

Imprecise Computing in Datapath Design: An Overclocking Approach

Kan Shi, David Boland, *Member, IEEE*, and George A. Constantinides, *Senior Member, IEEE*

Abstract—As process scaling introduces significant performance variations, releasing the tight accuracy requirement rather than covering all possible worst cases would potentially offer the greater freedom to create a design with better performance or energy efficiency. In this paper, we compare two different approaches that could trade accuracy for performance. One is the traditional approach where the precision used in the datapath is limited to meet a target latency. The other is a proposed new approach to boost the clock frequency by simply allowing the datapath to operate without timing closure. We demonstrate analytically and experimentally that on average, our approach obtains either smaller errors or faster operating frequencies in comparison to the traditional approach, since the worst case caused by timing violations only happens rarely, while the precision loss results in errors to all data. We also show that for embedded applications where silicon area is also a limited factor, using the proposed approach on simple arithmetic primitives could achieve better accuracy or performance than using the traditional method on advanced operators.

Index Terms—Overclocking, Imprecise Computing, Numerical Analysis.

I. INTRODUCTION

CIRCUIT performance has increased tremendously over the past decades, with the continuous scaling of CMOS technology to the nanometer regime. However, the drastic variations introduced by higher integration densities is anticipated to be the major obstacle when designing reliable, high performance circuits. Although hardware designers traditionally employ safety margins to ensure a uniform functionality across a variety of possible working environments, there are two major problems if we keep using this approach nowadays. One is that the uniform performance should no longer be expected in the future, because the highly scaled CMOS devices would inevitably exhibit probabilistic or statistical behavior. The other problem is that covering all possible worst cases would become increasingly difficult, expensive and result in large yield loss because of the variation.

In addition, we also notice that the clock frequency is no longer scaling significantly with time, and the dark silicon becomes an important issue in circuit design since energy efficiency is considered to be the main limitation for performance improvement in embedded or mobile applications [1], [2]. In this case, designers tend to employ multi-core processors, heterogeneous systems and massive parallel accelerators to sustain the performance scaling.

Specifically, large volume of current studies has demonstrated that, significant performance gains can be achieved by FPGA-based accelerators over software designs across a wide range of applications. However, one of the major factors that limits the performance of these accelerators is that they typically run at much lower clock frequencies than general purpose processors (GPPs) or GPUs. Furthermore, timing analysis tools typically report a very conservative clock frequency to guarantee “safe” operations. This substantially limits the potential performance of the device to a further extent.

In order to boost the operating frequency of a datapath, the standard approaches are either to heavily pipeline the design, or reduce the precision throughout the datapath. For the former method, it should be noted that pipelining will not tend to reduce the circuit latency. Actually the latency in terms of clock cycles will increase. As a result, this method will not be applicable to many embedded applications, which typically have strict latency requirements, or in datapath containing feedback where C-slow retiming is inappropriate. For the second approach, reducing the datapath precision would subtract out the latency, but at the cost of introducing quantization errors into the design. For hardware platforms such as FPGAs or ASICs, research into exploiting the potential benefits of using the minimum precision necessary to satisfy a design specification has been an extensive topic [3], as it is possible to use a number system with customized precisions on these platforms. A brief review of existing approaches in this area will be discussed in Section II-A.

Unfortunately, neither of the conventional approaches tends to remove the conservative safety margin. In recent years, we have seen a growth of research that explores the potential power or performance benefits that can be obtained when relaxing the variation-induced guard-bands [4]. The detailed review of this research field will be presented in Section II-B.

In this paper, we describe an alternative circuit design methodology when considering tradeoffs between accuracy, performance and area. We suggest that for certain applications it is beneficial to move away from the traditional model of creating a conservative design that is guaranteed to avoid timing violations. Instead, it may be preferable to create a design in which timing violations may occur, under the knowledge that they only occur rarely because specific input patterns are required to generate the worst case errors. This paper elaborates on the prior work [xxx], with a further enrichment by incorporating silicon area as another evaluation metric for the trade-offs. In this work, we explicitly demonstrate the optimum design choice of the basic arithmetic operators (i.e.

K. Shi, D. Boland and G. A. Constantinides are with the Department of Electrical and Electronic Engineering, Imperial College London, London, SW7 2BT, U.K. (email: {k.shi11, david.boland03, g.constantinides}@imperial.ac.uk)

Manuscript received April 19, 2005; revised December 27, 2012.

binary adders) under various design constraints. A summary of the main contributions of this work are as follows:

The rest of this paper is organized as follows. It first reviews in detail the current literature regarding the existing approaches that are used in word-length optimization and approximate datapath design in Section II. The analytical error models for two different adder structures: RCA and CSA, are discussed in Section III and Section IV, respectively. The decision of the optimum adder structure under the accuracy-performance-area trade-offs, is then put forward in Section V. It is followed by a description of the probabilistic error models for CCM in Section VI. Section VII discusses a practical experimental setup to verify our models and test the proposed design methodology. The demonstration of the benefits of our proposed approach is then demonstrated in Section VIII. Finally, Section IX discusses the conclusions of this work and the possible research directions in the future.

II. BACKGROUND

A. Word-length Optimization

There exists a significant amount of work demonstrating that optimizing the precision used throughout the datapath would bring a substantial benefits on clock speed, silicon area and power consumption. As loosing precisions causes truncation or overflow errors, a main stream of research focus on analyzing errors generated from using the minimum precision and providing bounds of inputs to ensure the overall error are within the given accuracy specification. Current literatures report two major approaches in this area: the simulation-based approach [xxx] and the analytical-based approach such as interval arithmetic [5], affine arithmetic [6], Satisfiability-Modulo Theories [7] and polynomial representations [8]. A detailed summary of these methods can be found in [3], [9].

B. Imprecise Circuit Design Methodology

Note that, however, the choice of precision is not the only source of errors when designing a datapath. Recently we have seen a growth of parallel streams of research which aim at exploring alternative methods to trade accuracy for design efficiency. This strand of research is motivated by the drastic performance variations introduced by process scaling.

To address this problem, the international technology roadmap for semiconductors (ITRS07) pointed out that extra benefits of manufacturing, test, power and timing can be obtained if the tight requirement of absolute correctness is released for devices and interconnects [10]. This topic is expected to be of growing importance in the future, because we will face new challenges as the technology scales further.

According to the existing literatures, one way to operating circuits beyond the deterministic region is to relax the design constraints and the guard-bands that are conventionally used to avoid the worst cases. Research in this area is based on the observations that in practice, the worst case happens rarely. A series of work known as “Better-Than-Worst-Case (BTWC) Design” introduced a universal framework to push the circuit performance to its limits [11]. In general, a BTWC design is composed of cores and checkers. The cores are operated with

high performance by eliminating the guard-bands, meanwhile the possible timing errors are diagnosed by the checkers. Optionally the system can be recovered at the observation of errors. As an exemplary design, the Razor project [12], [13] was proposed to shave the conservative timing margins by overscaling the supply voltage and clock frequency, while monitoring the output error rates by utilizing a self-checking circuit. This work demonstrated that the benefits brought by removing the safe margin outweigh the cost of monitoring and recovering from errors. For example, 22% or over 30% power consumption can be saved with $\sim 0.01\%$ or $\sim 1\%$ error rates at the output, respectively. Related work includes a similar frequency overscaling technique by operating circuits slightly slower than the critical path delay with dedicated checker circuits to ensure that timing errors will not occur [14], or developing timing analysis tools that decide the optimum operating frequencies in the non-deterministic region due to process variation [15].

Another key observation that has been made from existing work is that the worst-case normally happens with specific input patterns. Meanwhile, many studies have shown that errors can be potentially tolerated in various applications, e.g. DSP applications concerning human perceptions. In this case, an even greater performance/power benefits can be achieved, as the design overhead of the checkers can be eliminated. Current research in this area is focusing on designing probabilistic or imprecise circuits that employ techniques from either software level or hardware level. For instance at the programming language level, a tool was proposed [16], [17] to divide the program into the precise parts and the approximate parts, both of which are mapped to different hardware with different speed-grades, supply voltages etc. This technique enables a relatively high service quality, in the meantime energy reduction can be obtained due to approximation. Similar ideas have also been applied on hardware directly. As an example, Palem et al. [18] described a non-uniform voltage scaling technique for the ripple carry adder. In this study, multiple voltage regions are employed for different bits along a carry chain. That is, higher voltage would be applied for computations generating most significant bits, and vice versa.

While the aforementioned literatures take advantage of the fact that only specific input patterns could cause timing errors, research on imprecise architectures take this one step further by designing simplified circuits for performance/energy efficiency. For example, Lu et al. proposed a “shrinking” datapath that can be utilized to mimic and speculate the original logic functions [19]. Kulkarni et.al. described an underdesigned 2×2 multiplier unit, of which the worst case was replaced by a normal case based on the straight-forward Karnaugh-Map analysis [20]. In both cases, reduction of area and power consumption can be achieved with the cost of accuracy. In addition, Gupta et al. developed approximate adders at the transistor level and compared the energy efficiency of their proposed architectures over truncation of input word-length of conventional structures [21].

