is performed. For 32-bit error characterization,  $2^{32}$ ,  $\sim 4.3$  billion input pairs, uniformly distributed over the whole 32-bit interval have been

evaluated in Monte Carlo simulations<sup>2</sup>. Post-implementation results are summarized in Table III. The following conclusions can be made based on the results:

- **RAPID multipliers and dividers versus DSP-based and accurate Vivado IPs:** the results corroborate prior studies [9, 10] stating that DSPs are able to be efficiently-utilized, only for large bit-width precision. The reason behind this degraded performance is that DSP48E1 hosts a  $25 \times 18$  hard-wired multiplier and is not optimized for smaller multipliers. This has been the main motive that FPGA designers have recommended utilizing of soft multipliers should be used for implementing lower bit-widths [9, 10, 13, 12]. In particular, compared to LUT-based implementation, DSP-based dividers are less energy-efficient, for both 8- and 16-bit precision. Targeting higher-order precision of 32-bit, the proposed RAPID multiplier and divider have lower

<sup>1</sup>Similar to prior works, only dynamic power is reported from Xilinx Vivado, as the static power analysis of Vivado is for the entire FPGA [6, 54].

<sup>2</sup>Simulated on Rack Server: Intel Xeon E5-2667 @3.2 GHz, 512GB RAM

latency than DSP-based counterparts. In terms of energy-efficiency (relative energy/instruction), DSPs are better than RAPID, only for multiplication. In contrast, RAPID dividers are significantly better than DSP-based implementation (especially in 32-bit precision). Compared to LUT-based IPs, the results exhibit the efficiency of RAPID divider in all performance metrics, when targeting division of any precision. In case of multiplication, although accurate IP renders better performance in 8-bit over RAPID, the area-and/or energy-saving of RAPID over accurate Vivado IP becomes substantial in 16-bit. Moreover, the latency of RAPID also becomes smaller when targeting 32-bit. Overall, the improvements become more pronounced, when architectures are implemented for higher bit-width. This is due to: first, transforming the 2D array-based structure of Mul/Div to 1D Add/Sub through Mitchell's algorithms, as discussed earlier. Second, the cost of error coefficients gets amortized in higher-order designs.

- **RAPID multipliers and dividers render better performance over hierarchical structures:** comparing Mitchell-based RAPID with modular counterpart AFM (structured by incorporating smaller inexact instances) demonstrates three points: first, thanks to transforming the 2D array structure of multiplication to 1D addition in the logarithmic representation, the area of Mitchell-based unit grows by the factor of  $\sim 2.6$ , less aggressive compared to  $\sim 3.7$  for the array-based counterparts. This further accentuates the efficiency of RAPID, particularly in larger bit-width. Second, comparing the results of 8- versus 16-bit asserts that approximation applied on hierarchical approaches is beneficial in accuracy-resource trade-off, only when it is done from scratch for each multiplier size, otherwise accumulated error in larger designs significantly sacrifices output accuracy to achieve resource efficiency. On the other hand, accuracy metrics in Mitchell-based designs do not undergo notable changes. Finally, as already discussed in earlier work [17], some modular architectures sacrifice delay for LUT saving [7].
- **RAPID designs establish better resource-accuracy trade-off than leading-one based truncated counterparts:** the higher accuracy levels of DRUM and AAXD come at the cost of higher resource consumption than logarithmic counterparts (due to using an *accurate* core unit). Furthermore, employing an accurate instance of a divider still results in a long latency, multiple times that of a same-size multiplier. RAPID, on the other hand, has reduced the high latency of accurate divider, nearly to latency of its same-size multiplier. Overall, our proposed architectures yield better resource-accuracy trade-off, especially in higher bit-width. Finally, there are many cases with an error near or equal to 100% in the truncation-based AAXD divider. Such high error cases can result in false positive peaks in heartbeat and corner detection, as will be discussed later.
- **The proposed error-reduction outperforms other Mitchell-based architectures:** the error-reduction strategy of RAPID surpasses SIMDive/REALM in three aspects. First, the proposed error-reduction scheme achieves a higher accuracy level, even with fewer coefficients/LUTs. This is due to the better partitioning scheme of RAPID, while it still considers 4 MSBs of fractional parts. Second, the exponential increase in the number of coefficients (256 for 4 MSBs) of REALM/SIMDive not only poses a noticeable resource penalty, but also would nullify the gain, when realized in a LUT-based implementation. Overall, RAPID enables a more cost-effective strategy for partitioning the squarish zone. Third, in addition to error-reduction strategy, RAPID can achieve a higher throughput, energy, and throughput per Watt, compared to SIMDive, in non-pipeline configuration (except for 8- and 16-bit division). Albeit, as noted previously, the higher throughput

enabled by pipelining RAPID comes at the cost of increased throughput per Watt, when compared to SIMDive.

Please note although comparing the SIMD and pipelined architectures is not in the scope of this paper, herein some of the key differences between RAPID and SIMDive are discussed. Supporting higher throughput through SIMDive has added to the complexity of the sub-modules. For example,  $4 \times 3 = 12$  bit is used for LOD in 32-bit SIMDive or 2- and 4-MUX units are used to select the functionality and sub-word length in the intermediate adders. Moreover, integration of a multiplier and divider into a hybrid design marginally affects the complexity of the circuit through including 2-MUX units to select the mode in sub-modules. Finally, in order to support simultaneous scalings for sub-word length in SIMD mode, the LOD, error-reduction, and final shifter become more complex. For example, in a 32-bit LOD, instead of 5 bit for the SISD mode,  $4 \times 3 = 12$  bit is used to also support simultaneous 8-bit leading one detection for the SIMD mode (the complexity overhead is also posed to integer/fractional part adder).

