# A Review, Classification and Comparative Evaluation of Approximate Arithmetic Circuits

HONGLAN JIANG, University of Alberta CONG LIU, University of Alberta LEIBO LIU, Tsinghua University FABRIZIO LOMBARDI, Northeastern University JIE HAN, University of Alberta

Often as the most important arithmetic modules in a processor, adders, multipliers and dividers determine the performance and the energy efficiency of many computing tasks. The demand of higher speed and power efficiency, as well as the feature of error resilience in many applications (e.g., multimedia, recognition and data analytics), have driven the development of approximate arithmetic design. In this article, a review and classification are presented for the current designs of approximate arithmetic circuits including adders, multipliers and dividers. A comprehensive and comparative evaluation of their error and circuit characteristics is performed for understanding the features of various designs. By using approximate multipliers and adders, the circuit for an image processing application consumes as little as 47% of the power and 36% of the power-delay product of an accurate design while achieving a similar image processing quality. Improvements in delay, power and area are obtained for the detection of differences in images by using approximate dividers.

CCS Concepts: •General and reference  $\rightarrow$  Surveys and overviews; Evaluation; Measurement; •Hardware  $\rightarrow$  Arithmetic and datapath circuits; Combinational circuits;

Additional Key Words and Phrases: Approximate computing, approximate circuit, adder, multiplier, divider, hardware, accuracy, image processing.

### 1. INTRODUCTION

While computational errors are in general not desirable, applications such as multimedia, wireless communication, recognition, and data mining are tolerant of the occurrence of some errors [Han and Orshansky 2013]. Due to the perceptual limitations of humans, these errors do not make an obvious difference in applications such as image, audio and video processing. Moreover, in many digital signal processing (DSP) systems, inputs from the outside world are noisy, so there is a limit in precision or accuracy in the computed results. Many applications are based on statistical or probabilistic computation, such as classification and recognition algorithms. Due to the nature of these applications, trivial errors in computation do not result in a significant performance degradation. Therefore, approximate computing is applicable in many applications that can tolerate the loss of certain accuracy [Venkatesan et al. 2010].

© 2017 ACM. 1550-4832/2017/07-ART60 \$15.00 D01: http://dx.doi.org/10.1145/3094124

This work was partly supported by the University of Alberta and the Natural Sciences and Engineering Research Council (NSERC) of Canada.

Author's addresses: H. Jiang, C. Liu and J. Han, Department of Electrical and Computer Engineering, University of Alberta, Edmonton, AB T6G 1H9, Canada; email: {honglan, cong4, jhan8}@ualberta.ca; L. Liu, Institute of Microelectronics, Tsinghua University, Beijing 100084, China; email: liulb@tsinghua.edu.cn; F. Lombardi, Department of Electrical and Computer Engineering, Northeastern University, Boston, MA 02115, USA; email: lombardi@ece.neu.edu.

ACM acknowledges that this contribution was co-authored by an affiliate of the national government of Canada. As such, the Crown in Right of Canada retains an equal interest in the copyright. Reprints must include clear attribution to ACM and the author's government agency affiliation. Permission to make digital or hard copies for personal or classroom use is granted. Copies must bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. To copy otherwise, distribute, republish, or post, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Past research on approximate computing has spanned from circuits to programming languages [Han 2016]. In [Datla et al. 2009], an approximate squaring circuit is proposed. A new logic synthesis approach is introduced to reduce the area of a synthesized circuit for a given threshold of error rate in [Shin 2010]. In [Chippa et al. 2010], a scalable-effort design approach is proposed to implement highly efficient hardware for error-resilient applications. Automated design processes have been proposed for approximate digital circuit design using Cartesian genetic programming [Vasicek and Sekanina 2015; Mrazek et al. 2016]. Sampson et al. have developed the EnerJ language, an extension to Java [Sampson et al. 2011]. This language supports approximate data types for low-power computation. Various computing and memory architectures have been proposed for supporting approximate computing applications [Esmaeilzadeh et al. 2012; Miguel et al. 2015]. In this article, we focus on approximate circuit designs and particularly approximate arithmetic circuits of adders, multipliers and dividers.

Design metrics and analytical approaches have been proposed for the evaluation of approximate adders [Liang et al. 2013; Huang et al. 2012; Miao et al. 2012; Venkatesan et al. 2010; Liu et al. 2015; Mazahir et al. 2017]. Monte Carlo simulation has been employed to acquire data for analysis. In this article, the error rate (ER), the error distance (ED) and the average error are used to evaluate the error characteristics of the approximate designs. Hardware related figures of merit including critical path delay, circuit area and power dissipation, as well as compound metrics including the power-delay product (PDP) and area-delay product (ADP), are utilized to assess the circuit characteristics of these designs.

Image processing has been essential in diverse applications including multimedia, biomedical imaging and pattern recognition [Acharya and Ray 2005]. Taking advantage of its inherent error resilience, image processing can be efficiently implemented by using approximate arithmetic circuits. Therefore, image sharpening and change detection are considered for further evaluation of the approximate circuits in addition to the evaluation using design metrics. The simulation results show that the image sharpening circuit using approximate adders and multipliers saves as much as 53% of power and 58% of area compared to an accurate design with a similar accuracy. The change detection circuit using approximate dividers achieves as much as 40% improvement in speed and 25% improvement in power compared with an accurate design at a similar accuracy. While some preliminary results have been presented in [Jiang et al. 2015] and [Jiang et al. 2016b], this article makes the following new contributions.

1. A larger set of approximate adders that include more recent designs are considered in this article. Moreover, a new error metric (average error) is used to measure the output bias of an approximate design, which is very important in an accumulative operation.

2. Approximate Booth multipliers for signed multiplication are included in the evaluation, while only unsigned approximate multipliers are considered in [Jiang et al. 2016b].

3. Current approximate dividers are reviewed and their features are discussed with respect to performance, accuracy and hardware consumption.

4. The STM CMOS 28nm process is used in the circuit synthesis throughout the work in this article, while an older 65nm technology was used in [Jiang et al. 2015].

5. The considered approximate arithmetic circuits of adders, multipliers and dividers are applied to two image processing applications for further evaluation. The accuracy and circuit characteristics are obtained by simulation and synthesis.



Fig. 1. The *n*-bit ripple-carry adder (RCA). FA: a 1-bit full adder.



Fig. 2. The *n*-bit carry lookahead adder (CLA). SPG: the cell used to produce the sum, generate  $(g_i = a_i b_i)$  and propagate  $(p_i = a_i + b_i)$  signals.

## 2. APPROXIMATE ADDERS

An adder performs the addition of two binary numbers. Two basic adders are the ripple-carry adder (RCA) (Figure 1) and the carry lookahead adder (CLA) (Figure 2). In an *n*-bit RCA, the carry of each full adder (FA) is propagated to the next FA, thus the delay and circuit complexity increase proportionally with n (or O(n)). An *n*-bit CLA consists of n units that operate in parallel to produce the sum and the generate  $(g_i = a_i b_i)$  and propagate  $(p_i = a_i + b_i)$  signals for generating the lookahead carries. The delay of CLA is logarithmic in n (or O(log(n))), thus significantly shorter than for RCA. However, a CLA requires a larger circuit area (in O(nlog(n))), incurring a higher power dissipation.

Many approximation schemes have been proposed by reducing the critical path and hardware complexity of an accurate adder. An early methodology is based on a speculative operation [Lu 2004; Verma et al. 2008]. In an *n*-bit speculative adder, each sum bit is predicted by its previous k less significant bits (LSBs) (k < n). As the carry chain is shorter than n, a speculative adder is faster than a conventional design. A segmented adder is implemented by several smaller adders operating in parallel [Mohapatra et al. 2011; Zhu et al. 2009; Kahng and Kang 2012; Yang et al. 2016]. Hence, the carry propagation chain is truncated into shorter segments. Segmentation is also utilized in [Du et al. 2012; Kim et al. 2013; Ye et al. 2013; Lin et al. 2015; Li and Zhou 2014; Hu and Qian 2015; Miao et al. 2012; Camus et al. 2015; 2016], but the carry input for each sub-adder is selected differently. This type of adder is referred to as a carry select adder. Another method for reducing the critical path delay and power dissipation is by approximating a full adder [Mahdiani et al. 2010; Gupta et al. 2013; Yang et al. 2013; Cai et al. 2016; Angizi et al. 2017]. The approximate full adder is then used to implement the LSBs in an accurate adder. Thus, approximate adders are divided into four categories, as briefly summarized below.



Fig. 3. The almost correct adder (ACA).  $\Box$ : the carry propagation path of the sum.



Fig. 4. The equal segmentation adder (ESA). k: the maximum carry chain length; l: the size of the first sub-adder ( $l \le k$ ).

#### 2.1. Classification of Approximate Adders

2.1.1. Speculative adders. The almost correct adder (ACA) [Verma et al. 2008] is based on the speculative adder design of [Lu 2004]. In an *n*-bit ACA, k LSBs are used to predict the carry for each sum bit (n > k), as shown in Figure 3. Therefore, the critical path delay is reduced to O(log(k)) (for a parallel implementation such as CLA, the same below). The design in [Lu 2004] requires (n - k) k-bit sub-carry generators in an *n*-bit adder and thus, the hardware consumption is rather high (in O((n - k)klog(k)))). This overhead is reduced in [Verma et al. 2008] by sharing some components among the sub-carry generators.

2.1.2. Segmented adders. The equal segmentation adder (ESA) divides an *n*-bit adder into a number of smaller *k*-bit sub-adders operating in parallel with fixed carry inputs, so no carry is propagated among the sub-adders (Figure 4) [Mohapatra et al. 2011]. The delay of ESA is O(log(k)) and the circuit complexity is O(nlog(k)). Its hardware overhead is significantly lower than ACA.

The error-tolerant adder type II (ETAII) consists of parallel carry generators and sum generators [Zhu et al. 2009], as shown in Figure 5. The carry signal from the previous carry generator propagates to the next sum generator. Therefore, ETAII utilizes more information to predict the carry and thus it is more accurate than ESA for the same k. The circuit of ETAII is more complex than that of ESA, and its delay is larger due to the longer critical path (2k).

In an *n*-bit accuracy-configurable approximate adder (ACAA),  $\left\lceil \frac{n}{k} - 1 \right\rceil 2k$ -bit subadders are required [Kahng and Kang 2012]. Each sub-adder adds 2k consecutive bits



Fig. 5. The error-tolerant adder type II (ETAII) [Zhu et al. 2009]: the carry propagates through the two shaded blocks.



Fig. 6. The speculative carry selection adder (SCSA).

with an overlap of k bits and all 2k-bit sub-adders operate in parallel to reduce the delay to O(log(k)). In each sub-adder, half of the most significant sum bits is selected as the partial sum. The accuracy of ACAA can be configured at runtime. Moreover, ACAA has the same carry propagation path as ETAII for each sum, so they are equally accurate for the same k.

The dithering adder divides an adder into an accurate, more significant sub-adder and a less significant sub-adder with upper and lower bounding modules [Miao et al. 2012]. The output of the less significant sub-adder is conditionally selected. An effective "Dither Control" enables a smaller variance in the overall error.

To reduce the error distance, an error control and compensation method is proposed for a segmented adder in [Yang et al. 2016]. This method employs a multistage latency to compensate the carry prediction error in a more significant segmentation, thus trading off computing efficiency for an improved accuracy.

The delays of the segmented adders are O(log(k)) and the circuit complexities are O(nlog(k)) for ESA and ETAII, and O((n-k)log(k)) for ACAA.