However, the limitations of the existing work should not be neglected. Firstly, some techniques, such as employing several voltages regions within a ripple carry adder, are not practical

to achieve in real situations, and the overhead of applying this technique is not discussed. Secondly, many studies propose alternative architectures, which means the adaptability of this kind of technique is limited on current hardware platforms. More importantly, correct results would never be obtained as long as the approximate architectures were employed. Thirdly, the link between the probability of output correctness and energy savings or performance improvements are rarely analyzed from current research.

In this work, we propose an alternative methodology to remove these limitations. It is argued in this paper that our method can be easily applied onto existing hardware arithmetic operators such as RCAs and CCMs with almost no hardware overhead. We support this hypothesis experimentally and analytically by presenting theoretical probabilistic error models.

III. RIPPLE CARRY ADDER

A. Adder Structures in FPGAs

Adders serve as a key building block for arithmetic operations. Generally speaking, the ripple carry adder (RCA) is the most straightforward and widely used adder structure. As such, the philosophy of our approach is first exemplified with the analysis of a RCA. We later describe how this methodology can be extended to other arithmetic operators in Section VI by discussing the CCM that is commonly used in DSP applications and numerical algorithms.

Typically the maximum frequency of a RCA is determined by the longest carry propagation. Consequently, modern FPGAs offer built-in architectures for very fast ripple carry addition. For instance, the Altera Cyclone series uses fast tables [22] while the Xilinx Virtex series employs dedicated multiplexers and encoders for the fast carry logic [23]. Figure 1 illustrates the structure of an n -bit RCA, which is composed of n serial-connected full adders (FAs) and utilizes the internal fast carry logic of the Virtex-6 FPGA.

While the fast carry logic reduces the time of each individual carry-propagation delay, the overall delay of carry-propagation will eventually overwhelm the delay of sum generation of each LUT with increasing operand word-lengths. For our initial analysis, we assume that the carry propagation delay of each FA is a constant value μ , which is a combination of logic delay and routing delay, and hence the critical path delay of the RCA is $\mu_{RCA} = n\mu$, as shown in Figure 1. For an n -bit RCA, it follows that if the sampling period T_S is greater than μ_{RCA} , correct results will be sampled. If, however, $T_S < \mu_{RCA}$, intermediate results will be sampled, potentially generating errors.

In the following sections, we consider two methods that would allow the circuit to run at a frequency higher than $1/T_S$. The first is a traditional circuit design approach where operations occur without timing violations. To this end, the operand word-length is truncated in order to meet the timing requirement. This process results in truncation or roundoff error. In our proposed new scenario, circuits are implemented with greater word-length, but are clocked beyond the safe region so that timing violations sometimes occur. This process generates "overclocking error".

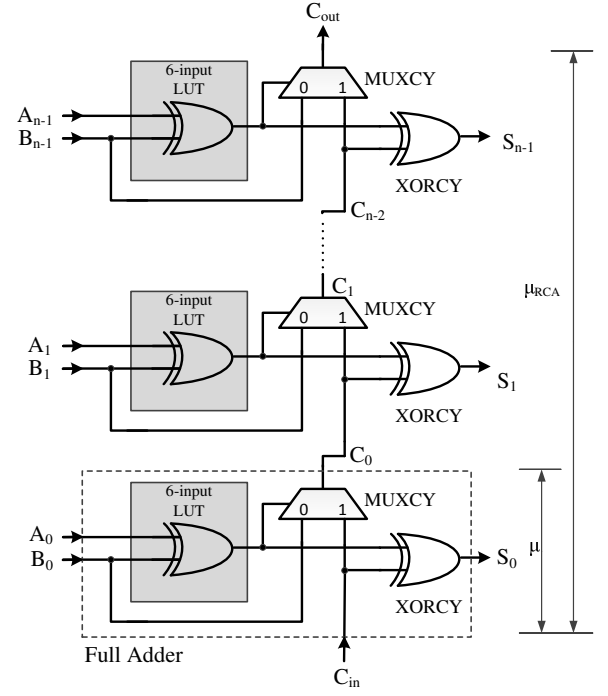


Fig. 1. An n -bit ripple carry adder in the Xilinx Virtex-6 FPGA.

B. Probabilistic Model of Truncation Error

For ease of discussion, we assume that the input to our circuit is a fixed point number scaled to lie in the range $[-1, 1)$. For our initial analysis, we assume every bit of each input is uniformly and independently generated. However, this assumption will be relaxed in Section VIII where the predictions are verified using real image data. The errors at the output are evaluated in terms of the absolute value and the probability of their occurring. These two metrics are combined as the error expectation.

If the input signal of a circuit is k bits, truncation error occurs when the input signal is truncated from k bits to n bits. Under this premise, the mean value of the truncated bits at signal input (E_{Tin}) is given by (1).

$$E_{Tin} = \frac{1}{2} \sum_{i=n+1}^k 2^{-i} = 2^{-n-1} - 2^{-k-1} \quad (1)$$

Since we assume there are two mutually independent inputs to the RCA, the overall expectation of truncation error for the RCA is given by (2).

$$E_T = \begin{cases} 2^{-n} - 2^{-k}, & \text{if } n < k \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

C. Probabilistic Model of Overclocking Error

1) *Generation of Overclocking Error:* For a given T_S , the maximum length of error-free carry propagation is described by (3), where f_S denotes the sampling frequency.

$$b := \left\lceil \frac{T_S}{\mu} \right\rceil = \left\lceil \frac{1}{\mu \cdot f_S} \right\rceil \quad (3)$$

However, since the length of an actual carry chain during execution is dependent upon input patterns, in general, the worst case may occur rarely. To determine when this timing constraint is not met and the size of the error in this case, we expand standard results [24] to the following statements, which examine carry generation, propagation and annihilation, as well as the corresponding summation results of a single bit i , according to the relationship between its input patterns A_i and B_i :

- If $A_i = B_i = 1$, a new carry chain is generated at bit i , and $S_i = C_{i-1}$;
- If $A_i \neq B_i$, the carry propagates for this carry chain at bit i , and $S_i = 0$;
- If $A_i = B_i$, the current carry chain annihilates at bit i , and $S_i = 1$.

2) *Absolute Value of Overclocking Error*: For an n -bit RCA, let C_{tm} denote the carry chain generated at bit S_t with the length of m bits. For a certain f_S , the maximum length of error-free carry propagation, b , is determined through (3). The presence of overclocking error requires $m > b$. Since the length of carry chain cannot be greater than n , parameters t and m are bounded by (4) and (5):

$$0 \leq t \leq n - b \quad (4)$$

$$b < m \leq n + 1 - t \quad (5)$$

For C_{tm} , correct results will be generated from bit S_t to bit S_{t+b-1} . Hence the absolute value of error seen at the output, normalized to the MSB (2^n), is given by (6), where \hat{S}_i and S_i denote the actual and error-free output of bit i respectively.

$$e_{tm} = \frac{\left| \sum_{i=t+b}^n (S_i - \hat{S}_i) \cdot 2^i \right|}{2^n} \quad (6)$$

S_i and \hat{S}_i can be determined using the equations from the previous statements in Section III-C1. In the error-free case, the carry will propagate from bit S_t to bit S_{t+m-1} , and we will obtain $S_{t+b} = S_{t+b+1} = \dots = S_{t+m-2} = 0$ for carry propagation, and $S_{t+m-1} = 1$ for carry annihilation. However, when a timing violation occurs, the carry will not propagate through all these bits. Substituting these values into (6) yields (7). Interestingly, the value of overclocking error has no dependence on the length of carry chain m .

$$e_{tm} = \frac{|2^{t+m-1} - 2^{t+m-2} - \dots - 2^{t+b}|}{2^n} = 2^{t+b-n} \quad (7)$$

3) *Probability of Overclocking Error*: The carry chain C_{tm} occurs when there is a carry generated at bit t , a carry annihilated at bit $t+m-1$ and the carry propagates in between. Consequently, its probability P_{tm} is given by (8).

$$P_{tm} = P_{(A_t=B_t=1)} P_{(A_{t+m-1}=B_{t+m-1})} \cdot \prod_{i=t+1}^{t+m-2} P_{(A_i \neq B_i)} \quad (8)$$

Under the assumption that A and B are mutually independent and uniformly distributed, we have $P_{(A_i=B_i=1)} = 1/4$, $P_{(A_i \neq B_i)} = 1/2$ and $P_{(A_i=B_i)} = 1/2$, so P_{tm} can be obtained

by (9). Note that (9) takes into account the carry annihilation always occurs when $t + m - 1 = n$.

$$P_{tm} = \begin{cases} (1/2)^{m+1} & \text{if } t + m - 1 < n \\ (1/2)^m & \text{if } t + m - 1 = n \end{cases} \quad (9)$$