- **Comparing latency, throughput and throughput per Watt in pipelined and non-pipelined designs:** overall, operating at a higher frequency and producing one-operation-per-cycle through pipelining has resulted in additional flip-flops, increased end-to-end latency, and higher dynamic power dissipation. However, pipelining enables significant improvement in throughput. In fact, higher throughput enabled by increasing the number of pipeline stages comes at the cost of lower *Throughput per Watt* as well. This descending trend is mainly due to the increase in the number of FFs and end-to-end latency in such designs. In particular, the analysis of latency and throughput per Watt highlights the efficacy of pipelining for RAPID circuits from three perspectives. First, comparing accurate- and RAPID-pipelined designs demonstrates that, even with increasing the number of pipelining stages, the end-to-end latency for each x-stage based RAPID remains smaller than its x-stage based accurate counterpart (except for 8-bit multiplication). Second, the x-stage based pipelined RAPID enables higher throughput per Watt versus its x-stage based *accurate* counterpart. Third, while increasing the number of pipeline stages results in a descending trend in throughput per Watt for RAPID multipliers and accurate IPs (both multiplier and divider), this trend is ascending for RAPID dividers. This means that increasing the pipeline depth is highly beneficial for dividers. Moreover, the *relative throughput per Watt* ratio for RAPID dividers are higher than 1, meaning that pipelined RAPID dividers achieve a better throughput per Watt than their accurate counterparts (even when compared to the non-pipelined mode). On the other hand, as discussed in Section II the reciprocal-based dividers are not suited for LUT-based platforms. In fact, the poor performance of pipelined SAADI-EC is due to two reasons: first, its datapath is divided into three non-uniform stages (normalizing block, multiplier, and error-correction accumulator based on Taylor iterations). Second, generating the reciprocal of the divisor (even through utilizing reciprocal IP of Vivado, as adopted herein) is a costly operation for LUT-based designs.

It is worth noting that Xilinx Vivado offers different implementation strategies to reduce the dynamic power, including `Power_DefaultOpt`, `Power_ExploreArea` (in which sequential resources are combined) and system-level power reduction techniques such as voltage scaling [62]. Overall, it is reported that such optimizations can enable up to 30% improvement in dynamic power [62]. However, considering the trade-off between area/power and performance, exploiting such techniques will increase the critical path delay. As pipelining has increased the critical path delay by itself, we have refrained from applying further directives, but users can utilize such options w.r.t. their constraints.



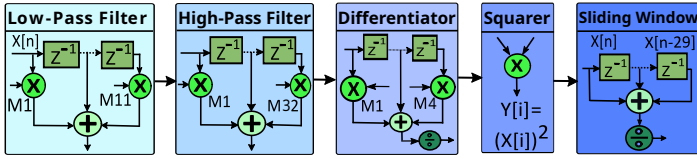


Fig. 5: The structure of Pan Tompkins QRS detection application [30]

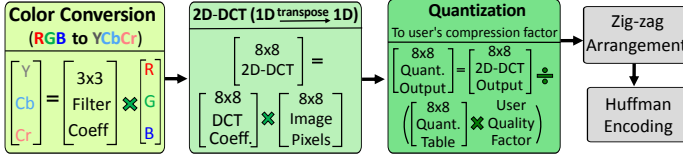


Fig. 6: The structure of JPEG Compression application [30]

### B. Evaluation of RAPID in three multi-kernel applications

The efficacy of RAPID over accurate and SoA multipliers<sup>3</sup> and dividers has been also appraised by deploying them in the end-to-end implementation of three multi-kernel applications. The application, shown in figures 5, 6, and 7 include JPEG compression, heart-beat detection through Pan-Tompkins algorithm, and HCD (the corners are further employed to generate movement vectors, which is used in object tracking programs).

**Hardware implementation and performance analysis of applications:** the source-code of applications is synthesized with Xilinx Vivado for the usually-adopted kernel configuration of 16-bit. The implementation of Pan-Tompkins algorithm is adopted from [63]. The base of JPEG compression is adopted from AxBench [64] and further optimized for a resource-efficient implementation on FPGA, by e.g., transforming 2D-DCT calculations to the butterfly-based 1D-DCT approach [64, 65]. We have developed both JPEG and HCD in C++ and synthesized them through Vivado High-Level Synthesis (HLS) and disabled DSPs. HLS has two key advantages: first, it facilitates applying various directives, in a system-level implementation as well as the process for generating different, customized configurations. Second, it simplifies the high-level behavioral evaluation of the entire design. To efficiently reflect the optimizations of RAPID in the final HDL design and to overcome the resource gap between HLS generated and HDL, we have employed a three-step approach. In the first step we have coded each of accurate multiplication and division functions in the applications. This scheme has facilitated replacing each function with its optimized RAPID versions, later in step 3. Second, the compiler has been forced, via HLS inline pragmas, to generate an independent HDL file for each of multiplication and division functions. In the third step, the HDL description of the respective functions are replaced by HDL-optimized versions of RAPID modules. Finally, for the end-to-end performance analysis, the HLS-synthesized designs for all accurate and approximate applications have been further passed to the downstream implementation phase, placed and routed on Virtex-7.

**QoR analysis on real-world benchmarks:** for assessing the end-to-end accuracy of Pan-Tompkins, we measured the QRS and Peak Signal to Noise Ratio (PSNR) for 30k ECG samples from the MIT-BIH database [66] in MATLAB. For quality measurements on JPEG-compressed images PSNR is used, while the percentage of correct vectors is considered as the application-level metric for HCD. In fact, similar to [67], the extracted corners from HCD algorithm are passed