2.1.3. Carry select adders. In a carry select adder, several signals are commonly used.

For the *i*<sup>th</sup> block, generate  $g_{i,j} = a_{i,j}b_{i,j}$ , propagate  $p_{i,j} = a_{i,j} \oplus b_{i,j}$ , and  $P_i = \prod_{j=0}^{\kappa-1} p_{i,j}$ , where  $\alpha$  and  $b_j$  are the *i*<sup>th</sup> LSPs of the input energy  $P_i = 1$  indicates that all  $b_i$ 

where  $a_{i,j}$  and  $b_{i,j}$  are the  $j^{th}$  LSBs of the input operands.  $P_i = 1$  indicates that all k propagate signals in the  $i^{th}$  block are true.

An *n*-bit speculative carry select adder (SCSA) consists of  $m = \lfloor \frac{n}{k} \rfloor$  sub-adders (or window adders) [Du et al. 2012]. Each sub-adder is made of two *k*-bit adders: adder0

with carry-in "0" and adder1 with carry-in "1". The carry-out of adder0 is connected to a multiplexer to select the addition result as part of the final result, as shown in Figure 6. SCSA and ETAII achieve the same accuracy for the same value of k due to the same carry predict function, while SCSA uses an additional adder and multiplexer in each block.

Similar to SCSA, an *n*-bit adder is divided into  $\lceil \frac{n}{k} \rceil$  blocks in the carry skip adder (CSA) [Kim et al. 2013]. Each block in CSA consists of a sub-carry generator and a sub-adder. The carry-in of the  $(i+1)^{th}$  sub-adder is determined by the propagate signals of the  $i^{th}$  block: it is the carry-out of the  $(i-1)^{th}$  sub-carry generator when all propagate signals are true ( $P_i = 1$ ), otherwise it is the carry-out of the  $i^{th}$  sub-carry generator. Therefore, the critical path delay of CSA is O(log(k)). This carry select scheme improves the carry prediction accuracy.

Different from SCSA, the carry speculative adder (CSPA) in [Lin et al. 2015] contains one sum generator, two internal carry generators (one with carry-0 and one with carry-1) and one carry predictor in each block. The output of the  $i^{th}$  carry predictor is used to select carry signals for the  $(i + 1)^{th}$  sum generator. l input bits (rather than k, l < k) in a block are used in a carry predictor. Therefore, the hardware overhead is reduced compared to SCSA.

The consistent carry approximate adder (CCA) is similar to SCSA in that each block of CCA consists of adders with carry-0 and carry-1 [Li and Zhou 2014]. The select signal of a multiplexer is determined by the propagate signal, i.e.,  $S_i = (P_i + P_{i-1})SC + (P_i + P_{i-1})C_{i-1}$ , where  $C_{i-1}$  is the carry-out of the  $(i-1)^{th}$  adder0 and SC is a global speculative carry. In CCA, the carry prediction depends not only on its LSBs, but also on the higher bits; its critical path delay is similar to that of SCSA.

The generate signals-exploited carry speculation adder (GCSA) has a similar structure as CSA and uses the generate signals for carry speculation [Hu and Qian 2015]. The difference between them lies in the carry selection; the carry-in for the  $(i + 1)^{th}$  sub-adder is selected by its own propagate signals rather than its previous block. The carry-in is the most significant generate signal  $g_{i,k-1}$  of the  $i^{th}$  block if  $P_i = 1$ , or else it is the carry-out of the  $i^{th}$  sub-carry generator. This carry selection scheme effectively controls the maximum relative error.

In the gracefully-degrading accuracy-configurable adder (GDA) [Ye et al. 2013], the control signals are used to configure the accuracy by selecting an accurate or approximate carry-in signal using a multiplexer for each sub-adder. The delay of GDA is determined by the carry propagation and thus by the control signals to the multiplexers.

In the carry cut-back adder (CCBA) [Camus et al. 2016], the full carry propagation is prevented by a controlled multiplexer or an OR gate for a high-speed operation. The multiplexer is controlled by a carry propagate block at a higher-significance position to cut the carry propagation at a lower-significance position. The delay and accuracy of the CCBA largely depend on the distance between the propagate block and the cutting multiplexer, thus allowing a high accuracy with a marginal overhead.

The critical path delays of the carry select adders are given by O(log(k)), where k is the size of the sub-adder.

2.1.4. Approximate full adders. In this type of design, approximate full adders are implemented in the LSBs of a multibit adder. It includes the simple use of OR gates (and one AND gate for carry propagation) in the so-called lower-part-OR adder (LOA) (Figure 7) [Mahdiani et al. 2010], the approximate designs of the mirror adder (AMAs) [Gupta et al. 2013] and the approximate XOR/XNOR-based full adders (AXAs) [Yang et al. 2013]. Additionally, emerging technologies such as magnetic tunnel junctions



60:7

Fig. 7. The lower-part-OR adder (LOA).

have been considered for the design of approximate full adders for a shorter delay, a smaller area and a lower power consumption [Cai et al. 2016; Angizi et al. 2017].

The critical path of this type of adders depends on its approximation scheme. For LOA, it is approximately O(log(n - l)), where l is the number of bits in the lower part of an adder. In the evaluation, LOA is selected as the reference design because the other designs require customized layouts at the transistor level; hence, they are not comparable with the other types of approximate adders that are approximated at the logic gate level. Finally, an adder with the LSBs truncated is referred to as a truncated adder that works with a lower precision. It is considered as a baseline design.

#### 2.2. Evaluation of Approximate Adders

2.2.1. Error Characteristics. Monte Carlo simulation is performed to evaluate the accuracy of the approximate adders. The error distance (ED) and the relative error distance (RED) are calculated as: ED = |M' - M| and  $RED = \frac{ED}{M}$ , where M' is the approximate result and M is the accurate result [Liang et al. 2013]. The mean error distance (MED) is the mean of all possible EDs. The error rate (ER), the probability of producing an incorrect result), the normalized MED (NMED, the normalization of MED by the maximum output of the accurate design) and the mean relative error distance (MRED), the average value of all possible REDs) are used to assess the error characteristics of the approximate designs. Moreover, the average error (the mean of all possible errors (M' - M)) is used to evaluate the bias of an approximate arithmetic design.

The functions of 16-bit approximate adders are simulated by MATLAB using 10 million uniformly distributed random input combinations. Table I shows the simulation results. The size of the carry predictor for CSPA is  $\lceil k/2 \rceil$  in this evaluation. The global speculative carry *SC* for CCA is "0", which is proved to be more accurate than using "1". Additionally, the adder with *k* LSBs truncated (TruA-*k*) is simulated for comparison.

As shown in Table I, ETAII, ACAA and SCSA have the same error characteristics (in ER, NMED and MRED) due to the same carry propagation chain for each sum bit. The NMED and MRED show the same trend, so only MRED and ER are considered in the comparison, as shown in Figure 8. An equivalent carry propagation chain is selected for the considered approximate adders i.e., the parameter k for ACA, ESA, LOA and TruA is 8, while it is 4 for CSA, GCSA, ETAII, ACAA, SCSA, CCA and CSPA. These approximate adders are considered as equivalent approximate adders.

Adder Type	ER <b>(%)</b>	NMED (10 <sup>-3</sup> )	MRED (10 <sup>-3</sup> )	Average Error
		Speculative A	dders	
ACA-4	16.66	7.80	18.90	-1024.6
ACA-5	7.76	3.90	9.60	-511.7
		Segmented A	dders	
ESA-4	85.07	15.70	40.40	-2047.5
ESA-5	80.03	7.80	20.80	-1023.4
ETAII-4	5.85	0.97	2.60	-127.5
ETAII-5	2.28	0.24	0.65	-31.6
ACAA-4	5.85	0.97	2.60	-127.5
ACAA-5	2.29	0.24	0.65	-31.6
		Carry Select A	Adders	
SCSA-4	5.85	0.97	2.60	-127.5
SCSA-5	2.28	0.24	0.65	-31.6
CSA-4	0.18	0.06	0.15	-7.4
CSA-5	0.02	0.004	0.01	-0.5
CSPA-4	29.82	3.90	10.40	-511.4
CSPA-5	11.31	0.98	2.70	-128.3
CCA-4	8.71	0.98	2.00	-128.3
CCA-5	3.78	0.25	0.49	-32.2
GCSA-4	4.26	0.48	0.98	-63.2
GCSA-5	1.52	0.12	0.25	-16.1
		Approximate Fu	ll Adders	
LOA-6	82.19	0.09	0.25	0.2
LOA-8	89.99	0.37	1.00	0.2
		Truncated A	dders	
TruA-6	99.98	0.48	1.30	-63.0
TruA-8	100.0	1.95	5.40	-255.0

Table I. Simulation Results of the Error Characteristics for the Approximate Adders

*Note:* The number following the name of each approximate adder is the number of LSBs used for the carry speculation in the speculative adders, the length of the segmentation in the segmented adders, and the number of approximated and truncated LSBs in the approximate full adder-based and truncated adders.

Among these approximate adders, CSA is the most accurate, and GCSA is the second most accurate in terms of MRED. LOA has a different structure from the other approximate adders. Its more significant part is fully accurate, while the approximate part is less significant. Therefore, the MRED of LOA is rather small, but its ER is very large. For a similar reason, TruA has the highest ER and very large MRED. The information used to predict each carry in ESA and CSPA is rather limited, so the ERand the MRED of ESA and CSPA are larger than most of the other approximate designs. Compared with the other approximate adders, CCA, ETAII, SCSA and ACAA show moderate ER and MRED. In terms of average error, LOA has the lowest value because it produces both positive and negative errors that can compensate each other, while errors are accumulated for the other approximate adders since only negative errors are generated. Therefore, LOA is suitable for an accumulative operation.

In summary, the carry select adders and the speculative adder (ACA) are very accurate with small values of ER and MRED (except for CSPA using a small number of bits for carry prediction). Represented by LOA, an approximate full adder based adder has a moderate MRED, the lowest average error but a very large ER. The segmented







(b)

Fig. 8. A comparison of error characteristics of the approximate adders. (a) ER and (b) MRED. The parameter, k, is 4 for CSA, GCSA, ETAII, ACAA, SCSA, CCA and CSPA, and it is 8 for ACA, ESA, LOA and TruA for an equivalent carry propagation chain.

adders are not very accurate in terms of NMED and MRED. With very large values of ER and MRED, the truncated adder is the least accurate among the equivalent designs. Three different types of approximate adders, ETAII, ACAA and SCSA, have the same error characteristics.

2.2.2. Circuit Characteristics. To assess the circuit characteristics, 16-bit approximate adders and the accurate CLA are implemented in VHDL and synthesized using the Synopsys Design Compiler (DC) based on an STM CMOS 28 nm process with a supply voltage of 1.0 V at a temperature of 25°C. For a fair comparison, all designs use the same process, voltage and temperature with the same optimization option. Both the



Fig. 9. A comparison of delay and power of the approximate adders (a) delay, (b) power and (c) power-delay product.