4) *Expectation of Overclocking Error*: Expectation of overclocking error can be expressed by (10).

$$E_O = \sum_t \sum_m P_{tm} \cdot e_{tm} \quad (10)$$

Using P_{tm} and e_{tm} from (7) and (9) respectively, E_O can be obtained by (11).

$$E_O = \begin{cases} 2^{-b} - 2^{-n-1}, & \text{if } b \leq n \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

D. Comparison between Two Scenarios

In the traditional scenario, the word-length of RCA must be truncated, using $n = b - 1$ bits, in order to meet a given f_S . The error expectation is then given by (12).

$$E_{trad} = 2^{-b+1} - 2^{-k} \quad (12)$$

Overclocking errors are allowed to happen in the second scenario, therefore the word-length of RCA is set to be equal to the input word-length, that is, $n = k$. Hence we obtain (13) according to (11).

$$E_{new} = 2^{-b} - 2^{-k-1} \quad (13)$$

Comparing (13) and (12), we have (14). This equation indicates that by allowing timing violations, the overall error expectation of RCA outputs drops by a factor of 2 in comparison to traditional scenario. This provides the first hint that our approach is useful in practice.

$$\frac{E_{new}}{E_{trad}} = \frac{2^{-b} - 2^{-k-1}}{2^{-b+1} - 2^{-k}} = \frac{1}{2} \quad (14)$$

IV. CARRY SELECT ADDER

A. Introduction

Since the delay of RCA is determined by the length of carry chain, alternative adder architectures, such as carry select adder (CSA) have been proposed to boost performance. For CSA, the carry chain is divided into multiple overlapped sections to increase the operating speed, as shown in Fig. 2(a). There are multiple stages within a CSA. Each stage contains two RCAs and two multiplexers, as seen in Fig. 2(b). For a given input, two additions are performed simultaneously where the carry input is zero and one respectively. One of these two results is then selected according to the actual carry input. Although this structure brings performance benefits, it costs extra hardware resources compared to a standard RCA because the carry chain is duplicated. Furthermore, in FPGA technology, multiplexers are expensive. Due to this reason, we explore the trade-offs between silicon area, accuracy and performance of RCA and CSA in this section.

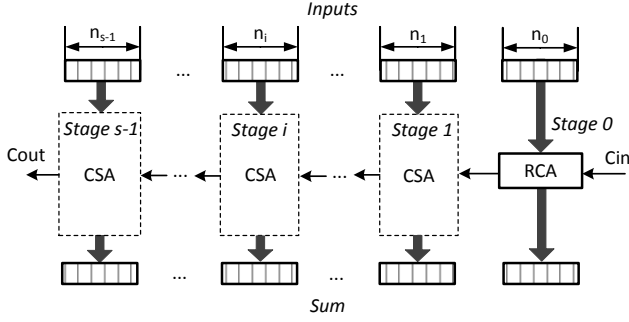
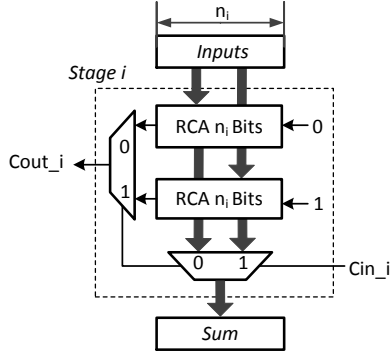
(a) The structure diagram of a CSA with s stages.(b) The structure of the i^{th} stage in the CSA.

Fig. 2. The structure diagram of CSA.

B. Timing Models for Carry Select Adder

We initially model the CSA in order to understand the relationship between the operating frequency and maximum word-length of the CSA. This information can then be employed to determine the truncation error based on the models presented in Section III-B.

In our analysis, the stage delay of stage i in the CSA refers to the combination of the i -bit carry propagation delay and the delay of multiplexing the carry output. For a CSA with s stages ($s \geq 2$), let the stage delay be denoted by d_{s-1}, \dots, d_0 , where d_{s-1} and d_0 represent the delay of the most significant and the least significant stages, respectively. We still follow the aforementioned assumption that the critical path delay of the CSA is due to carry propagation and multiplexing the carry output, instead of generating the sum outputs. It should be noted that unlike other stages, the least significant stage of the CSA is only built by one RCA without multiplexers, since it is directly driven by the carry input. Hence we can obtain the delay of the i^{th} stage as presented in (15), where μ_c , μ_{mux} and n_i denote the delay of 1-bit carry propagation, the delay of multiplexing and the word-length of the i^{th} stage of the CSA, respectively.

$$d_i = \begin{cases} n_i \cdot \mu_c + (s - i) \cdot \mu_{mux}, & \text{if } i \in [1, s - 1] \\ n_0 \cdot \mu_c + (s - 1) \cdot \mu_{mux}, & \text{if } i = 0 \end{cases} \quad (15)$$

Under the timing-driven design environment, the delay of each stage of CSA is equalized in order to achieve the fastest

operation, as presented in (16).

$$d_{s-1} = d_{s-2} = \dots = d_0 \quad (16)$$

In this case, combining (15) and (16) yields n_i , which is represented by the word-length of the most significant stage n_{s-1} , in (17).

$$n_i = \begin{cases} n_{s-1} - (s - 1 - i) \cdot \frac{\mu_{mux}}{\mu_c}, & \text{if } i \in [1, s - 1] \\ n_{s-1} - (s - 2) \cdot \frac{\mu_{mux}}{\mu_c}, & \text{if } i = 0 \end{cases} \quad (17)$$

Sum up n_i to give the total word-length of the CSA in (18).

$$\begin{aligned} n_{CSA} &= \sum_{i=0}^{s-1} n_i \\ &= s \cdot n_{s-1} - \frac{\mu_{mux}}{\mu_c} \cdot \frac{(s+1)(s-2)}{2} \end{aligned} \quad (18)$$

Conventionally under a given frequency constraint, the word-length of RCA is truncated by using $n_{RCA} = b - 1$ bits, where b is determined by (3). Similarly for CSA, the word-length of each stage should be selected in order to satisfy (19).

$$\forall i \in [0, s - 1] \quad , \quad d_i \leq \frac{1}{f_s} \quad (19)$$

Hence for the same frequency requirement, we can form the relationship between the delay of the most significant stage of CSA and the delay of RCA, as given by (20),

$$\mu_c \cdot (b - 1) = n_{s-1} \cdot \mu_c + \mu_{mux} \quad (20)$$

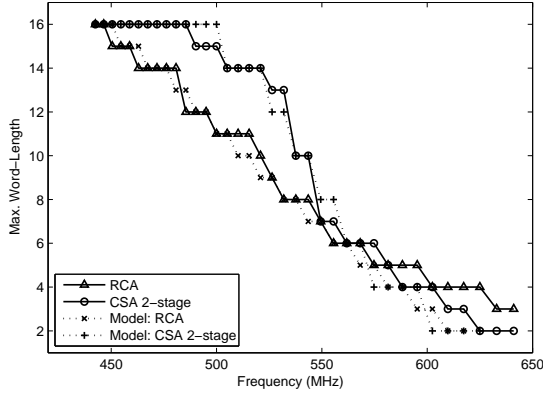
Substitute (20) into (18) to replace n_{s-1} , we derive the representation of the word-length of CSA in terms of b , as presented in (21).

$$n_{CSA} = s \cdot (b - 1) - \frac{\mu_{mux}}{\mu_c} \cdot \frac{(s+2)(s-1)}{2} \quad (21)$$

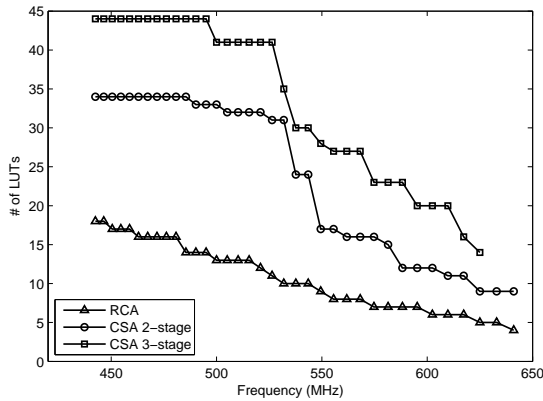
C. Accuracy Benefits and Area Overhead in CSA

We first verify the models for the RCA and CSA in terms of the maximum word-length under the given operating frequencies. For the CSA, the ratio μ_{mux}/μ_c can be computed experimentally. We perform post place-and-route simulations on the CSA with 2 stages using Xilinx Virtex-6 FPGA. The delay of the i^{th} stage d_i in (15) is recorded with respect to different word-lengths. The total word-length of CSA can be predicted through (21). In addition, the maximum word-lengths of the 2-stage CSA and RCA are obtained experimentally by increasing n_{CSA} and n_{RCA} respectively until errors are observed at the output. The comparison between the modeled value and the empirical results is illustrated in Fig. 3(a). It can be seen that our models for both RCA and CSA match well with the experimental results.