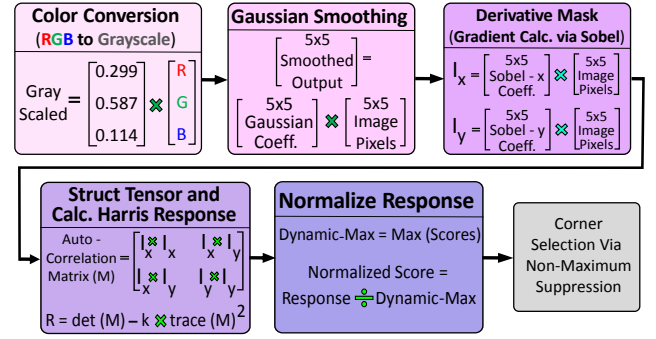


Fig. 7: The structure of Harris Corner Detection application [30]

to MATLAB for generating the motion/terrain vectors<sup>4</sup>. The accuracy changes in these applications are measured by conducting analysis on 50 images from three aerial imagery datasets [68, 69, 70]. To remain inline with industrial standards, we refrained from applying approximations to *zigzag* and *Huffman* kernels in JPEG which implement a re-arrangement and encoding scheme, respectively. Moreover, the *corner selection via non-maximum suppression* in HCD have also been remained accurate, as it is mainly comprised of comparison operations and has moderately low resource footprint. In this article we have considered the PSNR of, at least, 28 dB for ECG and JPEG compression and 90% correct vectors for HCD application (reported to be an acceptable confidence level for moving object tracking [67]). Please note, although a relatively low PSNR of 11-19 dB, and thereby, 95-100% detection ratio has been allowed in XBioSip [63], we adhere to a high PSNR in this article, that is also expected to render 100% detection accuracy [63].

**Post-implementation results:** Fig. 10 shows the obtained area, delay, and Area-Delay-Product (ADP) improvements of approximate designs, in the three applications<sup>5</sup>. The comparison is among three configurations: adopting with RAPID multipliers and dividers, adopting SoA SIMDive multiplier and divider, or dynamically truncated designs, i.e., DRUM-6 multiplier together with AAXD-8/4 divider. The following inferences are highlighted on the efficacy of RAPID designs, after appraising the end-to-end performance gain:

First, referring to Fig. 8 and Fig. 9 the quality drop by designs having a near-unbiased error characteristic (both RAPID and SIMDive) is smaller than when truncated counterparts are adopted. As discussed before, the biased error of truncated designs, i.e., DRUM and AAXD, results in accumulation of errors in consecutive kernels. In fact, our profiling has revealed that the near-zero biased errors of the RAPID multiplier and divider have been able to cancel out each other in consecutive operations/kernels and prevent a drastic error accumulation in the aggregation-based (mostly Add/Mul) structure of kernels. Such observations also corroborate [71, 72] in that error-bias feature plays a pivotal role in approximation of consecutive kernels having an aggregation-based structure, e.g., neural networks. Second, the negligible accuracy loss after deployment of RAPID is partly due to the small average and peak error of the proposed multiplier and divider, as compared to high peak error in DRUM/AAXD, (up to 100% in hundreds of input cases). In fact, the analysis shows that a high peak error is also reflected as incorrect vectors in the HCD application (especially due to the presence of division in the last stage of the HCD algorithm). This also holds true for false positive heartbeat (QRS peaks), as by utilizing both DRUM and AAXD the detection accuracy of heartbeat has dropped by ~1%.

<sup>4</sup>In [67], only DCT *adders* are approximated (DCT is used as a pre-processing compression step in UAV programs).

<sup>5</sup>Comprehensive energy analysis for multi-kernel programs on large image/ECG datasets is for follow-up track

<sup>3</sup>Herein the SISD mode of SIMDive has been analyzed. Detailed comparison of SIMD architectures with pipelining for multi-kernel applications and resolving instruction dependencies are targeted as interesting future works.



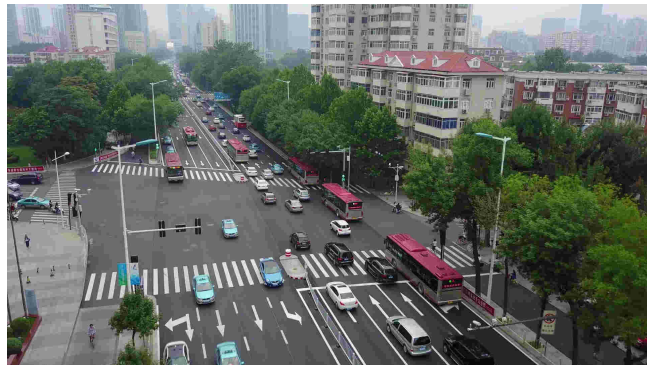
(a) Accurate multiplication and division (PSNR = 30.9)



(b) RAPID multiplier-10 and divider-9 (PSNR = 28.7)

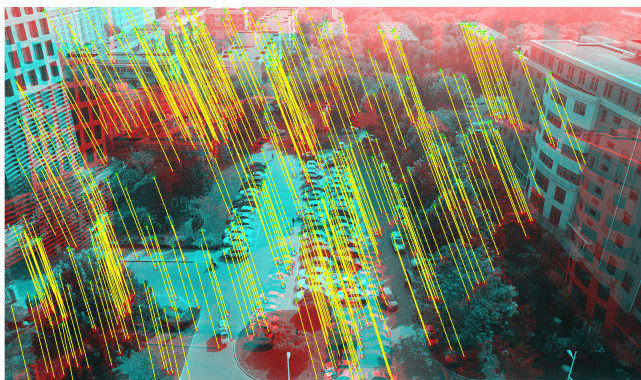


(c) SIMDive multiplication and division (PSNR = 29.3)

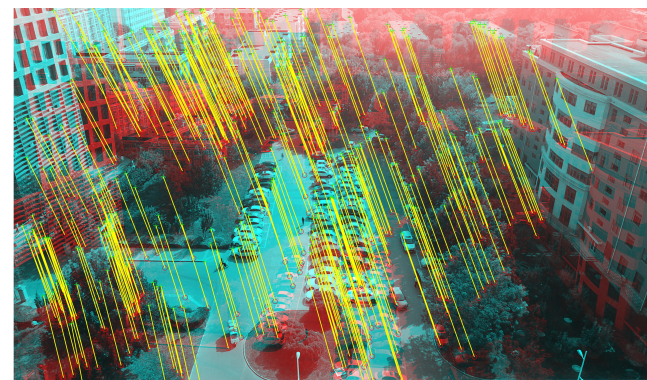


(d) DRUM-6 multiplication and AAXD-8/4 division (PSNR = 24.4)