Addam Truna	Delay	Power	PDP	Area	ADP
Adder Type	( <i>ps</i> )	(uW)	(fJ)	$(um^2)$	$(um^2.ns)$
CLAC	1000	65.9	65.9	60.7	60.7
CLAG	570	105.4	60.1	84.2	48.0
	$\operatorname{Spe}$	eculative	Adders		
ACA-4	250	118.4	29.6	73.8	18.5
ACA-5	270	119.4	32.2	71.8	19.4
	Seg	gmented	Adders		
ESA-4	260	47.0	12.2	49.9	13.0
ESA-5	310	50.6	15.7	51.7	16.0
ETAII-4	550	80.6	44.3	71.6	39.4
ETAII-5	670	78.5	52.6	70.2	47.0
ACAA-4	550	80.9	44.5	70.8	38.9
ACAA-5	650	87.3	56.8	74.6	48.5
	Car	ry Select	Adders	5	
SCSA-4	320	134.5	43.0	109.2	34.9
SCSA-5	400	163.0	65.2	126.2	50.5
CSA-4	390	97.8	38.1	142.5	55.6
CSA-5	420	94.3	39.6	131.2	55.1
CSPA-4	300	89.2	26.8	83.7	25.1
CSPA-5	370	117.6	43.5	100.7	37.3
CCA-4	320	172.6	55.2	131.4	42.0
CCA-5	420	209.5	88.0	155.0	65.1
GCSA-4	380	109.7	41.7	74.3	28.2
GCSA-5	460	113.6	52.3	73.3	33.7
	Approx	ximate F	ull Add	ers	
LOA-6	440	75.1	33.0	58.8	25.9
LOA-8	390	66.9	26.1	53.2	20.8
-	Tr	uncated A	Adders		
TruA-6	390	67.9	26.5	52.4	20.4
$TruA-8^1$	350	64.2	22.5	46.2	16.2

Table II. Circuit Characteristics of the Approximate Adders

<sup>1</sup> TruA-8 is synthesized at a medium mapping effort, different from the high mapping effort used for the other designs. In this case, TruA-8 attains a shorter critical path delay, but a similar PDP and ADP as the results synthesized using the high mapping effort.

P-channel and the N-channel transistors used in the designs have a typical design corner with a regular threshold voltage. The critical path delay and area are reported by the Synopsys DC. Power dissipation is measured by the PrimeTime-PX tool at 1 ns clock period with 10 million random input combinations. All adders and sub-adders are implemented as CLA in this article, unless otherwise noted.

Table II reports the results for the delay, power dissipation, power-delay product (PDP), circuit area and area-delay product (ADP) of the considered adders. Two structures of the accurate CLA are implemented: CLAC is realized by four cascaded 4-bit CLAs, while CLAG is realized by four parallel 4-bit CLAs and a carry look-ahead generator. Among ETAII, SCSA and ACAA (with the same error characteristics when the same value of k is selected), SCSA, albeit being the fastest, incurs the largest power dissipation and area because two sub-adders and one multiplexer are utilized in each



Fig. 10. A comprehensive comparison of the approximate adders: (a) ER and PDP and (b) MRED and PDP. The parameter k for LOA and TruA ranges from 9 down to 3 from left to right, it is 3 to 8 for ESA and ACA, and it is from 3 to 6 for the other adders from left to right. The adders marked by circles are equivalent in terms of carry propagation and are thus representatives of different designs.

block. ACAA is very slow due to its long critical path. The block of ETAII (a carry generator and a sum generator) is significantly simpler than those of SCSA and ACAA. Therefore, ETAII consumes less power and requires a smaller area than SCSA and ACAA.

In general a circuit with a larger area is likely to consume more power except for CSA with a relatively low power dissipation but a large area. This is due to its short critical

path and enhanced carry select that results in a complex wiring. Figure 9 shows the delay, power and PDP of the equivalent adders. As expected, the accurate CLAC has the longest delay among all adders, but not the highest power dissipation. Compared to CLAC, CLAG is significantly faster and consumes more power and area. TruA is not very fast, but it is the most power and area efficient design. LOA is also very power and area efficient compared with the other approximate adders. ESA is the slowest, but it is very power and area efficient due to its simple segmentation structure. CCA is very fast but is the most power and area consuming design due to its complex speculative circuit. Both CSPA and GCSA have moderate power dissipations, but CSPA is faster and GCSA uses a smaller area. Both the speed and power dissipation of CSA are in the medium range. In terms of PDP and ADP (shown in Table II), they show very similar trend. TruA, LOA and CSPA have very small values of PDP and ADP, while these values are relatively large for ACA and CCA (shown in Figure 9(c)).

In conclusion, the carry select adders are likely to have large values of power dissipation and area at a moderate performance. The segmented adders are power and area efficient. A speculative adder is very fast, but it is also very power consuming with a moderate area. Conversely, the approximate full adder based adder is slow, but it consumes a low power and area. The approximate full adders are very efficient in PDP and ADP, while the speculative adders are not. The truncated adder is very power and area efficient, but with a relatively long delay.

2.2.3. Discussion. Since the ADP shows a similar trend as the PDP, the PDP is considered for a comprehensive comparison of the approximate adders, as shown in the two-dimensional (2-D) plots of Figure 10. CSA-6 is accurate due to the precise carry generated for every block, so the ER and MRED of CSA-6 are 0. Therefore, they are not shown in Figure 10. The equivalent adders are marked by circles. Among adders with the same accuracy (ETAII, SCSA and ACAA), ETAII is the most efficient in terms of delay, power and area. Thus, it is shown as a representative in Figure 10. Compared with the other approximate adders, CCA has the largest PDP and moderate ER and MRED. Among the schemes with moderate PDPs (CSPA, GCSA and ETAII), ETAII and GCSA have moderate MREDs and ERs, while CSPA shows slightly higher values of these measures. ESA has a rather small PDP, but a considerably large ER and MRED. ACA has a larger PDP than ESA, but it has both lower ER and MRED. Among all approximate adders, CSA shows the best performance with very small PDP, ER and MRED values.

With the highest ERs, LOA and TruA show the smallest PDPs for a similar MRED due to their low power dissipation. In fact, these approximate adders show a decent tradeoff in error distance and hardware efficiency. In particular, they are useful in applications in which hardware efficiency is of the utmost importance.

## 3. APPROXIMATE MULTIPLIERS

## 3.1. Classification of Approximate Multipliers

Generally, a multiplier consists of stages of partial product generation, accumulation and a final addition, as shown in Figure 11 for a  $4 \times 4$  unsigned multiplication. Let  $A_i$ and  $B_j$  be the  $i^{th}$  and  $j^{th}$  least significant bits of inputs A and B respectively, a partial product  $P_{j,i}$  is usually generated by an AND gate (i.e.,  $P_{j,i} = A_i B_j$ ). The commonly used partial product accumulation structures include the Wallace, Dadda trees and a carry-save adder array. The Wallace tree for a  $4 \times 4$  unsigned multiplier is shown in the dotted box of Figure 11. The adders in each layer operate in parallel without carry propagation, and the same operation repeats until two rows of partial products are left. For an *n*-bit multiplier, log(n) layers are required in a Wallace tree. Therefore, the delay of the partial product accumulation stage is O(log(n)). Moreover, the adders in Figure



Fig. 11. The basic arithmetic process of a  $4 \times 4$  unsigned multiplication with possible truncations to a limited width.  $\bullet$ : an input, a partial product or an output product;  $\bigcirc$ : a truncated bit;  $\Box$ : a full adder or a half adder.



Fig. 12. Partial product accumulation of a  $4 \times 4$  unsigned multiplier using a carry-save adder array. HA: a half adder; FA: a full adder.

11 can be considered as a (3:2) compressor and can be replaced by other counters or compressors (e.g. a (4:2) compressor) to further reduce the delay. The Dadda tree has a similar structure as the Wallace tree, but it uses as few adders as possible.

A carry-save adder array is shown in Figure 12; the carry and sum signals generated by the adders in a row are passed on to the adders in the next row. Adders in a column operate in series. Hence the partial product accumulation delay of an *n*-bit multiplier is approximately O(n), longer than that of the Wallace tree. However, an array requires a smaller area due to the simple and symmetric structure.

Three main methodologies are used for the approximate design of a multiplier: i) approximation in generating the partial products [Kulkarni et al. 2011], ii) approximation (including truncation) in the partial product tree [Mahdiani et al. 2010; Kyaw et al. 2010; Bhardwaj et al. 2014], and iii) using approximate designs of adders [Liu et al. 2014], counters [Lin and Lin 2013] or compressors [Ma et al. 2013; Momeni et al. 2015] to accumulate the partial products. For a signed integer operation, Booth multi-

$M_2 M_1 M_0$		$B_1B_0$					
		00	01	11	10		
	00	000	000	000	000		
$A_1A_0$	01	000	001	011	010		
	11	000	011	111	110		
	10	000	010	110	100		

Table III. K-Map for the  $2 \times 2$  Underdesigned Multiplier Block.

pliers have been widely used due to the fast operation on a reduced number of partial products. Some recent designs use shifting and addition to obtain the final product by rounding the inputs to a form of  $2^m$  (*m* is a positive integer) [Hashemi et al. 2015; Zendegani et al. 2017; Mitchell Jr 1962].

Based on the different schemes in approximation, approximate multipliers are classified into three unsigned types and signed Booth multipliers. Following this classification, existing designs of approximate multipliers are briefly reviewed next.

3.1.1. Approximation in generating partial products. The underdesigned multiplier (UDM) utilizes an approximate  $2 \times 2$  multiplier obtained by altering a single entry in the Karnaugh Map (K-Map) of its function (as highlighted in Table III) [Kulkarni et al. 2011]. Table III shows the K-Map of the approximate  $2 \times 2$  multiplier, where  $A_1A_0$  and  $B_1B_0$  are the two 2-bit inputs, and  $M_2M_1M_0$  is the 3-bit output. In this approximation, the accurate multiplication result "1001" is simplified to "111" to save one output bit when both the inputs are "11". Assuming the value of each input bit is equally likely, the error rate of the  $2 \times 2$  multiplier is then  $(\frac{1}{2})^4 = \frac{1}{16}$ . Larger multipliers can be designed based on the  $2 \times 2$  multiplier. This multiplier introduces an error when generating the partial products, however the adder tree remains accurate.

3.1.2. Approximation in the partial product tree. A bio-inspired imprecise multiplier referred to as a broken-array multiplier (BAM) is proposed in [Mahdiani et al. 2010]. The BAM operates by omitting some carry-save adders in an array multiplier in both horizontal and vertical directions (Figure 13). A more straightforward approach to truncation is to truncate some LSBs on the input operands and thus, a smaller multiplier is sufficient for the remaining MSBs. This truncated multiplier (TruM) is considered as a baseline design.

The error tolerant multiplier (ETM) is divided into a multiplication section for the MSBs and a non-multiplication section for the LSBs [Kyaw et al. 2010]. Figure 14 shows the architecture of a 16-bit ETM. A NOR gate based control block is used to deal with the following two cases: i) if the product of the MSBs is zero, then the upper accurate 8-bit multiplier is activated to multiply the LSBs without any approximation, and ii) if the product of the MSBs is nonzero, the non-multiplication section is used as an approximate multiplier to process the LSBs, while the multiplication section is activated to accurately multiply the MSBs.

The static segment multiplier (SSM) was proposed using a similar partition scheme [Narayanamoorthy et al. 2015]. Different from ETM, no approximation is applied to the LSBs in the SSM. Either the MSBs or the LSBs of the operands are accurately multiplied depending on whether its MSBs are all zeros. [Liu et al. 2017a] has shown that only small improvements in accuracy and hardware are achieved compared to ETM, thus this design is not further considered in this comparison study.



Fig. 13. The broken-array multiplier (BAM) with 4 vertical lines and 2 horizontal lines omitted [Mahdiani et al. 2010].  $\Box$ : a carry-save adder cell.



Fig. 14. The 16-bit error-tolerant multiplier (ETM) of [Kyaw et al. 2010].