Fig. 3(a) also highlights that in comparison to the RCA, the CSA achieves greater word-length when frequency is initially increased. The RCA only outperforms than the CSA when very high frequency is applied. This is because at low frequencies, although the multiplexer delay limits the word-length of each stage in CSA when compared to RCA, the stage parallelism in CSA enables a greater word-length. However when frequency increases, the multiplexer delay becomes comparable to the



(a) The modeled value and the experimental results of the maximum word-length of RCA and CSA.



(b) Hardware resource usage for an RCA and a CSA with 2 and 3 stages.

Fig. 3. A comparison between RCA and CSA in terms of the maximum word-length of input signal and the area consumption.

delay of the carry chain, and this inhibits the benefits of parallelism.

However, the accuracy benefits brought by CSA comes at the cost of a large area overhead. Fig. 3(b) depicts the number of Look-Up Tables (LUTs) in the FPGA used for an RCA, a 2-stage CSA and a 3-stage CSA. It can be seen that in order to meet a given frequency, the 3-stage CSA consumes $2.4 \times \sim 3.7 \times$ area than RCA, while the 2-stage CSA requires $1.7 \times \sim 3.1 \times$ extra area. This finding poses a question of which is the best adder structure for a specific area budget.

V. CHOOSING THE OPTIMUM ADDER STRUCTURE

In Section III, we discussed two design scenarios for the RCA when considering timing constraints. In this section, we expand our analysis by incorporating silicon area as another evaluation metric, and investigate the accuracy, performance and area trade-offs for different adder structures. In the conventional design scenario, the word-length of RCA and CSA is limited by the given frequency constraint, or/and the available hardware resources. The precision loss potentially generates large errors even without timing violations. However in the new design scenario, we use the RCA with the maximum possible word-length under the given area budget, and the

timing constraints are allowed to be violated. This process might result in timing errors as well as truncation errors due to area limitation. We compare these two scenarios with different design goals with the aim of finding the optimum design methodology under each situation.

A. Determination of the Optimum Adder Structure for Given Frequency Requirements

First, suppose the algorithm designer wished to create a circuit that can run at a given frequency with the minimum achievable output errors and minimum resource usage, i.e. a pair of $\{Area, Frequency\}$ constraint is applied. In this case, the optimum design method is selected according to the following criteria:

- Design with the minimum mean value of errors at the output is the optimum design;
- If multiple designs achieve the same accuracy, then the design with minimum area is the optimum design;
- If both the accuracy and area are identical for multiple designs, they are all treated as the optimum design.

For instance, in Fig. 4 we record the mean value of error at outputs with respect to different operating frequencies for both design scenarios when the number of the available LUTs is set to 25. In this graph, we have labeled the optimum adder structure. Note that in this and all the following experiments within this section, in order to apply our models the input data is randomly generated following the uniform distribution.

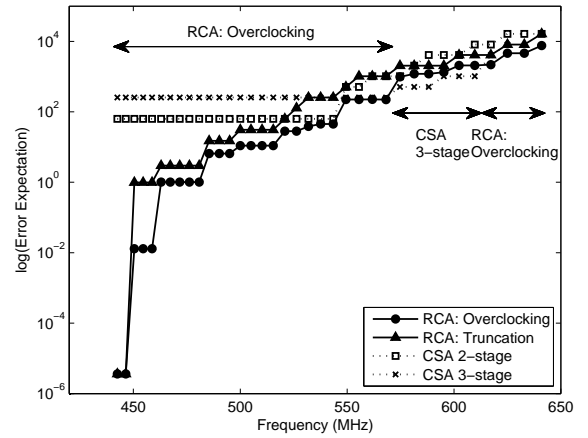


Fig. 4. A comparison between two design scenarios when the number of available LUT is 25. The RCA and CSA are investigated in the conventional scenario, while the RCA is explored in the proposed new scenario. The results are obtained from post place-and-route simulations on Xilinx Virtex-6 FPGAs.

We first notice that for all frequency values, the overclocked RCA achieves smaller error expectation than the RCA with truncated operand word-lengths, as predicted by the models for the RCA in Section III-D. It can also be observed that the CSA cannot be implemented with the original word-length due to the limited area budget and this leads to large truncation errors. Although the CSA with 3 stages outperforms temporarily when frequency is further increased, the overclocked RCA is still the optimum design when high operating frequencies are applied.

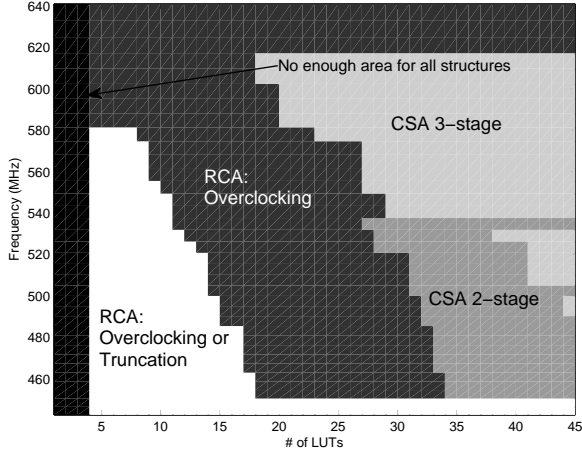


Fig. 5. A demonstration of the optimum design methodology which achieves the minimum error at outputs with respect to a variety of frequency and area constraints.

We then perform similar experiments with a variety of area constraints. The optimum design methods with respect to different operating frequencies and area consumptions are demonstrated in Fig. 5. From this figure, several observations can be made. Firstly, if the available area is large enough to implement a CSA in full precision, it will be the optimum design. This is expected from our earlier analysis in Fig. 3(a). Secondly, the 2-stage CSA is better than the 3-stage CSA when frequency is initially increased, as it consumes less area although both achieve the same error expectation. Thirdly, only part of the CSA can be implemented under a tighter area budget, whereas the RCA still keeps full precision. In this case, area becomes the dominate factor and precision is lost for the CSA, meaning the RCA with overclocking is the optimum design method across almost the whole frequency domain in this situation. Last but not the least, for a more stringent area constraint, the word-length of RCA is also limited. This results in truncation errors initially for all design scenarios. However, the RCA with overclocking can still be employed as the optimum design, because it loses less precision than the CSA under the same area constraint.

B. Determination of the Optimum Adder Structure for Given Accuracy Requirements

If the design goal is to operate the circuit as fast as possible with the minimum area whilst a certain error budget can be tolerated, the optimum design methodology can be decided as illustrated in Fig. 6. In this situation, the error specifications are evaluated in terms of the mean relative error (MRE), as given by (22), where E_{error} and E_{out} refer to the mean value of error and the mean value of outputs, respectively.

$$MRE = \left| \frac{E_{error}}{E_{out}} \right| \times 100\% \quad (22)$$

In our experiments, MRE is set ranging from 0.001% to 50%. For a certain value of MRE, the design with the maximum operating frequency is selected as the optimum

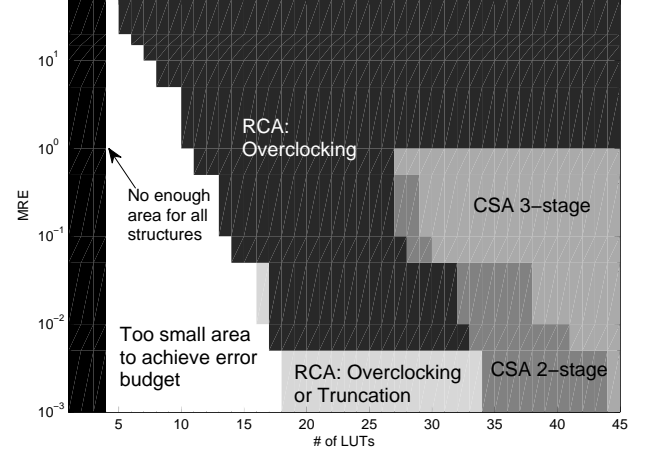


Fig. 6. A demonstration of the optimum design methodology which runs at the fastest frequency with respect to a variety of accuracy and area constraints.

design. Moreover, the smallest design is the optimum one if multiple structures operate at the same frequency, with a certain area requirement. Fig. 6 can thus be obtained based on these criteria.

For a tight accuracy requirement, i.e. $MRE < 0.005\%$, CSA serves as the optimum design choice with respect to large accessible area, as it intrinsically operates faster than RCA. Once again, when the area budget shrinks, the RCA performs best because the precision of the CSA is limited. Similarly to the results in Fig. 5, we see that the overclocked RCA achieves the fastest operating frequencies under most area constraints when the accuracy requirement is released.

C. Design Guidance

To sum up, the experiments reveal that for both design goals, CSA is the best option only when there is enough hardware resources. However for the remainder of this paper, we are interested in cases where there is a limited area budget, as may be the case in the embedded applications. Due to this reason, we will focus on the RCA in our following analysis and experiments.