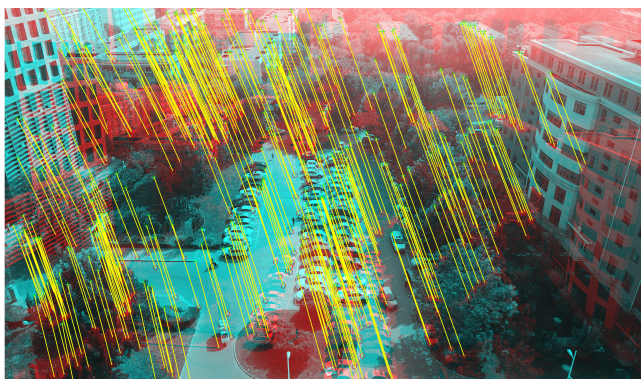
Fig. 8: Comparison of JPEG compression on aerial images with accurate and different approximate multipliers and dividers (16-bit).



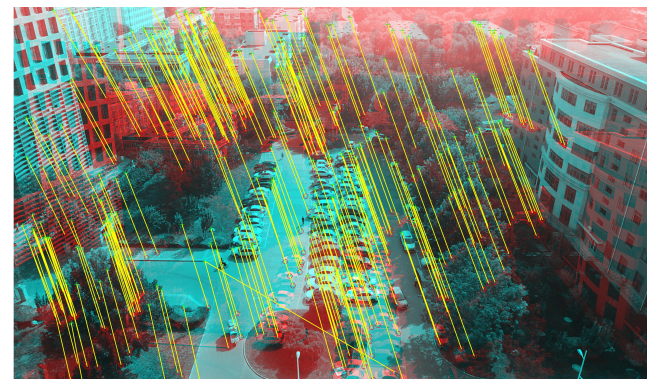
(a) Accurate multiplication and division (baseline) = 100%



(b) RAPID multiplier-10 and divider-9 = 94%



(c) SIMDive multiplication and division = 97%



(d) DRUM-6 multiplication and AAXD-8/4 division = 83%

Fig. 9: Tracking via Harris Corner Detection: changes in Harris score range also enables detection of new vectors in the threshold-based selection of tracking algorithm in MATLAB. Average of false positive vectors is less than 5%.

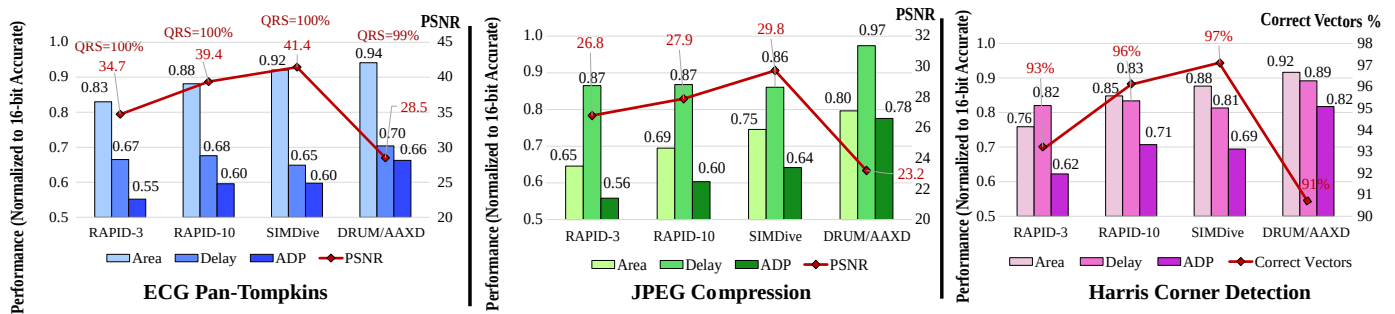


Fig. 10: End-to-end performance of applications utilizing RAPID multiplier and divider, compared to accurate and SoA approximate designs.

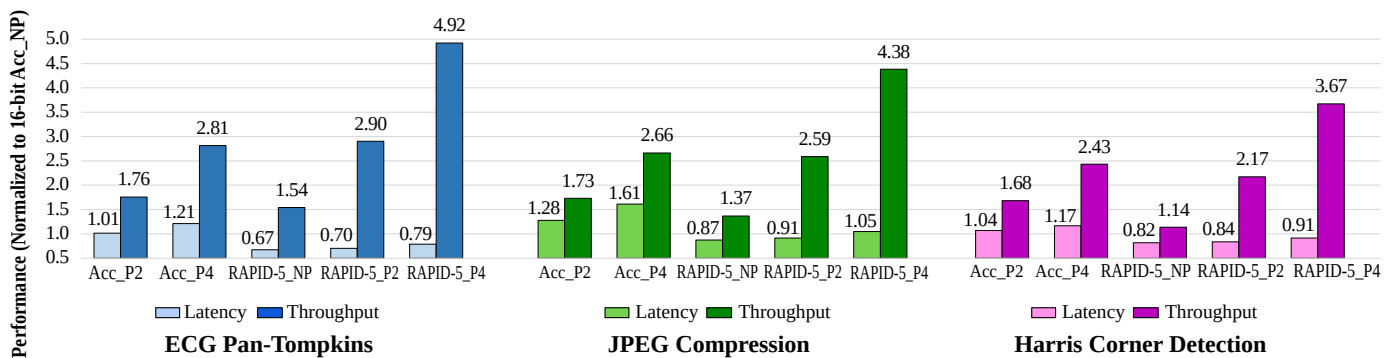


Fig. 11: The end-to-end latency and throughput of applications utilizing RAPID multiplier and divider, compared to accurate IP counterparts, in pipelined and non-pipelined format.