A power and area-efficient approximate Wallace tree multiplier (AWTM) is based on a bit-width aware approximate multiplication and a carry-in prediction method [Bhardwaj et al. 2014]. An *n*-bit AWTM is implemented by four n/2-bit sub-multipliers, as shown in Figure 15, where the most significant sub-multiplier  $A_HB_H$  is further implemented by four n/4-bit sub-multipliers. The AWTM is configured into four different modes by the number of approximate n/4-bit sub-multipliers in the most significant n/2-bit sub-multiplier, while the other three multipliers ( $A_HB_L$ ,  $A_LB_H$  and  $A_LB_L$ ) are approximate. The approximate partial products are then accumulated by a Wallace tree.

3.1.3. Using approximate counters or compressors in the partial product tree. An approximate (4:2) counter is proposed in [Lin and Lin 2013] for an inaccurate 4-bit Wallace multiplier. Table IV shows the K-Map of the approximate (4:2) counter, where  $X_1 \cdots X_4$  are the four input signals of a (4:2) counter (i.e., the partial products in the partial product tree of a multiplier), C and S are the carry and sum, respectively. The values of CS in the box are approximated as "10" for "100" in the approximate counter when all input signals are "1." As the probability of obtaining a partial product of "1" is  $\frac{1}{4}$ , the error rate of the approximate (4:2) counter is  $(\frac{1}{4})^4 = \frac{1}{256}$ . The inaccurate 4-bit multiplier is then used to construct larger multipliers with error detection and correction circuits.

In the compressor based multiplier, accurate (3:2) and (4:2) compressors are improved to speed up the partial product accumulation [Baran et al. 2010]. By using the improved compressors, better energy and delay characteristics are obtained for a multiplier. To further reduce delay and power, two approximate (4:2) compressor designs



Fig. 15. The approximate Wallace tree multiplier (AWTM) [Bhardwaj et al. 2014]. n is the width of the multiplier,  $A_HB_H$ ,  $A_LB_H$ ,  $A_HB_L$  and  $A_LB_L$  are partial products generated by the n/2-bit sub-multipliers,  $A_{HH}B_{HH}$ ,  $A_{HL}B_{HH}$ ,  $A_{HH}B_{HL}$  and  $A_{HL}B_{HL}$  are partial products generated by the n/4-bit sub-multipliers.

Table IV. K-Map for the 4 : 2 Approximate Counter

CS		$X_1 \overline{X_0}$					
	00	01	11	10			
$X_{3}X_{4}$	00	00	01	10	01		
	01	01	10	11	10		
	11	10	11	10	11		
	10	01	10	11	10		

(AC1 and AC2) are presented in [Momeni et al. 2015]; these compressors are used in a Dadda multiplier with four different schemes. Approximate counters in which the more significant output bits are ignored, are presented and evaluated in [Kelly et al. 2009]. Several signed multipliers are then implemented using these approximate counters. The more accurate schemes 3 and 4 of the approximate compressor based multiplier (referred to as ACM-3 and ACM-4) in [Momeni et al. 2015] are considered for comparison.

In the approximate multiplier with configurable error recovery, the partial products are accumulated by a novel approximate adder (Figure 16) [Liu et al. 2014]. The approximate adder utilizes two adjacent inputs to generate a sum and an error bit. The adder processes data in parallel, thus no carry propagation is required. Two approximate error accumulation schemes are then proposed to alleviate the error of the approximate multiplier (due to the approximate adder). OR gates are used in the first error accumulation stage in scheme 1 (AM1), while in scheme 2 (AM2), both OR gates and the approximate adders are used. The truncation of 16 LSBs in the partial products in AM1 and AM2 results in TAM1 and TAM2 respectively [Liu 2014].

3.1.4. Approximate Booth multipliers. The Booth recoding algorithm handles binary numbers in 2's complement. The modified (or radix-4) Booth algorithm is commonly used due to its ease in generating partial products. Little work has been reported for



Fig. 16. The approximate adder cell.  $S_i$ : the sum bit;  $E_i$ : the error bit [Liu et al. 2014].



Fig. 17. The partial products for an 8-bit fixed-width modified Booth multiplier [Cho et al. 2004].  $P_{i,j}$  is the  $j^{th}$  partial product in the  $i^{th}$  partial product vector and  $n_i$  is the sign of the  $i^{th}$  partial product vector.

approximate Booth multipliers, whereas the fixed-width Booth multiplier utilizes a truncation-based approach has been studied for more than a decade. The conventional post-truncated fixed-width multiplier generates an output with the same width as the input operand by truncating the lower half of the product. Truncation of half of the partial products is widely used because the post-truncated scheme does not achieve a significant circuit advantage over the accurate multiplier. A direct truncation of partial products incurs a large error, so many error compensation schemes have been proposed [Cho et al. 2004; Min-An et al. 2007; Wang et al. 2011; Chen and Chang 2012]. Another approach is to use an approximate Booth encoder with a simple circuit [Liu et al. 2017b]. Most of the approximate Booth multipliers are based on the modified Booth algorithm; the partial products are accumulated by an array structure in [Cho et al. 2004; Min-An et al. 2007; Wang et al. 2011; Farshchi et al. 2013] while a parallel carry-save-adder tree is used in [Chen and Chang 2012; Liu et al. 2017b]. The approximate Booth multiplier in [Jiang et al. 2016a] is based on the radix-8 Booth algorithm.

Figure 17 shows the partial products of an 8-bit fixed-width modified Booth multiplier with error compensation [Cho et al. 2004]. The final product is the addition of the main part (MP) and the carry signals generated in the truncation part (TP). The carry signals are approximated by the output of Booth encoders. The approximate carry  $\sigma$  is

 $\sigma = \left[2^{-1}\left(\sum_{i=0}^{n/2-2} \overline{zero_i} + 1\right)\right], \text{ where } n \text{ is the multiplier width, and } \overline{zero_i} \text{ is "1" if the } i^{th}$ 

partial product vector is not zero or  $\overline{zero_i} = 0$  otherwise. This multiplier is referred to as BM04.

The multiplier in [Min-An et al. 2007] can adaptively compensate the quantization error by keeping different numbers of the most significant columns of the partial products ( $\omega$  ( $\omega \ge 0$ )). Two types of binary thresholding are proposed for error compensation. Different from BM04, *n* rather than (*n* - 1) columns of partial products are truncated

for an *n*-bit multiplier. The error compensation for each type of binary thresholding varies with the value of  $\omega$  and the partial products of the  $\omega^{th}$  column in the truncation part (from left to right). This multiplier is denoted as BM07.

The multiplier presented in [Wang et al. 2011] uses n columns of partial products in the truncation part for an *n*-bit multiplier; the most significant one column in the truncation part is reserved for error compensation. The error compensation using a simplified sorting network significantly reduces the mean and mean-squared errors by making the error symmetric as well as centralizing the error distribution around zero. This design is referred to as BM11.

A fixed-width Booth multiplier was designed based on a probabilistic estimation bias in [Chen and Chang 2012]. Therefore, this multiplier is referred to as PEBM. The number of columns of the accumulated partial products varies in accordance with the desired trade-off between hardware and accuracy. The error compensation formula is derived from a probability analysis rather than a time-consuming exhaustive simulation. The carry generated by the truncation part is approximated by

 $\sigma = \left[ 2^{-1} \left( \sum_{i=0}^{n/2-1-\lfloor \omega/2 \rfloor} z_i - 1 \right) \right], \text{ where } z_i = P_{0,n/2-1} + n_{n/2-1} \text{ when } i \text{ is } n/2 - 1 \text{ and}$ 

 $z_i = \overline{zero_i}$  otherwise.

Based on BAM, the broken Booth multiplier (BBM) uses a modified Booth algorithm to generate partial products and omits carry-save adders to the right of a vertical line [Farshchi et al. 2013]. BBM has a smaller power-delay product (PDP) for the same mean-squared error compared to BAM.

An approximate recoding adder is proposed in [Jiang et al. 2016a] for calculating the triple multiplicands to reduce the additional delay encountered in a radix-8 Booth multiplier. A Wallace tree and a truncation technique are then utilized for partial product accumulation to reduce power and delay. The most efficient fixed-width multiplier ABM2\_R15 is considered in this comparison and referred to as ABM2 in this article for simplicity. In [Liu et al. 2017b], two approximate Radix-4 Booth encoders are designed for the partial product generation by simplifying the exact K-Map. The generated partial products are then accumulated by using exact 4-2 compressors.

## 3.2. Evaluation of Approximate Multipliers

3.2.1. Error Characteristics. The considered ( $16 \times 16$ ) approximate multipliers are simulated by MATLAB with 10 million uniformly distributed random input combinations. The ER, NMED, MRED and average error are obtained and shown in Table V. TruMk represents the truncated multiplier with k LSBs truncated in the input operands.

According to Table V, most of the designs, especially those with truncation, have large *ER*s close to 100%. However, ICM has a relatively low *ER* of 5.45%, because it uses just one approximate counter in a  $4 \times 4$  sub-multiplier with an error rate of only  $\frac{1}{256}$ . UDM also shows a lower ER than the other approximate multipliers. In terms of the average error, ACMs have the smallest value, while the average errors for all the other approximate unsigned multipliers show the same trend with the NMED. This is because ACMs produce both positive and negative errors, but the other approximate unsigned multipliers produce either negative or positive errors.

Figure 18 shows the NMEDs and MREDs of the equivalent approximate multipliers that are configured to have 16-bit accurate MSBs (except for ICM and UDM that have only one configuration). Thus, the truncated LSBs in the partial product is 16 for BAM, the number of MSBs used for error compensation is 16 for AM1, AM2, TAM1 and TAM2, the size of the accurate sub-multiplier is 8 for ETM, 8 LSBs are truncated for TruM, and the mode number of ACM and AWTM is 4. Among the unsigned approximate multipliers, UDM has the largest *NMED* while ACM has the smallest.

Multiplier Type	ER (%)	NMED (10 <sup>-3</sup> )	MRED (%)	Average Error $(10^4)$
Multiplie	rs with Ar	proximation in C	enerating Par	rtial Products
UDM	80.99	13.92	3.33	-5974.9
Multipl	iers with	Approximation in	the Partial P	roduct Tree
BAM-16	99.99	0.06	0.21	-24.6
<b>BAM-17</b>	99.99	0.11	0.36	-49.2
<b>BAM-18</b>	99.99	0.22	0.63	-95.0
<b>BAM-20</b>	100.00	0.79	1.79	-337.5
ETM-7	99.99	0.97	1.56	-413.4
ETM-8	100.00	1.94	2.85	-825.1
AWTM-1	100.00	2.70	75.00	1159.6
AWTM-2	100.00	1.88	39.37	807.4
AWTM-3	99.99	0.12	2.51	51.5
AWTM-4	99.94	0.02	0.33	8.6
Multi	pliers usin	g Approximate C	ounters or Co	mpressors
ICM	5.45	0.29	0.06	-124.2
ACM-3	99.99	0.01	0.29	3.95
ACM-4	99.97	0.01	0.26	1.44
AM1-13	99.38	0.90	0.54	-385.5
AM1-16	98.22	0.81	0.34	-347.9
AM2-10	99.64	0.88	1.20	-379.8
AM2-13	99.36	0.35	0.34	-148.9
AM2-16	97.96	0.27	0.13	-115.1
TAM1-13	99.99	1.14	0.77	-488.5
TAM1-16	99.99	1.06	0.58	-457.1
TAM2-10	99.99	0.90	1.27	-384.2
TAM2-13	99.99	0.36	0.41	-153.3
TAM2-16	97.99	0.28	0.22	-121.0
	Tru	ncated Unsigned	Multipliers	
TruM-4	99.61	0.11	0.23	-49.2
TruM-8	100.0	1.94	2.85	-834.1
	Apj	proximate Booth	Multipliers	
$\mathbf{PEBM}$	99.99	0.023	0.27	-1.02
BBM	100.00	0.092	0.57	-9.83
BM11	99.99	0.022	0.18	-0.003
BM07	99.99	0.024	0.16	-1.56
BM04	99.99	0.027	0.48	-2.66
ABM2	99.99	0.034	0.44	-0.614

Table V. Error Characteristics of the Approximate Multipliers

*Note:* The parameter k follows the acronym of each approximate multiplier. For AM1, AM2, TAM1 and TAM2, this parameter refers to the number of MSBs used for error reduction and for ETM, the number of LSBs in the inaccurate part. It is the mode number in AWTM and ACM, and the vertical broken length (VBL) for BAM.