VI. CONSTANT COEFFICIENT MULTIPLIER

As another key primitive of arithmetic operations, CCM can be implemented using RCA and shifters. For example, operation $B = 9A$ is equivalent to $B = A + 8A = A + (A << 3)$, which can be built using one RCA and one shifter. We first focus on a single RCA and single shifter structure. We describe how more complex structures consisting of multiple RCAs and multiple shifters can be built in accordance with this baseline structure in Section VI-C.

In this CCM structure, let the two inputs of the RCA be denoted by A_S and A_O respectively, which are both two's complement numbers. A_S denotes the "shifted signal", with zeros padded after LSB, while A_O denotes the "original signal" with MSB sign extension. For an n -bit input signal, it should be noted that an n -bit RCA is sufficient for this

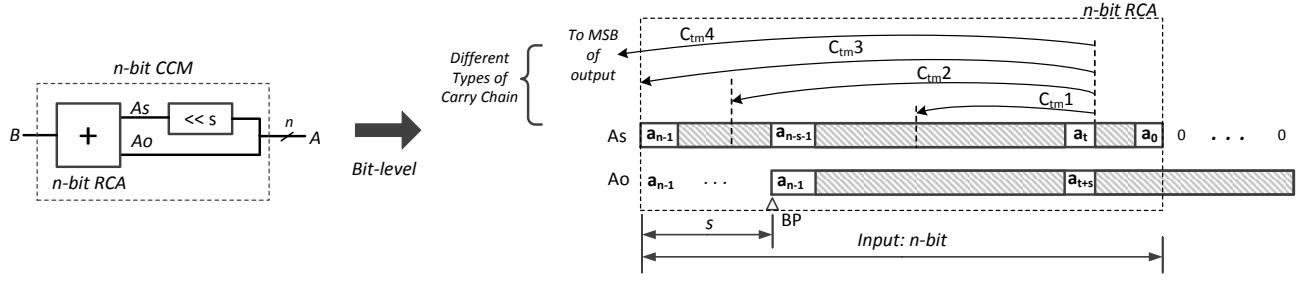


Fig. 7. Four possible carry chain types in a constant coefficient multiplier with n -bit inputs. The notion s denotes the shifted bits and BP denotes the binary point.

operation, because no carry will be generated or propagated when adding with zeros, as shown in Fig. 7.

A. Probabilistic Model of Truncation Error

Let E_{Tin} and E_{Tout} denote the expectation of truncation error at the input and output of CCM respectively. We then have (23), where coe denotes the coefficient value of the CCM, and E_{Tin} can be obtained according to (2).

$$E_{Tout} = |coe| \cdot E_{Tin} \quad (23)$$

B. Probabilistic Model of Overclocking Error

1) *Absolute Value of Overclocking Error*: The absolute value of overclocking error of carry chain C_{tm} is increased by a factor of 2^s due to shifting, compared to RCA. Hence e_{tm} in CCM can be modified from (7) to give (24).

$$e_{tm} = 2^{t+b-n+s} \quad (24)$$

2) *Probability of Overclocking Error*: Due to the dependencies in a CCM, carry generation requires $a_t = a_{t-s} = 1$, propagation and annihilation of a carry chain is best considered separately for four types of carry chain generated at bit t . We label these by C_{tm1} to C_{tm4} in Figure 7, defined by the end region of the carry chain. For C_{tm1} , we have:

- Carry propagation: $a_i \neq a_{i+s}$ where $i \in [t+1, n-s-2]$;
- Carry annihilation: $a_j = a_{j+s}$ where $j \in [t+1, n-s-1]$.

Similarly for C_{tm2} , we have:

- Carry propagation: $a_i \neq a_{n-1}$ where $i \in [n-s-1, n-3]$; or $a_i \neq a_{i+s}$ where $i \in [t+1, n-s-2]$;
- Carry annihilation: $a_j = a_{n-1}$ where $j \in [n-s-1, n-2]$.

For the first two types of carry chain C_{tm1} and C_{tm2} , the probability of carry propagation and annihilation is $1/2$ and the probability of carry generation is $1/4$, under the premise that all bits of input signal are mutually independent. Therefore (25) can be obtained by substituting this into (8).

$$P_{tm} = (1/2)^{m+1}, \quad \text{if } t+m-1 \leq n-2 \quad (25)$$

For carry annihilation of C_{tm3} , $a_{n-1} = a_{n-1}$, which is always true. Thus the probability of C_{tm3} is given by (26).

$$P_{tm} = (1/2)^m, \quad \text{if } t+m-1 = n-1 \quad (26)$$

C_{tm4} represents carry chain annihilates over a_{n-1} , therefore carry propagation requires $a_{n-1} \neq a_{n-1}$. This means C_{tm4} never occurs in a CCM.

Altogether, P_{tm} for a CCM is given by (27).

$$P_{tm} = \begin{cases} (1/2)^{m+1} & \text{if } t+m-1 < n-1 \\ (1/2)^m & \text{if } t+m-1 = n-1 \end{cases} \quad (27)$$

3) *Expectation of Overclocking Error*: Since the carry chain of a CCM will not propagate over a_{n-1} , the upper bound of parameter t and m should be modified from (4) and (5) to give (28) and (29).

$$0 \leq t \leq n-b-1 \quad (28)$$

$$b < m \leq n-t \quad (29)$$

Finally, by substituting (27) and (24) with modified bounds of t and m into (10), we obtain the expectation of overclocking error for a CCM to be given by (30).

$$E_O = \begin{cases} 2^{s-b-1} - 2^{s-n-1}, & \text{if } b \leq n-1 \\ 0, & \text{otherwise} \end{cases} \quad (30)$$

C. CCM with Multiple RCAs and Shifters

In the case where a CCM is composed of two shifters and one RCA, such as operation $B = 20A = (A \ll 2) + (A \ll 4)$, let the shifted bits be denoted as s_1 and s_2 respectively. Hence the equivalent s in (30) can be obtained through (31).

$$s = |s_1 - s_2| \quad (31)$$

For those operations such as $B = 37A = (A \ll 5) + (A \ll 2) + (A \ll 1)$, the CCM can be built using a tree structure. Each root node is the baseline CCM and the errors are propagated through an adder tree, of which the error can be determined based on our previous RCA model.

VII. TEST PLATFORM

In our experiments, we compare two design perspectives. In the first scenario, the word-length of the input signal is truncated before propagating through the datapath in order to meet a given latency. In our proposed overclocking scenario, the circuit is overclocked while keeping the original operand word-length. The benefits of the proposed methodology are demonstrated over a set of DSP example designs, which are implemented on the Xilinx ML605 board with a Virtex-6 FPGA.

A. Experimental Setup

We initially build up a test framework on an FPGA. The general architecture is depicted in Fig. 8. The main body of the test framework consists of the circuit under test (CUT), the test frequency generator and the control logic, as shown in the dotted box in Fig. 8. The I/Os of the CUT are registered by the launch registers (LRs) and the sample registers (SRs), which are all triggered by the test clock. Input test vectors are stored in the on-chip memory during initialization. The results are sampled using Xilinx ChipScope. Finally, we perform an offline comparison of the output of the original circuit at the rated frequency with the output of the overclocked as well as the truncated designs using the same input vectors.

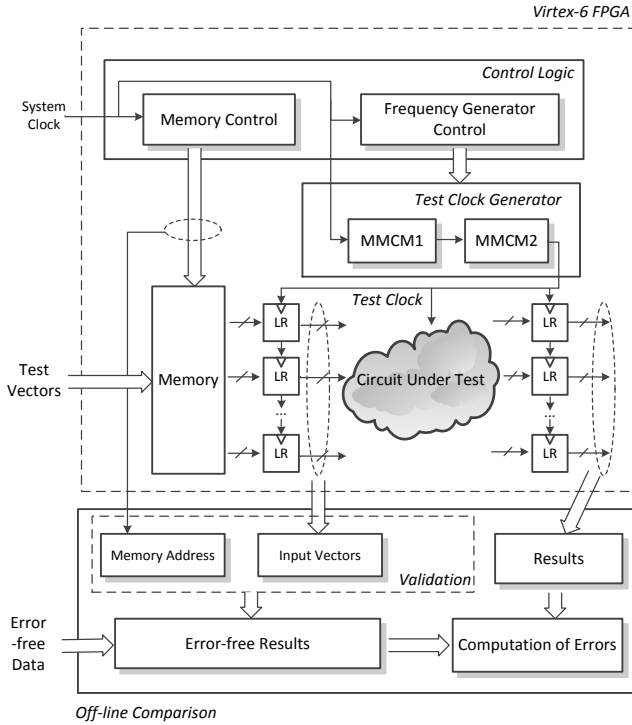


Fig. 8. Test framework, which is composed of a measurement architecture (the dotted box) on an FPGA and an off-line comparator using software. Note that the error-free data are obtained by either pre-computation or initial run with low frequencies.