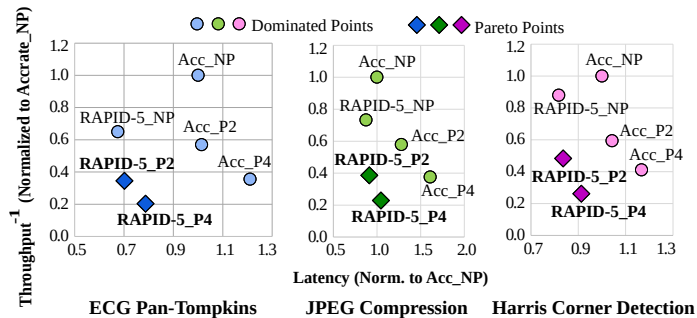


Fig. 12: Comparing the latency and throughput trade-off in applications exploiting pipelined and non-pipelined version of RAPID and accurate multiplier and divider.

Second, RAPID-configured applications also have better area- and ADP- gains, compared to both SIMDive and truncated counterparts (see Fig. 10). SIMDive on the other hand has marginally better end-to-end latency (by at most 3%), since its error-reduction circuitry is customized for performance-efficiency by directly configuring LUTs. Nonetheless, the difference of QoR and performance-gain between RAPID and SIMDive is not significant at application-level.

To compare the throughput of applications in pipelined and non-pipelined configuration, we have also deployed the 2- and 4-stage pipelined versions of the RAPID multiplier and divider (along with their accurate counterparts). For a fair comparison, we have avoided user-specified optimizations such as *function pipelining* pragmas (e.g., on matrix multiplication) and the applications are also implemented on the basis of streaming approach. Fig. 11 compares the end-to-end latency and throughput of both pipelined and non-pipelined configurations. Please note, the area difference by only replacing the non-pipelined RAPID multipliers and dividers with their pipelined versions is not significant (see Fig. 10). The throughput is estimated as the inverse of the clock period for the applications,

especially as they are constantly fed with bulk of data. As can be seen in Fig. 11, the throughput of the applications increases after applying the pipelining, but this comes with the cost of increase in the end-to-end latency. Nevertheless, the latency overhead of converting RAPID\_P2 configuration to RAPID\_P4 is less than of the overhead from converting Acc\_P2 to Acc\_P4. Similar observation also holds true for the improvement in the throughput. Moreover, application configurations having the RAPID\_P2 reach smaller latency along with the higher throughput than when incorporating Acc\_NP or Acc\_P2. Finally, Fig. 12 illustrates the Pareto points in the trade-off between Latency and throughput for pipelined and non-pipelined architectures. As can be observed, RAPID\_P2 and RAPID\_P4 render the Pareto points having better trade-off in latency and throughput than other configurations.

*Discussion:* it should be highlighted that although more complicated object tracking or heart arrhythmia programs utilize machine-learning techniques for feature extraction, continuously offloading the complete video stream (for the former) or patient bio-signal data (for the latter) to the insecure/untrustworthy network, or process/store it on the third-party cloud can pose performance bottleneck for a real-time processing and deplete the battery in a short interval. Therefore, enabling extraction of some features at the edge is highly desired.

## VI. CONCLUSIONS AND FUTURE WORK

In this paper, we have proposed RAPID, the first fine-grain pipelined architecture for approximate multiplier and divider. The proposed error-reduction schemes of RAPID enable 99-99.4% accuracy with smaller cost, compared to the existing approaches. In particular, pipelined RAPID multipliers (dividers) enable up to  $3.3\times$  ( $8.1\times$ ) higher throughput,  $2.3\times$  ( $6.8\times$ ) higher throughput/Watt, and 56% (36%) savings of LUTs, over pipelined accurate IPs. The end-to-end evaluations of RAPID in heartbeat detection, JPEG compression, and Harris corner detection demonstrate up to 35%, 33%, and 45% improvements in area, latency, and Area-Delay-Product,

respectively, over accurate configuration, with negligible loss in QoR. RAPID pipelined designs are interesting candidates to speed up the execution of a wide domain of stream-based applications that are constantly fed with a bulk of data. 35%, 33%, and 45% improvements in area, latency, and Area-Delay-Product (ADP), respectively, over accurate kernels, with negligible loss in QoR.

For future work, we intend to assess the efficacy of the RAPID pipeline mode in different application domains, e.g., Neural networks which offers both SIMD and pipelining opportunities. The challenge is resolving data dependencies in consecutive instructions. Such challenges are usually partially addressed through out-of-order execution in processors. However, this technique cannot fully utilize the pipelining potentials. Therefore, we target providing specialized versions of the pipelined multiplier and divider, which will be able to support internal data forwarding [55, 56] and able to resolve data dependencies. It should be noted that bypassing would be faster and posed with smaller overhead, when implemented through an intra-unit granularity.

Furthermore, we plan to design an approximate Arithmetic Logic Unit (ALU) and assess its applicability in the data-path of soft processors such as RISC-V. In fact, RAPID bears a great potential to be deployed in the mantissa multiplier/divider which consume more than 95% of the total area and power in the floating point unit (in which division latency is up to  $35\times$  of addition operation) [20, 73]. Recently, this track has attracted noticeable attention, especially due to the ever-growing usage of 3D computer graphics [74, 75].

#### ACKNOWLEDGEMENT

This research is co-funded by the projects *X-ReAp: Cross(X)-Layer Runtime Reconfigurable Approximate Architecture* (Number 380524764), funded by the German research foundation *Deutsche Forschungsgemeinschaft (DFG)* and *Re-learning: Self-learning and flexible electronics through inherent component reconfiguration* (Number 100382146), funded by the *European Social Fund (ESF)*.