ICM, AM2 and TAM2 have similar values of NMED, however ICM has the smallest MRED, while the MRED of TAM2 is the largest. Therefore ICM has the highest accuracy in terms of MRED, while TAM2 is the least accurate among these three approximate multipliers. This indicates that multipliers with simple truncation tend to have larger MREDs when their NMEDs are similar. BAM has moderate values of NMED and MRED, while ETM and TruM have both large MRED and NMED.



(b)

Fig. 18. Comparison of *NMED* and *MRED* of the approximate multipliers with increasing (a) *NMED* and (b) *MRED*. ACM and AWTM represent ACM-4 and AWTM-4, respectively. The truncated number of LSBs in the partial product is 16 for BAM, the number of MSBs used for error compensation is 16 for AM1, AM2, TAM1 and TAM2, and 8 LSBs are truncated for TruM. ETM is ETM-8 in Table V.

Hence, ICM is the most accurate design with the lowest *ER*, *MRED* and a moderate *NMED*. ACM, AWTM, BAM, AM2 and TAM2 also show good accuracy among all considered approximate multipliers with both low *NMED*s and *MRED*s. ETM, TruM and UDM are not very accurate in terms of these metrics.

For the approximate Booth multipliers in Table V, a column of the most significant partial products in the truncation part (adjacent to the MP part) is kept for PEBM, BM07 and BM04. 15-bit columns of partial products are truncated in BBM and ABM2

to keep the same width of the output as the other designs. The ERs of all approximate multipliers are close to 100% due to truncation. Most designs have similar NMEDs except for BBM. BBM has the largest NMED and MRED, because there is no error compensation. BM07 and BM11 have very small MRED values, while PEBM has a slightly larger value. BM11 has the lowest average error and the average error of ABM2 is also very small.

In summary, as a multiplier approximated in generating the partial products, UDM has very large values of NMED and MRED, and a relatively small ER. The multipliers approximated in the partial product tree mostly have moderate NMEDs and very large MREDs (except for BAMs with fewer than or equal to 18 truncated bits and AWTM-4). The multipliers approximated using approximate counters or compressors have small values of both NMED and MRED, while the multiplier truncated on the input operands have large values of both metrics (when the truncated number of LSBs is larger than 4). Among the considered approximate Booth multipliers, BBM shows the lowest accuracy in terms of both NMED and MRED. BM11 has the smallest average error. The other approximate Booth multipliers show similar NMEDs and various MREDs.

3.2.2. Circuit Characteristics. The  $16 \times 16$  approximate multipliers are implemented in VHDL and synthesized using the same tool and process as in the simulation of approximate adders. The only difference is that the clock period is 4 ns for the power estimation of the multipliers because of a longer critical path delay. The accurate Wallace multiplier (WallaceM) optimized for speed [Oklobdzija et al. 1996] and array multiplier (ArrayM) are also simulated for comparison. To reduce the effect of the final addition, the same multi-bit adder in the tool library is utilized in all approximate multiplier designs as the final adder. Table VI shows the critical path delay, area, power, PDP and ADP of the considered multipliers. TruMA and TruMW are the truncated array and Wallace multipliers, respectively.

Figure 19 shows the comparison of delay, power and PDP of the equivalent approximate multipliers. The accurate array multiplier (ArrayM) is the slowest and the Wallace multiplier (WallaceM) consumes more area (as per Table VI); this is consistent with the theoretical analysis. Due to the expressively fast carry-ignored operation, AM1/TAM1, AM2/TAM2 have smaller delays compared to most of the other designs. BAM is significantly slower due to its array structure. AWTM, UDM, ICM and ACM have larger delays than the other approximate multipliers. BAM consumes a very low power, the power consumptions of AWTM and ACM are in the medium range, while UDM and ICM incur a relatively high power consumption. TruMA, TruMW and ETM have both a short delay and a low power dissipation.

A multiplier with a higher power dissipation usually has a larger area and thus larger PDP and ADP. In terms of power and area, TruMA, TruMW, ETM, TAM1/TAM2 and BAM are among the best designs. A common feature of these designs is that they all use truncation, which can significantly affect the *MRED* while the *NMED* may not be significantly changed. If most of the inputs have large values, the error introduced by truncation can be tolerated; thus truncation is a useful scheme to save area and power. Otherwise, truncation-based designs may yield unacceptably inaccurate results. Without truncation, a multiplier whose design is approximated in generating partial product (e.g. UDM) tends to have a large delay, power and area. These measures for the multipliers approximated in the partial product tree (e.g. AWTM) are moderate, while the multipliers using approximate counters or compressors (ICM, ACM, AM1, AM2) require higher power and area.

M. 14 . 1	Delay	Power	PDP	Area	ADP
Multiplier Type	( <i>ns</i> )	(uW)	(fJ)	$(um^2)$	$(um^2.ns)$
ArrayM	2.58	477.4	1,231.7	921	2,375.7
WallaceM	2.03	461.3	936.4	934	1,896.0
Multipliers with	Approxi	mation in	Generati	ing Partia	al Products
ŪDM	2.01	352.7	708.9	829	1666.7
Multipliers wit	h Appro	ximation	in the Pa	rtial Proc	duct Tree
BAM-16	2.34	221.3	517.8	441	1,031.9
<b>BAM-17</b>	2.17	189.5	411.2	384	833.3
BAM-18	1.99	161.1	320.6	331	658.7
BAM-20	1.65	111.0	183.2	237	390.4
ETM-7	1.57	140.5	220.6	349	547.6
ETM-8	1.50	108.5	162.8	288	431.4
AWTM-1	1.69	247.8	418.8	640	1.081.8
AWTM-2	1.69	259.4	438.4	665	1.123.7
AWTM-3	1.69	270.3	456.8	690	1,165.6
AWTM-4	1.74	280.0	478.2	715	1.243.2
Multipliers us	sing App	roximate	Counters	s or Com	pressors
ICM	1.87	367.4	687.0	937	1.751.4
ACM-3	1.97	279.5	550.6	738	1 454 5
ACM-4	2.00	284 1	568.2	724	1,447.0
AM1-13	1.38	355.4	490.5	819	1 128 8
AM1-16	$1.00 \\ 1.57$	380.6	597.5	878	1,120.0 1,378.5
AM2-10	1.01	336.8	434.5	816	1,052.6
$\Delta M_{2-13}$	1.20 1.53	364.2	557 9	919	1,002.0
$\Delta M_{2-16}$	1.00 1 71	400 4	684.7	1 051	1,400.1 1 797 9
TAM1-13	1.71	192.0	251.5	460	602.6
TAM1-16	1.01 1 45	214 G	311.9	516	748 2
TAM2-10	1.10	180.2	991.6	477	586 7
TAM2-10 $T\Delta M2-13$	1.20	219.2 219.5	221.0 314.5	581	864 3
TAM2-15 TAM9 16	1.40	212.0	314.5	603	1 199 7
TAM2-10	1.04	244.5	od Multin	liorg	1,122.1
$T_{m_1}M\Lambda A$	1 20	2 Onsigne 949 5	460 9	502	050 6
$T_{m_1}M\Lambda$	1.09	240.0 09.1	400.2 100.6	000 011	950.0 950.5
TrunA-0	1.19	94.1	109.0	561	200.0
Truivi vv-4	1.02	202.4	420.1	001 001	900.4
1 ruivi vv-0	1.10 	90.4 Loto Doot	100.4 h Multinl	209	203.0
	1 00			FOO	066.9
	1.00	204.3	403.7	020 407	900.2
	1.91	200.3	478.1	401	930.2
BWIII DM07	1.90	200.1	505.9	410	931.0
BMU7	2.03	270.4	548.9	528	1071.8
BM04	2.05	249.8	512.1	44'/	916.4
ABM2	2.25	208.0	468.0	424	954.0

Table VI. Circuit Characteristics of the Approximate Multipliers











(c)

Fig. 19. A comparison of delay, power and PDP of the approximate multipliers. (a) delay, (b) power and (c) power-delay product.



Fig. 20. *MRED* and PDP of the approximate unsigned multipliers. The parameter k for TruMA and TruMW is from 8 down to 2 from left to right; it is 21 down to 13 for BAM, and 10 to 16 for TAM1 and TAM2. The multipliers marked by circles are equivalent in terms of the number of accurate MSBs and are thus representatives of different designs.

In terms of PDP (Figure 19(c)) and ADP, TruMA, TruMW, ETM, TAM1 and TAM2 have very small values, while ICM, UDM and AM2 are the opposite. The values of PDP and ADP for AM1, ACM, BAM and AWTM are in the medium range.

For the approximate Booth multipliers, PEBM is the fastest, but it is very power and area consuming due to the use of a carry save adder tree for the parallel accumulation. ABM2 is the slowest but the most power and area efficient design due to the smaller number of partial products generated by the time-consuming recoding adder in the radix-8 algorithm. Thus, it has the smallest PDP and a moderate value of ADP. With no error compensation, BBM shows a small delay, low power dissipation and small circuit area, and thus smaller PDP and ADP compared with most of the other designs (except ABM2 and BM04). BM11 and BM04 have similar values for all circuit metrics. BM07 has a similar delay, but with a higher power and area and thus a larger PDP and ADP.

3.2.3. Discussion. *MRED* and PDP are jointly considered next for an overall evaluation of the approximate multipliers, as shown in Figure 20 and Figure 21.

Figure 20 shows that TruMW has a smaller PDP than TruMA when the same number of LSBs is truncated. Among the truncation-based designs, the truncated unsigned multipliers (TruMA and TruMW) are slightly more accurate (in MRED) than BAM and ETM for a similar PDP. TruMW has a smaller MRED than most other approximate designs (except TAM1, TAM2 and ICM).

In Figure 20, TAM1-13, TAM1-16, TAM2-13, TruMA-6, TruMW-6 and BAM-18 have both small PDPs and *MREDs*. Most of the other designs have at least one major shortcoming. ICM and ACM incur a very low error, but their PDPs are very high. Other than the truncated designs, ETM-8 has the smallest PDP but with a rather large *MRED*. UDM shows a poor performance in both PDP and *MRED*. Even though some BAM



Fig. 21. *MRED* and PDP of the approximate Booth multipliers.

configurations have small PDPs, their delays are generally large (Figure 19(a)); moreover, some BAM configurations have low accuracies. AWTMs have large PDPs and only AWTM-4 has a high accuracy.

As for the approximate Booth multipliers (Figure 21), ABM2 shows the lowest PDP and a moderate accuracy. BM11 and BM07 are very accurate in terms of MRED but with relatively poor PDPs. PEBM shows both a moderate PDP and MRED.

## 4. APPROXIMATE DIVIDERS

#### 4.1. Classification

The divider is a less frequently used arithmetic module compared to the adder and multiplier; therefore, less research has been pursued on an approximate design.