The test frequency generator is implemented using two cascaded mixed-mode clock managers (MMCMs), created using Xilinx Core Generator [25]. Besides the outputs, the corresponding input vectors and memory addresses are also recorded into the comparator, as can be seen in Fig. 8, in order to ensure that the recorded errors arise from overclocking the CUT rather than the surrounding circuitry when high test frequencies are applied.

B. Benchmark Circuits

Three types of DSP designs are tested: digital filters (FIR, IIR and Butterworth), a Sobel edge detector and a direct implementation of a Discrete Cosine Transformation (DCT). The filter parameters are generated through MATLAB filter

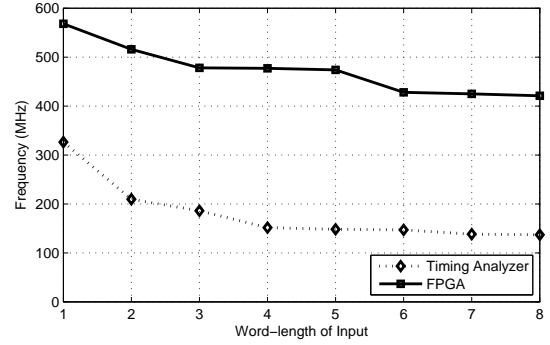


Fig. 9. The maximum operating frequencies for different input word-lengths of an FIR filter. The dotted line depicts the rated frequency reported by the timing analysis tool. The solid line is obtained through real FPGA tests using our platform.

design toolbox, and they are normalized to integers for implementation. Table I summarizes the operating frequency of each implemented design in Xilinx ISE14.1 when the word-length of input signal is 8-bit.

TABLE I
RATED FREQUENCIES OF EXAMPLE DESIGNS.

Design	Frequency (MHz)	Description
FIR Filter	126.2	5 th order
Sobel Edge Detector	196.7	3 × 3
IIR Filter	140.3	7 th order
Butterworth Filter	117.1	9 th order
DCT	176.7	4-point

The input data are generated from two sources. One is called “uniform independent inputs”, which are randomly sampled from a uniform distribution of 8-bit numbers. The other is referred to as “real inputs”, which denote 8-bit pixel values of the 512×512 Lena image.

C. Exploring the Conservative Timing Margin

Generally, the operating frequency provided by EDA tools tends to be conservative to ensure the correct functionality under a wide range of operating environments and workloads. In a practical situation, this may result in a large gap between the predicted frequency and the actual frequency under which the correct operation is maintained [26].

For example, the predicted frequencies and the actual frequencies of a 5th order FIR filter using different word-lengths are depicted in Fig. 9. The “actual” maximum frequencies are computed by increasing the operating frequency from the rated value until errors are observed at the output; the maximum operating frequency with correct output is recorded for the current word-length. As can be seen in Fig. 9, the circuit can operate without errors at a much higher frequency in practice than predicted according to our experiments. A maximum speed differential of 3.2× is obtained when the input signal is 5-bit.

In our experiments in Section VIII, the conservative timing margin is removed in the traditional scenario for a fairer

comparison to the overclocking scenario. To do this, for each truncated word-length, we select the maximum frequency at which we see no overclocking error on the FPGA board in our lab. For example, in Fig. 9, the operating frequency of the design when the word-lengths are truncated to 8, 5 and 2 bits are 400MHz, 450MHz and 500MHz respectively.

Fig. 9 also demonstrates that when the circuit is truncated, it allows the circuit to operate at a higher frequency than the frequency of full precision implementation. However, a non-uniform period change can be observed for both results. For instance, the maximum operating frequency keeps almost constant when the operand word-length reduces from 8 to 6 or from 5 to 3 in both the experimental results and those of timing analyzer. This will cause a slight deviation between our analytical model which assumes that the single bit carry propagation delay to be a constant value, as discussed in (12) with expression $n = b - 1$. This deviation will be influenced by many factors including how the architecture has been packed onto LUTs and CLBs and process variation causing non-uniform interconnection delays [27]. However, we shall see that our model remains close to the true empirical results in Section VIII.

D. Computing Model Parameters

The accuracy of our proposed models is examined with practical results on Virtex-6 FPGA. We first determine the model parameters. There are two types of parameters in the models of overclocking error. The first is based on the circuit architecture. For example, the word-length of RCAs and CCMs (n), the shifted bits of the shifters in CCM (s), and the word-length of the input signal (k). This is determined through static analysis. The second depends on timing information, such as the single bit carry propagation delay μ . In order to keep consistency with the assumption made in models that μ is a fixed value, it is obtained according to the actual FPGA measurement results.

Initially the maximum error-free frequency f_0 is applied. In this case we have (32) where d_c is a constant value which denotes the interconnection delay, and $d_0 = 1/f_0$. The frequency is then increased such that (33) is obtained. This process repeats until the maximum frequency f_{n-1} is applied in (34). Based on these frequency values, μ can be determined.

$$d_0 = n\mu + d_c \quad (32)$$

$$d_1 = (n-1)\mu + d_c \quad (33)$$

...

$$d_{n-1} = \mu + d_c \quad (34)$$

VIII. RESULTS AND DISCUSSION

A. Case study: FIR filter

We first assess the accuracy of our proposed models of error. The modeled values of both overclocking error and truncation error of the FIR filter are presented in Fig. 10 (dotted lines), as well as the actual measurements on the FPGA (solid lines) with two types of input data. The results demonstrate that our

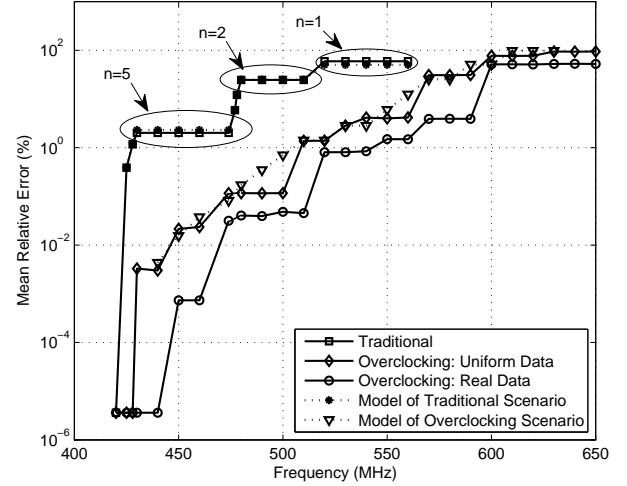


Fig. 10. A demonstration of two design perspectives with a 5th order FIR filter, which is implemented on Virtex-6 FPGA. The modeled values of both overclocking errors and truncation errors are presented as dotted lines. The actual FPGA measurements are depicted using solid lines. Two types of inputs are employed in the overclocking scenario: the uniformly distributed data and the real image data from Lena.

models match well with the practical results obtained using the uniform independent inputs.

According to Fig. 10, output errors are reduced in the overclocking scenario for both input types in comparison to the traditional scenario, as the analytical model validates. In addition, we see that using real data, more significant reduction of MRE are achieved, and that no errors are observed when frequency is initially increased. This is because for real data, long carry chains are typically generated with even smaller probabilities, and the longest carry chain rarely occurs.

The output images of the FIR filter for both of the two scenarios with increasing frequencies are presented in Fig. 11, from which we can clearly see the differences between the errors generated in these two scenarios. In the overclocking scenario, we observe errors in the MSBs for certain input patterns. This leads to “salt and pepper noise”, as shown on the images in the top row of Fig. 11. In the traditional scenario, truncation causes an overall degradation of the whole image, as can be seen in the bottom row of Fig. 11. Furthermore, it is difficult to recover from the latter type of error, since it is generated due to precision loss.

In addition, we record the probability distribution of different length of carry chains using an 8-bit RCA with both data types, as shown in Fig.xxx. As the input data is 8-bit, the longest possible carry chain length is 9-bit. It can be observed that when using Lena data, the probability is only higher for carry chain with a length of 3-bit. While in other situations, using uniform data leads to higher probability, especially for long carry chain length. This finding explicitly demonstrates that longer carry chains happens with an even smaller probability when utilizing real data.