#### REFERENCES

- [1] World Health Organisation. 2018. Cardiovascular diseases (CVDs). [https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-\(cvds\)](https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-(cvds)). (2018).
- [2] P. Kostic. 2017. Heart Disease and Early Heart Attack Care. [https://www.bnl.gov/hr/ocmed/hpp/linkable\\_files/pdf/EarlyHeartAttackSymptoms.pdf](https://www.bnl.gov/hr/ocmed/hpp/linkable_files/pdf/EarlyHeartAttackSymptoms.pdf). (2017).
- [3] Y. Yang et al. 2019. FPNNet: Customized Convolutional Neural Network for FPGA Platforms. In *IEEE International Conference on Field-Programmable Technology (ICFPT)*.
- [4] X. Gu et al. 2016. A Real-Time FPGA-Based Accelerator for ECG Analysis and Diagnosis Using Association-Rule Mining. *ACM Transactions on Embedded Computing Systems (TECS)*, 15, 2.
- [5] H.K. Chatterjee et al. 2015. Real-time detection of electrocardiogram wave features using template matching and implementation in FPGA. *International Journal of Biomedical Engineering and Technology (IJ-BET)*, 17, 3.
- [6] S. Ullah et al. 2021. High-Performance Accurate and Approximate Multipliers for FPGA-based Hardware Accelerators. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*.
- [7] S. Ullah et al. 2018. Area-Optimized Low-Latency Approximate Multipliers for FPGA-Based Hardware Accelerators. In *IEEE/ACM Design Automation Conference (DAC)*.
- [8] I. Kuon and J. Rose. 2007. Measuring the gap between fpgas and asics. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, 26, 2.
- [9] A. Boutros et al. 2018. Embracing Diversity: Enhanced DSP Blocks for Low-Precision Deep Learning on FPGAs. In *IEEE International Conference on Field Programmable Logic and Applications (FPL)*.
- [10] S. Lee et al. 2019. Double MAC on a DSP: Boosting the Performance of Convolutional Neural Networks on FPGAs. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, 38, 5.
- [11] Xilinx. 2015. LogiCORE IP multiplier v12.0. [https://www.xilinx.com/support/documentation/ip\\_documentation/mult\\_gen/v12\\_0/pg108-mult-gen.pdf](https://www.xilinx.com/support/documentation/ip_documentation/mult_gen/v12_0/pg108-mult-gen.pdf). (2015).
- [12] Xilinx. 2016. LogiCORE IP Divider v5.1. [https://www.xilinx.com/support/documentation/ip\\_documentation/div\\_gen/v5\\_1/pg151-div-gen.pdf](https://www.xilinx.com/support/documentation/ip_documentation/div_gen/v5_1/pg151-div-gen.pdf). (2016).
- [13] Xilinx. 2018. 7 Series DSP48E1 Slice. [https://www.xilinx.com/support/documentation/user\\_guides/ug479\\_7Series\\_DSP48E1.pdf](https://www.xilinx.com/support/documentation/user_guides/ug479_7Series_DSP48E1.pdf). (2018).
- [14] N. Van Toan and J. Lee. 2020. FPGA-Based Multi-Level Approximate Multipliers for High-Performance Error-Resilient Applications. *IEEE Access*, 8.
- [15] Z. Ebrahimi et al. 2020. SIMDive: Approximate SIMD Soft Multiplier-Divider for FPGAs with Tunable Accuracy. In *ACM Great Lakes Symposium on VLSI (GLSVLSI)*.
- [16] H. Saadat et al. 2019. Approximate Integer and Floating-Point Dividers with Near-Zero Error Bias. In *IEEE/ACM Design Automation Conference (DAC)*.
- [17] Z. Ebrahimi et al. 2020. LeAp: Leading-one Detection-based Softcore Approximate Multipliers with Tunable Accuracy. In *Asia & South Pacific Design Automation Conference (ASP-DAC)*.
- [18] J. N. Mitchell. 1962. Computer Multiplication and Division using Binary Logarithms. *IRE Transactions on Electronic Computers (IRETEC)*, 11, 4.
- [19] A. R. Baranwal et al. 2020. ReLAccS: A Multi-level Approach to Accelerator Design for Reinforcement Learning on FPGA-based Systems. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*.
- [20] H. Saadat et al. 2018. Minimally Biased Multipliers for Approximate Integer and Floating-Point Multiplication. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, 37, 11.
- [21] H. Jiang et al. 2020. Approximate Arithmetic Circuits: A Survey, Characterization, and Recent Applications. *Proceedings of the IEEE*, 108, 12.
- [22] S. Rehman et al. 2016. Architectural-space exploration of approximate multipliers. In *IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*.
- [23] M. Wang et al. 2020. An Optimized Compression Strategy for Compressor-Based Approximate Multiplier. In *IEEE International Symposium on Circuits and Systems (ISCAS)*.
- [24] A. G. M. Strollo et al. 2020. Comparison and Extension of Approximate 4-2 Compressors for Low-Power Approximate Multipliers. *IEEE Transactions on Circuits and Systems I (TCAS-I): Regular Papers*, 67, 9.
- [25] P. J. Edavoor et al. 2020. Approximate Multiplier Design Using Novel Dual-Stage 4:2 Compressors. *IEEE Access*, 8.
- [26] D. Esposito et al. 2018. Approximate Multipliers Based on New Approximate Compressors. *IEEE Transactions on Circuits and Systems I (TCAS-I): Regular Papers*, 65, 12.
- [27] S. Venkatchalam and S. Ko. 2017. Design of Power and Area Efficient Approximate Multipliers. *IEEE Transactions on Very Large Scale Integration Systems (TVLSI)*, 25, 5.
- [28] O. Akbari et al. 2017. Dual-Quality 4:2 Compressors for Utilizing in Dynamic Accuracy Configurable Multipliers. *IEEE Transactions on Very Large Scale Integration Systems (TVLSI)*, 25, 4.
- [29] Y. Guo et al. 2020. Small-Area and Low-Power FPGA-Based Multipliers using Approximate Elementary Modules. In *Asia & South Pacific Design Automation Conference (ASP-DAC)*.
- [30] Z. Ebrahimi, et al. 2021. Plasticine: A Cross-Layer Approximation Methodology for Multi-Kernel Applications through Minimally Biased, High-Throughput, and Energy-Efficient SIMD Soft Multiplier-Divider. *ACM Transactions on Design Automation of Electronic Systems (TO-DAES)*, 27, 2.
- [31] M. S. Ansari et al. 2018. Low-Power Approximate Multipliers Using Encoded Partial Products and Approximate Compressors. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems (JETCAS)*, 8, 3.
- [32] V. Mrazek et al. 2017. EvoApprox8b: Library of Approximate Adders and Multipliers for Circuit Design and Benchmarking of Approximation Methods. In *Design, Automation & Test in Europe (DATE)*.
- [33] S. Ullah et al. 2018. SMApProxLib: Library of FPGA-based Approximate Multipliers. In *IEEE/ACM Design Automation Conference (DAC)*.
- [34] F. Frustaci et al. 2020. Approximate Multipliers With Dynamic Truncation for Energy Reduction via Graceful Quality Degradation. *IEEE Transactions on Circuits and Systems II (TCAS-II): Express Briefs*, 67, 12.

