

Evaluation of Design Trade-offs for Approximate Arithmetic Operators

Kan Shi and George A. Constantinides
 Department of Electrical and Electronic Engineering
 Imperial College London
 London, UK
 Email: k.shi11, g.constantinides@imperial.ac.uk

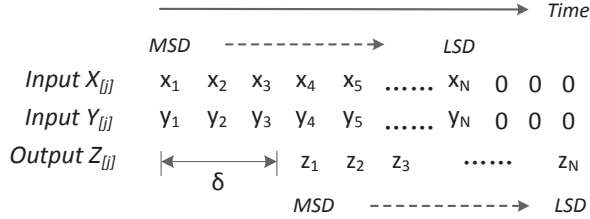


Fig. 1. Dataflow in digit-serial online arithmetic, in which both inputs and outputs are processed from the MSD to the LSD. δ denotes the online delay.

Abstract—The abstract goes here.

Keywords—approximate computing; datapath; computer arithmetic;

I. INTRODUCTION

II. BACKGROUND

A. Approximate Datapath Design

B. Online Arithmetic

As an alternative form of computer arithmetic, online arithmetic has been used in numerous applications such as signal processing and control algorithms [1], [2]. Online arithmetic was originally designed for digit-serial operation, as illustrated in Fig 1. It can be seen that in order to generate the first output digit, δ digits of inputs are required, where δ is called the “online delay”. δ is normally a small constant, which is independent of the precision. For ease of discussion, for the rest of this paper, the input data is assumed to be fixed point numbers in the range $(-1, 1)$. Based on this premise, the online representation of N -digit operands and result at iteration j are given by (1), where $j \in [-\delta, N - 1]$ and r denotes the radix [3].

$$X_{[j]} = \sum_{i=1}^{j+\delta} x_i r^{-i}, \quad Y_{[j]} = \sum_{i=1}^{j+\delta} y_i r^{-i}, \quad Z_{[j]} = \sum_{i=1}^j z_i r^{-i} \quad (1)$$

MSD-first operation is possible only if a redundant number system is used. Normally there are two most commonly used redundant number representations: carry-save (CS) [4] and signed-digit (SD) [5]. With SD representation, each digit is represented using a redundant digit set $\{-a, \dots, -1, 0, 1, \dots, a\}$, where $a \in [r/2, r - 1]$. In comparison, the standard non-redundant representation only uses a digit set $\{0, \dots, r - 1\}$. Thus a standard number corresponds to several possible redundant representations. For example, the binary number 0.011

can be represented in SD form as $0.1\bar{1}1$, $0.10\bar{1}$ or 0.011 among many other possible representations.

Due to the redundancy, the MSDs of the result can be calculated using partial information from both inputs. Then the value of the number can be revised using the subsequent digits, because each number has multiple representations.

III. COMPARISON OF DIFFERENT ADDER STRUCTURES

A. Ripple Carry Adder

Adders serve as a key building block for arithmetic operations. In general, the ripple carry adder (RCA) is the most straightforward and widely used adder structure. As such, in our previous work we proposed probabilistic models of overclocking errors for RCA.

For an n -bit RCA, it is composed of n serial-connected full adders (FA). Typically the maximum frequency of RCA is determined by the longest carry propagation. Under the assumption that the carry propagation delay of each FA is a constant value μ , the critical path of the RCA is: $\mu_{RCA} = n\mu$. For the sampling period T_S , it follows that if $T_S \geq \mu_{RCA}$, correct result will always be sampled; otherwise if $T_S < \mu_{RCA}$, intermediate be sampled and potentially generating errors. In the previous work, the mean value of overclocking error is modeled as given in (2), where coefficient b is given in (3) and it determines the maximum length of error-free carry propagation in bits.