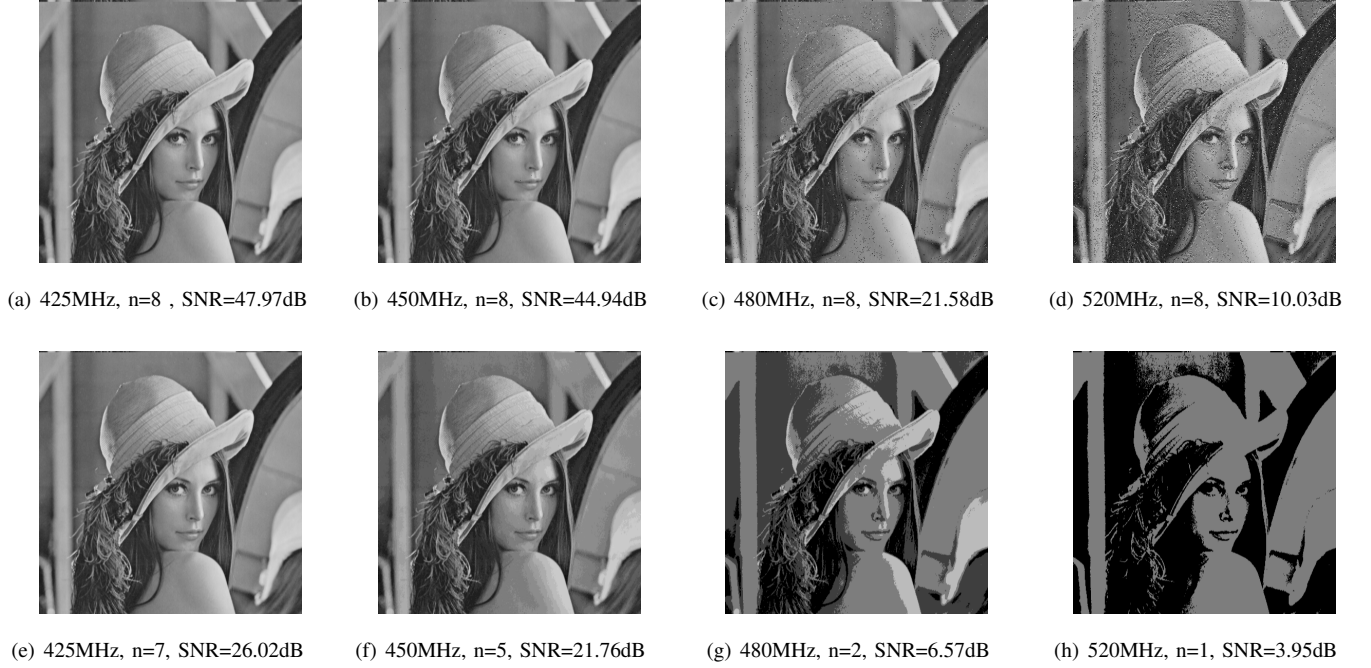


Fig. 11. Output images of the FIR filter for both overclocking scenario (top row) and traditional scenario (bottom row) under various operating frequencies.

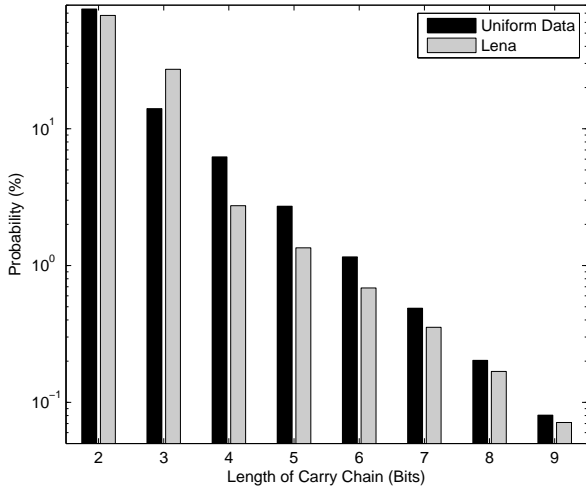


Fig. 12. Probability distribution of different length of carry chains in a 8-bit RCA. For the uniform data, two inputs of RCA are randomly sampled from a uniform distribution. For the image data Lena, one input of RCA uses the original data and the other uses the delayed original data for several clock cycles.

B. Potential Benefits in Datapath Design

As we mentioned in Section V, our results could be of interest to a circuit designer in two ways. Typically, either the designer will want to create a circuit that can run at a given frequency with the minimum possible MRE, or the algorithm designer will wish to run as fast as possible whilst maintaining a specific error tolerance. For the first design target, the experimental results for all five example designs on FPGA are summarized in Table IV in terms of the relative

reduction of MRE as given in (35) where MRE_{Trad} and MRE_{ovrc} denote the value obtained in the traditional scenario and in the overclocking scenario, respectively.

$$\frac{MRE_{Trad} - MRE_{ovrc}}{MRE_{Trad}} \times 100\% \quad (35)$$

In this table, the frequency is normalized to the maximum error-free frequency for each design when the input signal is 8-bit. The N/A in Table IV refers to the situations where a certain frequency simply cannot be achieved using the traditional scenario. It can be seen that a significant reduction of MRE can be achieved using the proposed overclocking scenario, and the geometric mean reduction varies from 67.9% to 95.4% using uniform input data. Even larger differences of MRE can be observed when testing with real image data for each design, ranging from 83.6% to 98.8%, as expected given the results shown in Fig. 10.

Table V illustrates the frequency speedups for each design when the specified error tolerance varies from 0.05% to 50%. For all designs, we see that the overclocking scenario still outperforms the traditional scenario for each MRE budget in terms of operating frequency. Likewise, the frequency speedup is higher for real image inputs than uniform inputs. The geometric mean of frequency speedups of 3.1% to 21.8% can be achieved by using uniform data, while 5.3% to 27.6% when using real image data.

C. Area Overhead

For both aforementioned design goals, it should be noted that an extra benefit of using the traditional approach is that reducing datapath precision would potentially leads to a smaller design, in comparison to our proposed overclocking

TABLE II
RELATIVE REDUCTION OF MRE IN OVERCLOCKING SCENARIO FOR VARIOUS NORMALIZED FREQUENCIES BASED ON (35).

Normalized Frequency	FIR		Sobel		IIR		Butterworth		DCT4		Geo.Mean	
	Uniform	Lena	Uniform	Lena	Uniform	Lena	Uniform	Lena	Uniform	Lena	Uniform	Lena
1.04	99.85%	100.00%	99.51%	99.74%	72.28%	90.09%	79.03%	100.00%	83.58%	98.06%	86.14%	97.50%
1.08	98.93%	99.97%	96.26%	93.75%	71.64%	90.50%	78.81%	100.00%	83.26%	98.45%	85.15%	96.46%
1.12	94.27%	98.82%	96.25%	93.62%	73.63%	88.25%	81.88%	84.87%	89.44%	99.56%	86.68%	92.84%
1.16	99.66%	99.91%	73.73%	93.62%	73.10%	89.92%	79.30%	84.23%	79.07%	99.44%	80.44%	93.23%
1.20	96.03%	99.90%	81.55%	81.52%	70.76%	75.67%	64.96%	84.50%	N/A*	N/A*	77.46%	84.95%
1.24	98.46%	99.32%	81.43%	81.67%	70.47%	76.12%	63.66%	84.23%	N/A*	N/A*	77.44%	84.92%
1.28	95.39%	99.29%	60.41%	78.24%	N/A*	N/A*	54.38%	75.15%	N/A*	N/A*	67.92%	83.58%
1.32	95.37%	98.75%	N/A*	N/A*	N/A*	N/A*	N/A*	N/A*	N/A*	N/A*	95.37%	98.75%

* Current frequency cannot be achieved in the traditional scenario. These points are excluded from the calculation of geometric means.

TABLE III
FREQUENCY SPEEDUPS IN OVERCLOCKING SCENARIO UNDER VARIOUS ERROR BUDGETS.

Error Budget %	FIR		Sobel		IIR		Butterworth		DCT4		Geo.Mean	
	Uniform	Lena	Uniform	Lena	Uniform	Lena	Uniform	Lena	Uniform	Lena	Uniform	Lena
0.05	4.76%	21.43%	6.82%	6.82%	0.95%	1.26%	12.40%	24.03%	0.72%	0.96%	3.07%	5.32%
0.5	19.05%	21.43%	13.64%	6.82%	0.63%	10.06%	24.03%	24.03%	0.48%	12.44%	4.52%	13.45%
1	19.05%	28.57%	13.64%	18.18%	10.06%	16.35%	24.03%	24.03%	0.48%	12.44%	7.86%	19.10%
5	19.15%	25.53%	18.18%	18.18%	0.54%	0.82%	0.63%	0.94%	7.66%	12.44%	3.91%	5.36%
10	19.15%	25.53%	10.64%	10.64%	0.54%	1.09%	6.92%	13.21%	4.91%	4.911%	5.19%	7.19%
20	19.15%	25.53%	5.77%	15.39%	8.70%	8.70%	3.26%	3.26%	6.70%	10.88%	7.32%	10.39%
50	15.69%	15.69%	10.53%	19.30%	42.50%	50.00%	46.74%	54.89%	15.06%	19.25%	21.8%	27.59%

scenario. Consequently, the area overheads of our approach are summarized in Table IV and Table V for both design goals. In our experiments, the area overhead is evaluated in terms of the number of LUTs and Slices in the FPGA technology, while similar metrics can also be employed if applying our approach into other hardware platforms. It can be seen that for a given frequency requirement with minimum possible errors, the geometric mean of area overhead ranges from $1.16\times$ and $1.15\times$ to $5.57\times$ and $4.46\times$ for LUTs and Slices, respectively.