- [35] S. Vahdat et al. 2019. TOSAM: An Energy-Efficient Truncation- and Rounding-Based Scalable Approximate Multiplier. *IEEE Transactions on Very Large Scale Integration Systems (TVLSI)*, 27, 5.
- [36] R. Zendegani et al. 2017. RoBA Multiplier: A Rounding-Based Approximate Multiplier for High-Speed yet Energy-Efficient Digital Signal Processing. *IEEE Transactions on Very Large Scale Integration Systems (TVLSI)*, 25, 2.
- [37] H. Jiang et al. 2019. Low-Power Unsigned Divider and Square Root Circuit Designs Using Adaptive Approximation. *IEEE Transactions on Computers (TC)*, 68, 11.
- [38] H. Jiang et al. 2018. Adaptive approximation in arithmetic circuits: A low-power unsigned divider design. In *Design, Automation & Test in Europe (DATE)*.
- [39] M. Vaeztourshizi et al. 2018. An Energy-Efficient, Yet Highly-Accurate, Approximate Non-Iterative Divider. In *International Symposium on Low Power Electronics and Design (ISLPED)*.
- [40] S. Vahdat et al. 2017. TruncApp: A truncation-based approximate divider for energy efficient DSP applications. In *Design, Automation & Test in Europe (DATE)*.
- [41] R. Zendegani et al. 2016. SEERAD: A high speed yet energy-efficient rounding-based approximate divider. In *Design, Automation & Test in Europe (DATE)*.
- [42] J. Melchert et al. 2019. SAADI-EC: A Quality-Configurable Approximate Divider for Energy Efficiency. *IEEE Transactions on Very Large Scale Integration Systems (TVLSI)*, 27, 11.
- [43] M. S. Ansari et al. 2019. A Hardware-Efficient Logarithmic Multiplier with Improved Accuracy. In *Design, Automation & Test in Europe (DATE)*.
- [44] W. Liu et al. 2018. Design and Evaluation of Approximate Logarithmic Multipliers for Low Power Error-Tolerant Applications. *IEEE Transactions on Circuits and Systems I (TCAS-I): Regular Papers*, 65, 9.
- [45] H. Saadat et al. 2020. REALM: Reduced-Error Approximate Log-based Integer Multiplier. In *Design, Automation & Test in Europe (DATE)*.
- [46] Zahra Ebrahimi and Akash Kumar. 2021. BioCare: An Energy-Efficient CGRA for Bio-Signal Processing at the Edge. In *IEEE International Symposium on Circuits and Systems (ISCAS)*.
- [47] S. Hashemi et al. 2015. DRUM: A Dynamic Range Unbiased Multiplier for approximate applications. In *IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*.
- [48] S. Hashemi et al. 2016. A low-power dynamic divider for approximate applications. In *IEEE/ACM Design Automation Conference (DAC)*.
- [49] E. Adams et al. 2020. Approximate Restoring Dividers Using Inexact Cells and Estimation From Partial Remainders. *IEEE Transactions on Computers (TC)*, 69, 4.
- [50] S. Venkatachalam et al. 2019. Design of Approximate Restoring Dividers. In *IEEE International Symposium on Circuits and Systems (ISCAS)*.
- [51] L. Chen et al. 2018. Design, Evaluation and Application of Approximate High-Radix Dividers. *IEEE Transactions on Multi-Scale Computing Systems (TMSCS)*, 4, 3.
- [52] L. Chen et al. 2016. On the Design of Approximate Restoring Dividers for Error-Tolerant Applications. *IEEE Transactions on Computers (TC)*, 65, 8.
- [53] S. Behroozi et al. 2019. SAADI: A Scalable Accuracy Approximate Divider for Dynamic Energy-Quality Scaling. In *Asia & South Pacific Design Automation Conference (ASP-DAC)*.
- [54] S. Ullah et al. 2020. Area-optimized Accurate and Approximate Soft-core Signed Multiplier Architectures. *IEEE Transactions on Computers (TC)*.
- [55] J. Seo and D. H. Kim. 2019. Dependency-Resolving Intra-Unit Pipeline Architecture for High-Throughput Multipliers. In *Design, Automation & Test in Europe (DATE)*.
- [56] J. Seo and D. H. Kim. 2019. High-Throughput Multiplier Architectures Enabled by Intra-Unit Fast Forwarding. In *IEEE International Symposium on Computer Arithmetic (ARITH)*.
- [57] S. Scarfone et al. 2021. Design and Analysis of a Leading One Detectorbased Approximate Multiplier on FPGA. In *International Conference on Synthesis, Modeling, Analysis and Simulation Methods, and Applications to Circuit Design (SMACD)*.
- [58] K. Chapman. 2013. Multiplexer design techniques for datapath performance with minimized routing resources (Xilinx Application Note). <https://docs.xilinx.com/v/u/en-US/xapp522-mux-design-techniques>. (2013).
- [59] Xilinx. 2013. Xilinx 7 Series FPGA Programmable Guide for HDL Designs. [https://www.xilinx.com/support/documentation/sw\\_manuals/xilinx14\\_7/7series\\_hdl.pdf](https://www.xilinx.com/support/documentation/sw_manuals/xilinx14_7/7series_hdl.pdf). (2013).
- [60] B. Parhami. 2010. *Computer Arithmetic: Algorithms and Hardware Designs*. Volume 20. Oxford university press.
- [61] ARM. 2021. Performance per Watt Is the New Moore's Law. <https://www.arm.com/blogs/blueprint/performance-per-watt>. (2021).
- [62] Xilinx. 2021. Vivado Design Suite User Guide: Power Analysis and Optimization. [https://www.xilinx.com/content/dam/xilinx/support/documents/sw\\_manuals/xilinx2021\\_2/ug907-vivado-power-analysis-optimization.pdf](https://www.xilinx.com/content/dam/xilinx/support/documents/sw_manuals/xilinx2021_2/ug907-vivado-power-analysis-optimization.pdf)<https://docs.xilinx.com/v/u/2019.1-English/ug907-vivado-power-analysis-optimizationversion2019andbetter>. (2021).
- [63] B. Prabhakaran et al. 2019. XBioSiP: A Methodology for Approximate Bio-Signal Processing at the Edge. In *IEEE/ACM Design Automation Conference (DAC)*.
- [64] A. Yazdanbakhsh et al. 2017. AxBench: A Multiplatform Benchmark Suite for Approximate Computing. *IEEE Design & Test*, 34, 2.
- [65] M. Jridi et al. 2013. Low complexity DCT engine for image and video compression. In *Real-Time Image and Video Processing (RTIVP)*. Volume 8656.
- [66] A. L. Goldberger et al. 2000. PhysioBank, PhysioToolkit, and PhysioNet: Components of a New Research Resource for Complex Physiologic Signals. *Circulation*.
- [67] T. Nomani et al. 2020. xUAVs: Towards Efficient Approximate Computing for UAVs—Low Power Approximate Adders With Single LUT Delay for FPGA-Based Aerial Imaging Optimization. *IEEE Access*, 8.
- [68] M. Mueller et al. 2016. A Benchmark and Simulator for UAV Tracking. In *European Conference on Computer Vision (ECCV)*.
- [69] H. Fan et al. 2020. VisDrone-SOT2020: The Vision Meets Drone Single Object Tracking Challenge Results. In *Workshops in European Conference on Computer Vision (ECCV)*.
- [70] Y. Lyu et al. 2020. UAVid: A semantic segmentation dataset for UAV imagery. *ISPRS Journal of Photogrammetry and Remote Sensing*, 165.
- [71] M. S. Ansari et al. 2020. Improving the Accuracy and Hardware Efficiency of Neural Networks Using Approximate Multipliers. *IEEE Transactions on Very Large Scale Integration Systems (TVLSI)*, 28, 2.
- [72] Z. G. Tasoulas et al. 2020. Weight-Oriented Approximation for Energy-Efficient Neural Network Inference Accelerators. *IEEE Transactions on Circuits and Systems I (TCAS-I): Regular Papers*, 67, 12.
- [73] S. Tamimi et al. 2019. An Efficient SRAM-Based Reconfigurable Architecture for Embedded Processors. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, 38, 3.
- [74] Juwon Yun et al. 2020. A latency-effective pipelined divider for double-precision floating-point numbers. *IEEE Access*, 8.
- [75] Yuheng Yang et al. 2021. An architecture of area-effective high radix floating-point divider with low-power consumption. *IEEE Access*, 9.