Two methodologies have been advocated for sequential division, the digit recurrent algorithm [Liu and Nannarelli 2012] and the functional iterative algorithm (e.g., using the Newton-Raphson algorithm [Flynn 1970]). A sequential divider has a low hardware complexity, however its delay is considerably longer than an adder and a multiplier, so it significantly affects the overall performance of a processor. Thus, dividers made of combinational logic circuits are discussed in this article. Like multiplication, division can also be implemented by an array structure, in which adder cells are replaced by subtractor cells. Several approximations are made on the array divider while retaining a low-power [Chen et al. 2015; 2016; Chen et al. 2017] and high-speed [Hashemi et al. 2016] operation. In addition, different approximate divider designs based on rounding [Zendegani et al. 2016] and curve fitting [Low and Jong 2013; Wu and Jong 2015] are also proposed.

4.1.1. Approximate array dividers. Four types of approximate unsigned integer nonrestoring divider (AXDnr) are presented in [Chen et al. 2015]. Three approximate subtractors (AXCSs) are designed for the array of an unsigned divider by simplifying the circuit of an exact subtractor cell. The AXCSs are then used to replace the exact subtractor cells at the least significant vertical, horizontal, square or triangle cells of the array divider. Moreover, a truncation scheme is utilized by discarding the approximate subtractors for comparison. Based on the same theory and design, four types of approximate restoring dividers (AXDr) are further proposed by using the proposed AXCSs [Chen et al. 2016]. It has been shown that AXDrs are slightly more accurate and consumes lower power than AXDnrs.

To make the remaining subtractor cells more efficient, a dynamic approximate divider (DAXD) is proposed by dynamically selecting the inputs of the subtractor cells [Hashemi et al. 2016]. DAXD selects a fixed number of bits in the input operands from the most significant non-zero bit, and then truncates the least significant bits. Therefore, it can be implemented by two leading-one detectors, two multiplexers, a smaller array divider and a barrel shifter.

To further improve the performance and power efficiency, an approximate signeddigit adder is proposed for high-radix dividers [Chen et al. 2017]. Compared to the conventional radix-2 design, the approximate radix-4 and radix-8 dividers show a higher speed and consume lower power, albeit with a slightly lower accuracy.

4.1.2. Curve fitting based approximate dividers. A widely used methodology to reduce the hardware overhead of a divider is based on binary logarithms, that is to obtain the antilogarithmic value of the difference between the logarithmic values of the dividend and the divisor . Mitchell et al. first developed the logarithm approximation of a binary number by shifting and counting [Mitchell Jr 1962]. Inspired by it, a new antilogarithmic algorithm was proposed using a piecewise linear approximation [Low and Jong 2013]. A high-speed divider (HSD) based on this algorithm was then presented. The antilogarithmic algorithm is directly approximated from the input operands, thus only look-up tables and multiplications are required, i.e., no logarithmic or subtraction operation is needed. HSD achieves a better accuracy and a much higher speed (but at a larger area) than the divider implemented directly by Mitchell's algorithm.

A similar curve fitting approach was used for the design of a floating-point divider (FPD) [Wu and Jong 2015]. FPD partitions the curved surfaces of the quotient into several square or triangular regions and linearly approximates each region by curve fitting. Finally, the division of the mantissas is implemented by a comparison module, a look-up table, shifters and adders. This approximate divider achieves a better accuracy than that in [Low and Jong 2013] with similar circuit characteristics.

4.1.3. Rounding based approximate dividers. A high-speed and energy-efficient roundingbased approximate divider (SEERAD) is presented in [Zendegani et al. 2016]. It transforms the division to a smaller multiplication by rounding the divisor B to a form of  $2^{K+L}/D$ , where K shows the bit position of the most significant "1" of B ( $K = \lfloor log_2 B \rfloor$ ), and L and D are constant integers found from an exhaustive simulation by the condition of obtaining the lowest mean relative error. Different accuracy levels are considered to improve the accuracy of SEERAD by varying D and L with combinations of the more relevant bits of B after the most significant "1". The multiplier in SEERAD is implemented by several shift units and an adder block. Thus, SEERAD is very fast.

## 4.2. Discussion

The approximate array dividers, AXDnr, AXDr and DAXD, are designed for the unsigned  $n/(\frac{n}{2})$  division, where n is the width of the dividend. As the borrow signal passes through all subtractor cells, the critical path of an array divider is  $O(n^2)$  [Parhami 2000]. Thus, the speed of the approximate array dividers is not very fast, but the area and power dissipation are relatively low among the approximate designs. The accuracy of the array dividers depends on the number of replaced cells (AXDnr and AXDr) and the size of the sub-divider (DAXD).

The approximate dividers based on curve fitting are very accurate, e.g, the maximum relative error distance of a 16/16 HSD and FPD are 0.20% and 0.14%, respectively. Although the curve fitting using software is complex, the hardware implementation consisting of look-up tables, smaller multipliers and adders is very simple. Thus, curve

fitting based dividers are usually fast. However, the look-up tables used for storing the constant coefficients make the approximate dividers area and power consuming.

The rounding based approximate divider has a maximum relative error distance of 6.25% (for the design with the highest accuracy level) in a 16/16 divider and hence, it is not very accurate. Its hardware implementation is simpler than those of the curve fitting based dividers. Therefore, the rounding based approximate divider has a very high speed and relatively large area and power dissipation due to the use of look-up tables.

## 5. IMAGE PROCESSING APPLICATION

Low power dissipation and small circuit area are basic requirements for consumer electronic products, especially for mobile devices with stringent battery restrictions. Therefore, approximate designs have been pervasively considered for implementations in image processing. Approximate multipliers have been utilized for image sharpening [Jiang et al. 2016b]. In this article, both the adder and multiplier in the image sharpening algorithm are replaced by approximate designs. Moreover, approximate dividers are used to detect the difference between two images to show the changes. This application is known as change detection.

## 5.1. Image Sharpening

The image sharpening algorithm computes R(x, y) = 2I(x, y) - S(x, y) [Lau et al. 2009], where *I* is the input image, *R* is the sharpened image, and *S* is given by

$$S(x,y) = \frac{1}{4368} \sum_{m=-2}^{2} \sum_{n=-2}^{2} G(m+3,n+3)I(x-m,y-n),$$
(1)

where G is a  $5 \times 5$  matrix given by

$$G = \begin{bmatrix} 16 & 64 & 112 & 64 & 16\\ 64 & 256 & 416 & 256 & 64\\ 112 & 416 & 656 & 416 & 112\\ 64 & 256 & 416 & 256 & 64\\ 16 & 64 & 112 & 64 & 16 \end{bmatrix}.$$
 (2)

Simulation results in [Jiang et al. 2016b] show that AM2-15, AM1-15, TAM2-16, TAM1-16, BAM-16, AM2-13, AM1-13, ACM-4, ACM-3, TAM2-13, TAM1-13, BAM-17, AWTM-4 and BAM-18 achieve visually acceptable image sharpening results. Among these multipliers, the ones with a moderate hardware overhead (AM1-13, TAM2-13, TAM1-16, TAM1-13, BAM-17 and BAM-18) are selected in this article for image sharpening. Likewise, the approximate adders LOA, CSA, ETAII and CSPA are selected. As the multiplication result of a  $16 \times 16$  multiplier is 32-bit wide, 32-bit approximate adders are used for image sharpening. The value of parameter k is 8 for CSA, ETAII and CSPA, and 16 for LOA.

The results for image sharpening using the selected approximate multipliers and adders are given in Table VII, while the accurate result is shown in Figure 22. The images sharpened using CSPA have unacceptable defects and some defects (white dots) can be seen in the image sharpened by AM1-13 and ETAII-8 when zooming into the images in Table VII. Other images processed by the approximate designs show similar quality with the accurate result. This is also confirmed by the peak signal-to-noise ratio (PSNR), as shown in Table VIII. The PSNRs of the images sharpened by a truncation based multiplier (i.e., TAM1-16, TAM2-13, TAM1-13, BAM-17 or BAM-18) are fixed as the adder is changed among LOA-16, CSA-8 and ETAII-8. This occurs because the lower 16 bits of the multiplication results generated by these multipliers are zeros and



Fig. 22. The image sharpened using an accurate multiplier and an accurate adder.

hence, only the higher half of an approximate adder (as an accurate 16-bit adder for LOA-16, CSA-8 or ETAII-8) is used.

The image sharpening algorithm is implemented in VHDL by using the selected approximate adders and multipliers. In the implementation, no pipelining or memory unit is used for exclusively showing the hardware characteristics of the approximate arithmetic circuits. It is synthesized by Synopsys DC using the same process, voltage and temperature as in the simulation of the approximate adders. The  $512 \times 512$  pixel values of the image shown in Table VII are used as inputs for assessing the power dissipation using the PrimeTime-PX tool at a clock period of 10 *ns*.

Table IX shows the circuit characteristics of image sharpening using approximate multipliers and adders. Using the same multiplier, the image sharpening implementations with LOA-16 and ETAII-8 show similar values for the delay, power and area (except for AM1-13), while the implementations using CSA-8 have relatively larger values for these metrics. Likewise, the image sharpening circuits have similar characteristics using the same adder except that AM1-13, BAM-17 and BAM-18 based schemes show slightly larger values.

Compared with the image sharpening circuit using accurate multipliers and adders, the approximate designs using CSA-8 or AM1-13 achieve small improvement in terms of delay and area because CSA and AM1 are less efficient in delay and area compared with the other approximate designs. By using LOA-16, ETAII-8, TAM2-13, BAM-17 or BAM-18, the circuit can be 23% faster and saves as much as 53% in power, 58% in area, 64% in PDP and 62% in ADP compared to the accurate design.

## 5.2. Change Detection

In the application of change detection, the ratio between two corresponding pixel values is calculated by a divider [Chen et al. 2016]. The changes in two images are then highlighted by normalizing the pixel ratios. In this section, 16/8 divider designs are used to calculate the division of two 8-bit gray-level images as shown in Figure 23(a) and (b). To ensure a higher accuracy, the pixel values of the first image are multiplied by 64. As HSD and FPD are designed for floating-point division, AXDr, DAXD and SEERAD are selected for the change detection. Among the four types of AXDr, the triangle replacement has been shown to achieve the best results for image processing [Chen et al. 2016]. Therefore, three designs of AXDr with the triangle replacement of depth 8, AXDr1 (using approximate subtractor 1), AXDr2 (using approximate subtractor 2) and AXDr3 (using approximate subtractor 3), are used for the change detection. For DAXD, 8/4 and 10/5 accurate dividers are utilized in DAXD8 and DAXD10, respectively. Four accuracy levels are considered in SEERAD; they are referred to as SEERAD1, SEERAD2, SEERAD3 and SEERAD4. Moreover, the accurate array divider (denoted as ArrayD) is simulated for comparison.

Figure 23 shows the change detection results by the dividers and the obtained P-SNR value is shown in the parentheses. It is clear that AXDr1, AXDr3 and SEERAD4

Approximate design	LOA-16	CSA-8	ETAII-8	CSPA-8
TAM1-16				
AM1-13				
TAM2-13				
TAM1-13				
BAM-17				
BAM-18				

Table VII. Images Sharpened using Different Approximate Adder and Multiplier Pairs

perform very well in the application of change detection, while the results by the other designs are of lower quality.