An even smaller area overhead can be found for a specified error budget with the maximum clock frequencies. As seen in Table IV, the geometric mean of area overhead in terms of LUTs and Slices is $1.00\times \sim 4.67\times$ and $1.00\times \sim 3.70\times$ respectively. Notice that an area overhead of $1\times$ means both the overclocking scenario and the traditional scenario use the original precision, however the former achieves frequency speedups as seen in Table V, because a specific error budget can be tolerated.

In general, we can observe a relatively small area overhead of our approach, especially when frequency is initially increased or with a small error budget (i.e. $\sim 10\%$).

IX. CONCLUSION

The conclusion goes here.

ACKNOWLEDGMENT

The authors would like to acknowledge the support of the EPSRC (Grants EP/I020557/1 and EP/I012036/1).

REFERENCES

[1] K. Olukotun and L. Hammond, "The future of microprocessors," *Communications of the ACM*, vol. 3, no. 7, pp. 26–29, 2009.

[2] H. Esmailzadeh, E. Blem, R. S. Amant, K. Sankaralingam, and B. Doug, "Dark silicon and the end of multicore scaling," in *Int. Symp. Computer Architecture*, 2011, pp. 365–376.

[3] G. Constantinides, N. Nicolici, and A. Kinsman, "Numerical data representations for FPGA-based scientific computing," *IEEE Design Test of Computers*, vol. 28, no. 4, pp. 8–17, 2011.

[4] B. Colwell, "We may need a new box," *Computer*, vol. 37, no. 3, pp. 40–41, 2004.

[5] R. Moore, *Interval analysis*. Prentice-Hall Englewood Cliffs, NJ, 1966, vol. 60.

[6] L. H. De Figueiredo and J. Stolfi, "Affine arithmetic: concepts and applications," *Numerical Algorithms*, vol. 37, no. 1–4, pp. 147–158, 2004.

[7] A. Kinsman and N. Nicolici, "Bit-width allocation for hardware accelerators for scientific computing using sat-modulo theory," *Computer-Aided Design of Integrated Circuits and Systems, IEEE Trans. on*, vol. 29, no. 3, pp. 405–413, 2010.

[8] D. Boland and G. A. Constantinides, "Bounding variable values and round-off effects using handelman representations," *Computer-Aided Design of Integrated Circuits and Systems, IEEE Trans on*, vol. 30, no. 11, pp. 1691–1704, 2011.

[9] D. Lee, A. Gaffar, R. Cheung, O. Mencer, W. Luk, and G. Constantinides, "Accuracy-guaranteed bit-width optimization," *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, vol. 25, no. 10, pp. 1990–2000, 2006.

[10] S. I. Association, "International technology roadmap for semiconductors (ITRS)," 2007.

[11] T. Austin, V. Bertacco, D. Blaauw, and T. Mudge, "Opportunities and challenges for better than worst-case design," 2005, pp. 2–7.

[12] D. Ernst, N. Kim, S. Das, S. Pant, R. Rao, T. Pham, C. Ziesler, D. Blaauw, T. Austin, K. Flautner, *et al.*, "Razor: A low-power pipeline based on circuit-level timing speculation," in *Int. Symp. on Microarchitecture*, 2003, pp. 7–18.

[13] T. Austin, D. Blaauw, T. Mudge, and K. Flautner, "Making typical silicon matter with razor," *Computer*, vol. 37, no. 3, pp. 57–65, 2004.

[14] A. Uht, "Going beyond worst-case specs with TEAtime," *Computer*, vol. 37, no. 3, pp. 51–56, 2004.

[15] K. Keutzer and M. Orshansky, "From blind certainty to informed uncertainty," in *Proc. Int. workshop on Timing Issues in the Specification and Synthesis of Digital Systems*, 2002, pp. 37–41.

[16] A. Sampson, W. Dietl, E. Fortuna, D. Gnanapragasam, L. Ceze, and D. Grossman, "Enerj: Approximate data types for safe and general low-power computation," in *ACM SIGPLAN Notices*, vol. 46, no. 6. ACM, 2011, pp. 164–174.

TABLE IV
AREA OVERHEAD OF OUR APPROACH WITH RESPECT TO GIVEN FREQUENCY REQUIREMENTS

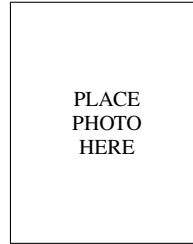
Normalized Frequency	FIR		Sobel		IIR		Butterworth		DCT4		Geo.Mean	
	LUTs	Slices	LUTs	Slices	LUTs	Slices	LUTs	Slices	LUTs	Slices	LUTs	Slices
1.04	1.24	1.05	1.12	1.22	1.26	1.42	1.16	1.11	1.13	1.00	1.18	1.15
1.08	1.24	1.05	1.03	1.08	1.26	1.42	1.16	1.11	1.11	1.15	1.16	1.16
1.12	1.24	1.05	1.03	1.08	1.26	1.42	1.37	1.54	2.30	2.17	1.38	1.40
1.16	2.95	2.56	1.03	1.08	1.61	1.71	1.99	2.03	4.60	2.89	2.14	1.95
1.20	2.95	2.56	1.66	1.39	2.77	3.48	3.17	3.50	N/A*	N/A*	2.56	2.57
1.24	4.42	3.15	4.00	3.55	7.87	6.73	6.92	5.25	N/A*	N/A*	5.57	4.46
1.28	4.42	3.15	4.00	3.55	N/A*	N/A*	6.92	5.25	N/A*	N/A*	4.97	3.89
1.32	4.42	3.15	N/A*	N/A*	N/A*	N/A*	N/A*	N/A*	N/A*	N/A*	4.42	3.15

* Current frequency cannot be achieved in the traditional scenario. These points are excluded from the calculation of geometric means.

TABLE V
AREA OVERHEAD OF OUR APPROACH WITH RESPECT TO GIVEN ERROR BUDGETS.

Error Budget %	FIR		Sobel		IIR		Butterworth		DCT4		Geo.Mean	
	LUTs	Slices	LUTs	Slices	LUTs	Slices	LUTs	Slices	LUTs	Slices	LUTs	Slices
0.05	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
0.5	1.02	1.15	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.03
1	1.02	1.15	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.03
5	1.24	1.05	1.00	1.00	1.26	1.42	1.65	1.11	1.00	1.00	1.13	1.11
10	1.24	1.05	1.12	1.22	1.26	1.42	1.65	1.11	1.11	1.15	1.18	1.18
20	1.24	1.05	1.03	1.08	1.26	1.42	6.92	5.25	3.45	2.17	2.07	1.79
50	2.95	2.56	4.00	3.55	7.87	6.73	6.92	5.25	3.45	2.12	4.67	3.70

- [17] H. Esmailzadeh, A. Sampson, L. Ceze, and D. Burger, "Architecture support for disciplined approximate programming," in *Proc. Int. Conf. Architectural Support for Programming Languages and Operating Systems*, 2012, pp. 301–312.
- [18] Z. Kedem, V. Mooney, K. Muntimadugu, and K. Palem, "An approach to energy-error tradeoffs in approximate ripple carry adders," in *Int. Symp. on Low Power Electronics and Design*, 2011, pp. 211–216.
- [19] S. Lu, "Speeding up processing with approximation circuits," *IEEE Computer*, vol. 37, no. 3, pp. 67–73, 2004.
- [20] P. Kulkarni, P. Gupta, and M. Ercegovac, "Trading accuracy for power with an underdesigned multiplier architecture," in *Int. Conf. on VLSI Design*, 2011, pp. 346–351.
- [21] V. Gupta, D. Mohapatra, A. Raghunathan, and K. Roy, "Low-power digital signal processing using approximate adders," *Computer-Aided Design of Integrated Circuits and Systems, IEEE Trans. on*, vol. 32, no. 1, pp. 124–137, 2013.
- [22] Altera, "Cyclone device handbook," 2008.
- [23] Xilinx, "Virtex-6 FPGA configurable logic block user guide," 2009.
- [24] J. Rabaey, A. Chandrakasan, and B. Nikolic, *Digital integrated circuits: a design perspective (2nd edition)*. Prentice-Hall, 2003.
- [25] Xilinx, "Virtex-6 FPGA clocking resources user guide," 2011.
- [26] B. Gojman, S. Nalmela, N. Mehta, N. Howarth, and A. DeHon, "GROK-LAB: generating real on-chip knowledge for intra-cluster delays using timing extraction," in *Proc. Int. Symp. on Field Programmable Gate Arrays*, 2013, pp. 81–90.
- [27] H. Wong, L. Cheng, Y. Lin, and L. He, "FPGA device and architecture evaluation considering process variations," in *Proc. Int. Conf. on Computer-Aided Design*, 2005, pp. 19–24.



Author Biography text here.

Author Biography text here.

Author Biography text here.