**Zahra Ebrahimi** received her B.Sc. and M.Sc. degrees in computer engineering at Sharif University of Technology (SUT), Iran, in 2014 and 2016, respectively. Meanwhile, she was also a Research Assistant at the Data Storage, Networks, and Processing Laboratory, at SUT. She started her Ph.D. at the Center for Advancing Electronics Dresden (cfaed), Technische Universität Dresden, Germany, in 2018. Her research interests include approximate computing, reconfigurable accelerator design, and energy-efficient edge computing.



**Muhammad Zaid** received his B.Sc. degree in Electrical Engineering from National University of Sciences and Technology (NUST), Islamabad, Pakistan, in 2017. He is currently pursuing the M.Sc. in Nanoelectronic Systems at Technische Universität Dresden. His research interest includes HW/SW co-design of embedded AI systems.



**Mark Wijnvliet** received his MS.c. and Ph.D. degrees at the Electronic Systems group at Eindhoven university of technology in 2011 and 2020, respectively. His Ph.D. topic focused at energy efficient reconfigurable processor architectures. Afterwards, he worked as a post-doctoral researcher at TU Dresden at the processor design chair in the field of hardware security. Currently he works as a researcher at ASMPT in the Netherlands. His interests include reconfigurable hardware, hardware security, processor optimization, chip design, and space applications.



**Akash Kumar** (SM'13) received the joint Ph.D. degree in electrical engineering and embedded systems from the Eindhoven University of Technology, Eindhoven, The Netherlands, and the National University of Singapore (NUS), Singapore, in 2009. From 2009 to 2015, he was with NUS. He is currently a Professor with Technische Universität Dresden, Dresden, Germany, where he is directing the Chair for Processor Design. His current research interests include the design, analysis, and resource management of low-power and fault-tolerant embedded multiprocessor systems.

ded multiprocessor systems.