The characteristics of the change detection circuits are obtained by using the same tool and process as in the simulation of approximate adders. The clock period is 6 ns,

Adder Multiplier	LOA-16	CSA-8	ETAII-8	CSPA-8
TAM1-16	46.97	46.97	46.97	25.01
AM1-13	45.21	45.06	36.86	24.20
TAM2-13	41.87	41.87	41.87	24.32
TAM1-13	41.42	41.42	41.42	24.35
BAM-17	40.09	40.09	40.09	25.19
BAM-18	33.99	33.99	33.99	24.21

Table VIII. PSNRs of the Sharpened Images (dB)

Table IX. Delay, Power, and Area of Image Sharpening using Approximate Multipliers and Adders

Multiplier	Adder	Delay (ns)	Power (mW)	PDP (pJ)	<b>Area</b> ( <i>um</i> <sup>2</sup> )	$\begin{array}{c} \textbf{ADP} \\ (um^2.ns) \end{array}$
ArrayM	CLAG	6.74	1.995	13.45	31,183.9	210,179.5
TAM1-16	LOA-16	5.36	0.9723	5.21	18,139.0	97,225.0
TAM1-16	CSA-8	7.45	1.032	7.69	$23,\!652.1$	176,208
TAM1-16	ETAII-8	5.34	0.9643	5.15	$18,\!056.8$	96,423.3
AM1-13	LOA-16	5.41	1.193	6.45	26,644.0	144,144
AM1-13	CSA-8	7.41	1.377	10.20	$30,\!586.5$	226,646
AM1-13	ETAII-8	6.40	1.369	8.76	$28,\!214.7$	180,574
TAM2-13	LOA-16	5.25	1.055	5.54	17,057.8	89,553.5
TAM2-13	CSA-8	6.43	1.053	6.77	$20,\!526.6$	131,986
TAM2-13	ETAII-8	5.22	1.041	5.43	$16,\!975.6$	88,612.6
TAM1-13	LOA-16	5.25	0.9467	4.97	$17,\!221.0$	90,410.3
TAM1-13	CSA-8	7.45	0.9942	7.41	22,734.1	169,369
TAM1-13	ETAII-8	5.34	0.9350	4.88	$17,\!138.8$	89,464.5
<b>BAM-17</b>	LOA-16	6.14	1.226	7.53	14,993.8	92,061.9
BAM-17	CSA-8	7.36	1.247	9.17	$16,\!533.0$	121,683
BAM-17	ETAII-8	6.13	1.211	7.42	$14,\!868.5$	91,143.9
BAM-18	LOA-16	5.97	1.097	6.55	$13,\!285.3$	79,313.2
BAM-18	CSA-8	6.89	1.117	7.70	16,901.8	116,453
BAM-18	ETAII-8	5.96	1.076	6.41	$13,\!156.0$	78,409.8

ACM Journal on Emerging Technologies in Computing Systems, Vol. 13, No. 4, Article 60, Pub. date: July 2017.



Fig. 23. Change detection using different approximate dividers.

and the input combinations for the power evaluation are the two images shown in Figure 23(a) and (b) with  $384 \times 507$  pixels. The synthesis results are shown in Table X. To be consistent with the other designs, AXDr1, AXDr2 and AXDr3 are implemented at the gate level rather than at the transistor level as in [Chen et al. 2016].

Table X shows that the array-based dividers (ArrayD, AXDrs and DAXD) are more power and area efficient than the rounding based approximate dividers (SEERADs). However, they are very slow except for DAXD that uses a smaller accurate array divider. SEERADs are very fast, but they consume more power and area due to the use of look-up tables and the large output size (including the fractional part). In terms of

PDP and ADP, DAXD10 has the smallest values followed by SEERAD3. Among the approximate dividers that produce excellent change detection results, SEERAD4 has the smallest PDP but the largest ADP. AXDr1 results in the largest PDP and the smallest ADP, while AXDr3 has both moderate PDP and ADP. To be more specific, the change detection using AXDr3 saves 25% of power and 12% of area compared with the accurate design. It is 40% faster with 18% reduction in PDP by using SEERAD4.

Divider	PSNR ( $dB$ )	Delay (ns)	Power (uW)	PDP (fJ)	<b>Area</b> ( <i>um</i> <sup>2</sup> )	ADP (um <sup>2</sup> .ns)
ArrayD	-	4.08	54.29	221.50	425.8	1,737.2
AXDr1	40.46	3.85	50.71	195.23	415.5	$1,\!599.7$
AXDr2	22.91	4.36	54.08	235.79	408.2	1,779.6
AXDr3	41.71	4.58	40.55	185.72	376.2	1,722.9
DAXD10	24.33	2.43	40.25	97.84	375.7	912.9
SEERAD3	33.17	1.83	60.27	110.29	615.8	1,126.8
SEERAD4	36.61	2.43	70.62	181.33	765.4	1,859.9

Table X. Accuracy and Circuit Characteristics of the Change Detection using Approximate Dividers.

## 6. CONCLUSION

In this article, designs of approximate arithmetic circuits are reviewed. Their error and circuit characteristics are evaluated using functional simulation and hardware synthesis with an industrial 28nm technology library.

Approximate Adders: In general, approximate speculative adders show high accuracy and relatively small PDPs. The approximate adders using approximate full adders in the LSBs are slow, but they are power efficient with high ERs (due to the approximate LSBs), low average error and moderate NMED and MRED values (due to the accurate MSBs). The error and circuit characteristics of the segmented and carry select adders vary with the prediction of the carry signals.

A truncated adder has a smaller MRED (an indicator of a smaller error magnitude) than most approximate designs at a similar PDP except for LOA and CSA. However, it has a lower performance and a significantly higher ER compared with the other approximate designs. As a result, a simple truncation of the LSBs in an adder causes a high ER and does not significantly improve the performance of the adder, though with a relatively small error distance.

Approximate Multipliers: For approximate multipliers, truncation of part of the partial products is an effective scheme to reduce hardware complexity, while preserving a moderate NMED and MRED. Similarly, truncating some LSBs of the input operands can efficiently reduce the hardware overhead of a multiplier and result in a moderate MRED (an indicator of the error magnitude) that is smaller than most other approximate designs, except for TAM1, TAM2 and ICM, for a similar PDP.

Albeit with a relatively low ER, UDM shows a low accuracy in terms of the error distance and a relatively high circuit overhead, because the  $2 \times 2$  approximate multi-

plier is used to compute the most significant bits and accurate adders are utilized to accumulate the generated partial products. ICM has the lowest ER among all designs. When truncation is not used, multipliers approximated in the partial product tree tend to have a poor accuracy (except AWTM-3 and AWTM-4) and moderate hardware consumption, while multipliers using approximate counters or compressors are usually very accurate with relatively high power dissipation and hardware consumption. The approximate Booth multipliers show different characteristics in hardware efficiency and accuracy.

Approximate Dividers: For the dividers, the approximate array dividers are slow, but they are hardware efficient with variable accuracy depending on the approximation parameters. The dividers based on curve fitting are very accurate and fast but they require a large area and high power dissipation due to the utilization of look-up tables. The rounding based approximate dividers have a very high speed, large area and power dissipation, with a relatively low accuracy.

*Application:* Image sharpening is implemented using the selected approximate multipliers and adders. The accuracy and circuit characteristics of the image sharpening obtained by simulation indicate significant savings in hardware, while producing similar sharpening quality as the accurate design. On average, the designs using approximate adders and multipliers with an acceptable accuracy consumes only 55% of power, 62% of area, 51% of PDP and 57% of ADP compared to the accurate design.

Approximate dividers are utilized in the change detection of images. The simulation and synthesis results show that the change detection circuits using the approximate array divider (AXDr1 and AXDr3) are power and area efficient but very slow, whereas the one using the rounding based approximate divider (SEERAD4) consumes more power and area with a high performance for an excellent detection accuracy.

#### ACKNOWLEDGMENT

The authors would like to thank Vincent Camus from EPFL for his contributions in proof reading the article.

## REFERENCES

Tinku Acharya and Ajoy K Ray. 2005. Image processing: principles and applications. John Wiley & Sons.

- Shaahin Angizi, Zhezhi He, Ronald F. DeMara, and Deliang Fan. 2017. Composite Spintronic Accuracy-Configurable Adder for Low Power Digital Signal Processingr. In *ISQED*. IEEE.
- Dursun Baran, Mustafa Aktan, and Vojin G Oklobdzija. 2010. Energy efficient implementation of parallel CMOS multipliers with improved compressors. In ACM/IEEE international symposium on Low power electronics and design. ACM, 147–152.
- Kartikeya Bhardwaj, Pravin S. Mane, and Jorg Henkel. 2014. Power- and area-efficient Approximate Wallace Tree Multiplier for error-resilient systems. In *ISQED*. 263–269.
- Hao Cai, You Wang, Lirida AB Naviner, Zhaohao Wang, and Weisheng Zhao. 2016. Approximate computing in MOS/spintronic non-volatile full-adder. In International Symposium on Nanoscale Architectures (NANOARCH). IEEE, 203–208.
- Vincent Camus, Jeremy Schlachter, and Christian Enz. 2015. Energy-efficient inexact speculative adder with high performance and accuracy control. In *IEEE International Symposium on Circuits and Systems* (ISCAS). 45–48.
- Vincent Camus, Jeremy Schlachter, and Christian Enz. 2016. A low-power carry cut-back approximate adder with fixed-point implementation and floating-point precision. In *Proceedings of the 53rd Annual Design Automation Conference*. ACM, 127.
- Linbin Chen, Jie Han, Weiqiang Liu, and Fabrizio Lombardi. 2015. Design of Approximate Unsigned Integer Non-restoring Divider for Inexact Computing. In Proceedings of the 25th edition on Great Lakes Symposium on VLSI. ACM, 51–56.
- Linbin Chen, Jie Han, Weiqiang Liu, and Fabrizio Lombardi. 2016. On the Design of Approximate Restoring Dividers for Error-Tolerant Applications. *IEEE Trans. Comput.* 65, 8 (2016), 2522–2533.

- Linbin Chen, Paolo Montuschi, Jie Han, Weiqiang Liu, and Fabrizio Lombardi. 2017. Design of Approximate High-Radix Dividers by Inexact Binary Signed-Digit Addition. In Proceedings of the 27th IEEE/ACM Great Lakes Symposium on VLSI.
- Yuan-Ho Chen and Tsin-Yuan Chang. 2012. A high-accuracy adaptive conditional-probability estimator for fixed-width Booth multipliers. *IEEE Trans. Circuits and Systems I: Regular Papers* 59, 3 (2012), 594– 603.
- Vinay K. Chippa, Debabrata Mohapatra, Anand Raghunathan, Kaushik Roy, and Srimat T. Chakradhar. 2010. Scalable effort hardware design: exploiting algorithmic resilience for energy efficiency. In DAC. ACM Press, New York, New York, USA, 555–560.
- Kyung-Ju Cho, Kwang-Chul Lee, Jin-Gyun Chung, and Keshab K Parhi. 2004. Design of low-error fixedwidth modified booth multiplier. *IEEE Transactions on VLSI Systems* 12, 5 (2004), 522–531.
- S.R. Datla, M.A Thornton, and D.W. Matula. 2009. A Low Power High Performance Radix-4 Approximate Squaring Circuit. In ASAP. 91–97.
- Kai Du, P. Varman, and K. Mohanram. 2012. High performance reliable variable latency carry select addition. In Design, Automation & Test in Europe Conference & Exhibition (DATE), 2012. 1257–1262.
- Hadi Esmaeilzadeh, Adrian Sampson, Luis Ceze, and Doug Burger. 2012. Architecture support for disciplined approximate programming. In ACM SIGPLAN Notices, Vol. 47. ACM, 301–312.
- Farzad Farshchi, Muhammad Saeed Abrishami, and Sied Mehdi Fakhraie. 2013. New approximate multiplier for low power digital signal processing. In CADS. IEEE, 25–30.
- Michael J Flynn. 1970. On division by functional iteration. IEEE Trans. Comput. 100, 8 (1970), 702-706.
- V. Gupta, D. Mohapatra, A. Raghunathan, and K. Roy. 2013. Low-Power Digital Signal Processing Using Approximate Adders. *IEEE Trans. CAD* 32, 1 (January 2013), 124–137.
- Jie Han. 2016. Introduction to Approximate Computing. In 2016 IEEE 34th VLSI Test Symposium (VTS). IEEE, 1–1.
- Jie Han and Michael Orshansky. 2013. Approximate computing: An emerging paradigm for energy-efficient design. In Test Symposium (ETS), 2013 18th IEEE European. IEEE, 1–6.
- Soheil Hashemi, R Bahar, and Sherief Reda. 2015. Drum: A dynamic range unbiased multiplier for approximate applications. In Proceedings of the IEEE/ACM International Conference on Computer-Aided Design. IEEE Press, 418–425.
- Soheil Hashemi, R Bahar, and Sherief Reda. 2016. A low-power dynamic divider for approximate applications. In *Proceedings of the 53rd Annual Design Automation Conference*. ACM, 105.
- Junjun Hu and Weikang Qian. 2015. A New Approximate Adder with Low Relative Error and Correct Sign Calculation. In *DATE*. 1449–1454.
- Jiawei Huang, John Lach, and Gabriel Robins. 2012. A methodology for energy-quality tradeoff using imprecise hardware. In Proceedings of the 49th ACM Annual Design Automation Conference. 504–509.
- Honglan Jiang, Jie Han, and Fabrizio Lombardi. 2015. A Comparative Review and Evaluation of Approximate Adders. In Proceedings of ACM Great Lakes Symposium on VLSI. 343–348.
- Honglan Jiang, Jie Han, and Fabrizio Lombardi. 2016a. Approximate Radix-8 Booth Multiplier for Low-Power and High-Performance Operation. *IEEE Trans. Comput.* 65, 8 (2016), 2638–2644.
- Honglan Jiang, Jie Han, and Fabrizio Lombardi. 2016b. A Comparative Evaluation of Approximate Multipliers. In IEEE / ACM International Symposium on Nanoscale Architectures.
- Andrew B Kahng and Seokhyeong Kang. 2012. Accuracy-configurable adder for approximate arithmetic designs. In Proceedings of the 49th ACM Annual Design Automation Conference. 820–825.
- D Kelly, B Phillips, and S Al-Sarawi. 2009. Approximate signed binary integer multipliers for arithmetic data value speculation. In *Conference on Design & Architectures For Signal And Image Processing*. Sophia Antipolis, France.
- Yongtae Kim, Yong Zhang, and Peng Li. 2013. An Energy Efficient Approximate Adder with Carry Skip for Error Resilient Neuromorphic VLSI Systems. In *ICCAD*. 130–137.
- Parag Kulkarni, Puneet Gupta, and Milos Ercegovac. 2011. Trading accuracy for power with an underdesigned multiplier architecture. In *International Conference on VLSI Design*. 346–351.
- Khaing Yin Kyaw, Wang Ling Goh, and Kiat Seng Yeo. 2010. Low-power high-speed multiplier for errortolerant application. In *EDSSC*. 1–4.
- Mark SK Lau, Keck-Voon Ling, and Yun-Chung Chu. 2009. Energy-aware probabilistic multiplier: design and analysis. In CASES. 281–290.
- Li Li and Hai Zhou. 2014. On Error Modeling and Analysis of Approximate Adders. In ICCAD. 511-518.
- Jinghang Liang, Jie Han, and F. Lombardi. 2013. New Metrics for the Reliability of Approximate and Probabilistic Adders. *IEEE Trans. Comput.* 62, 9 (September 2013), 1760–1771.

- Chia-Hao Lin and Ing-Chao Lin. 2013. High accuracy approximate multiplier with error correction. In *ICCD*. IEEE, 33–38.
- IngChao Lin, YiMing Yang, and ChengChian Lin. 2015. High-Performance Low-Power Carry Speculative Addition With Varible Latency. *IEEE Trans. VLSI Syst.* 23, 9 (2015), 1591–1603.
- Cong Liu. 2014. Design and Analysis of Approximate Adders and Multipliers. Master's thesis. University of Alberta, Canada.
- Cong Liu, Jie Han, and Fabrizio Lombardi. 2014. A Low-Power, High-Performance Approximate Multiplier with Configurable Partial Error Recovery. In *DATE*.
- Cong Liu, Jie Han, and Fabrizio Lombardi. 2015. An analytical framework for evaluating the error characteristics of approximate adders. *IEEE Trans. Comput.* 64, 5 (2015), 1268–1281.
- Cong Liu, Honglan Jiang, Fabrizio Lombardi, and Jie Han. 2017a. High-Performance Approximate Unsigned Multipliers with Configurable Error Recovery. *IEEE Transactions on Circuits and Systems I, under revision.* (2017).
- Wei Liu and Alberto Nannarelli. 2012. Power efficient division and square root unit. IEEE Trans. Comput. 61, 8 (2012), 1059–1070.
- Weiqiang Liu, Liangyu Qian, Chenghua Wang, Honglan Jiang, Jie Han, and Fabrizio Lombardi. 2017b. Design of Approximate Radix-4 Booth Multipliers for Error-Tolerant Computing. *IEEE Trans. Comput.* (2017).
- Joshua Yung Lih Low and Ching Chuen Jong. 2013. Non-iterative high speed division computation based on Mitchell logarithmic method. In *IEEE International Symposium on Circuits and Systems (ISCAS)*. 2219–2222.
- Shih-Lien Lu. 2004. Speeding up processing with approximation circuits. *Computer* 37, 3 (March 2004), 67–73.
- Jieming Ma, Ka Lok Man, Nan Zhang, Sheng-Uei Guan, and Taikyeong Ted Jeong. 2013. High-speed areaefficient and power-aware multiplier design using approximate compressors along with bottom-up tree topology. In *ICMV: Algorithms, Pattern Recognition, and Basic Technologies*.
- H R Mahdiani, A Ahmadi, S M Fakhraie, and C Lucas. 2010. Bio-Inspired Imprecise Computational Blocks for Efficient VLSI Implementation of Soft-Computing Applications. *IEEE Trans. Circuits and Systems* 57, 4 (April 2010), 850–862.
- Sana Mazahir, Osman Hasan, Rehan Hafiz, Muhammad Shafique, and Jörg Henkel. 2017. Probabilistic error modeling for approximate adders. *IEEE Trans. Comput.* 66, 3 (2017), 515–530.
- Jin Miao, Ku He, Andreas Gerstlauer, and Michael Orshansky. 2012. Modeling and synthesis of qualityenergy optimal approximate adders. In Proceedings of the ACM International Conference on Computer-Aided Design. 728–735.
- Joshua San Miguel, Jorge Albericio, Andreas Moshovos, and Natalie Enright Jerger. 2015. Doppelgänger: a cache for approximate computing. In *Proceedings of the 48th International Symposium on Microarchitecture*. ACM, 50–61.
- SONG Min-An, VAN Lan-Da, and KUO Sy-Yen. 2007. Adaptive low-error fixed-width Booth multipliers. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences* 90, 6 (2007), 1180–1187.
- John N Mitchell Jr. 1962. Computer multiplication and division using binary logarithms. *IRE Transactions* on Electronic Computers 4 (1962), 512–517.
- D. Mohapatra, V.K. Chippa, A Raghunathan, and K. Roy. 2011. Design of voltage-scalable meta-functions for approximate computing. In *Design, Automation & Test in Europe Conference & Exhibition (DATE)*. 1–6.
- Amir Momeni, Jie Han, Paolo Montuschi, and Fabrizio Lombardi. 2015. Design and Analysis of Approximate Compressors for Multiplication. *IEEE Trans. Comput.* 64, 4 (2015), 984–994.
- Vojtech Mrazek, Syed Shakib Sarwar, Lukas Sekanina, Zdenek Vasicek, and Kaushik Roy. 2016. Design of power-efficient approximate multipliers for approximate artificial neural networks. In *International Conference On Computer Aided Design (ICCAD)*. 7.
- Srinivasan Narayanamoorthy, Hadi Asghari Moghaddam, Zhenhong Liu, Taejoon Park, and Nam Sung Kim. 2015. Energy-efficient approximate multiplication for digital signal processing and classification applications. *IEEE Transactions on VLSI Systems* 23, 6 (2015), 1180–1184.
- Vojin G. Oklobdzija, David Villeger, and Simon S. Liu. 1996. A method for speed optimized partial product reduction and generation of fast parallel multipliers using an algorithmic approach. *IEEE Trans. Comput.* 45, 3 (1996), 294–306.

Behrooz Parhami. 2000. Computer arithmetic. Oxford university press.

Adrian Sampson, Werner Dietl, Emily Fortuna, Danushen Gnanapragasam, Luis Ceze, and Dan Grossman. 2011. EnerJ: Approximate data types for safe and general low-power computation. 46, 6 (June 2011), 164–174.

Doochul Shin. 2010. Approximate logic synthesis for error tolerant applications. In DATE. IEEE, 957-960.

- Zdenek Vasicek and Lukas Sekanina. 2015. Evolutionary approach to approximate digital circuits design. IEEE Transactions on Evolutionary Computation 19, 3 (2015), 432–444.
- Rangharajan Venkatesan, Amit Agarwal, Kaushik Roy, and Anand Raghunathan. 2010. MACACO: Modeling and analysis of circuits for approximate computing. In *ICCAD*. 667–673.
- Ajay K Verma, Philip Brisk, and Paolo Ienne. 2008. Variable latency speculative addition: A new paradigm for arithmetic circuit design. In *DATE*. 1250–1255.
- Jiun-Ping Wang, Shiann-Rong Kuang, and Shish-Chang Liang. 2011. High-accuracy fixed-width modified Booth multipliers for lossy applications. *IEEE Transactions on VLSI Systems* 19, 1 (2011), 52–60.
- Lei Wu and Ching Chuen Jong. 2015. A curve fitting approach for non-iterative divider design with accuracy and performance trade-off. In *IEEE 13th International New Circuits and Systems Conference* (NEWCAS). 1–4.
- Xinghua Yang, Yue Xing, Fei Qiao, Qi Wei, and Huazhong Yang. 2016. Approximate Adder with Hybrid Prediction and Error Compensation Technique. In IEEE Computer Society Annual Symposium on VLSI (ISVLSI). IEEE, 373–378.
- Zhixi Yang, Ajaypat Jain, Jinghang Liang, Jie Han, and Fabrizio Lombardi. 2013. Approximate XOR/XNORbased Adders for Inexact Computing. In *IEEE-NANO*. 690–693.
- Rong Ye, Ting Wang, Feng Yuan, Rakesh Kumar, and Qiang Xu. 2013. On Reconfiguration-Oriented Approximate Adder Design and Its Application. In *ICCAD*. 48–54.
- Reza Zendegani, Mehdi Kamal, Milad Bahadori, Ali Afzali-Kusha, and Massoud Pedram. 2017. RoBA multiplier: A rounding-based approximate multiplier for high-speed yet energy-efficient digital signal processing. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* 25, 2 (2017), 393–401.
- Reza Zendegani, Mehdi Kamal, Arash Fayyazi, Ali Afzali-Kusha, Saeed Safari, and Massoud Pedram. 2016. SEERAD: A high speed yet energy-efficient rounding-based approximate divider. In *Design, Automation* & Test in Europe Conference & Exhibition (DATE). IEEE, 1481–1484.
- Ning Zhu, Wang Ling Goh, and Kiat Seng Yeo. 2009. An enhanced low-power high-speed adder for errortolerant application. In *ISIC 2009.* 69–72.