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

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

On the other hand, a specific timing target can be met by truncating the word-length of RCA while timing violation is not permitted. This will result in truncation error for most data. The mean value of truncation error is also modeled as given in (4) [xxx], where parameters k and n denote the word-length of RCA before and after truncation, respectively.

$$E_O = \begin{cases} 2^{-n} - 2^{-k}, & \text{if } n < k \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

In this paper we consider two design scenarios for RCA: one is the overclocking scenario where timing violation is

allowed to happen while maintaining the original word-length; the other is the trad scenario by truncating RCA word-length to meet timing. For a given f_S , in the first design scenario we have $n = k$. In the second design scenario we have $n = b - 1$ to prevent timing violation while minimizing the value of E_O . Therefore the comparison between the two scenarios in mean error value is given in (5) by updating (2) and (4) respectively. As seen in (5), in theory the mean value of error in the overclocking scenario (E_{oc}) is two times smaller than that in the traditional scenario (E_{trad}).

$$\frac{E_{oc}}{E_{trad}} = \frac{2^{-b} - 2^{-n-1}|_{n=k}}{2^{-n} - 2^{-k}|_{n=b-1}} = \frac{1}{2} \quad (5)$$

B. Carry Select Adder

Although smaller value of mean error can be achieved in the overclocking scenario, it costs extra area because the full precision is maintained. Instead, alternative adder structures, such as carry select adder (CSA), are originally designed to trade silicon area for low latency. In a CSA, the carry chain is divided into multiple overlapped stages, and each stage contains two RCAs and two multiplexers, as shown in Fig. xxx(a) and Fig. xxx(b). For a given input, two additions are performed simultaneously within a single stage where the carry input is zero and one separately. 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, multiplexers are area-expensive to implement with FPGA technology.

C. Online Adder

In addition to standard binary arithmetic, alternative forms of computer arithmetic was studied to boost performance. For example, adder with online arithmetic is also designed for low latency at the cost of more silicon area. The structure of a digit-parallel online adder (OA) where all signals represented with SD numbers of digit set $\{-1, 0, 1\}$ is shown in Fig. 2. The module “3:2” denotes a 3:2 compressor, which takes three inputs and generates two outputs, and is logically equivalent to a full adder (FA). A major advantage of the redundant number system over the standard ripple-carry based arithmetic is that the propagation of carry is eliminated, resulting in a precision-independent computation time for addition. As labelled in Fig. 2, ideally the computation delay of this adder is only two FA delays for any operand word-length, at the expenses of one extra FA for each digit of operands.

D. Design Trade-offs for Adders

In this section, we initially demonstrate the benefits of CSA and OA in accuracy and performance, as well as the area overhead in comparison to RCA. As an example, the maximum word-length with respect to a variety of operating frequencies for RCA, CSA and OA is illustrated in Fig. 3. Notice that in this case, overclocking error is not allowed to happen and the conventional design scenario is evaluated for all structures.

Correspondingly, the area consumed by each adder structure in terms of look-up-tables (LUTs) in FPGA is depicted in Fig. 4. Clearly we see that for a given frequency, CSA costs up to $4.5\times$ (4 stages) and $3.1\times$ (2 stages) area than RCA,

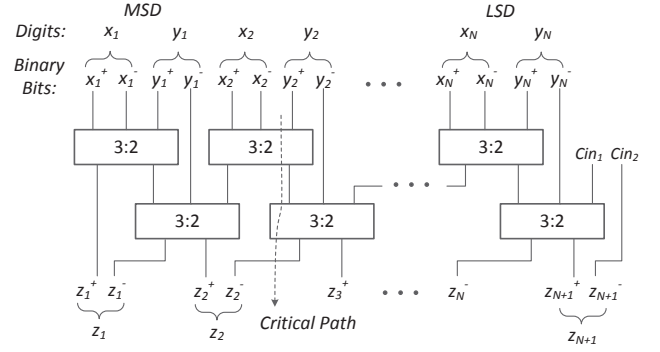


Fig. 2. An N -digit binary digit-parallel online adder. Both inputs and outputs are represented using SD representation. “3:2” denotes a 3:2 compressor.

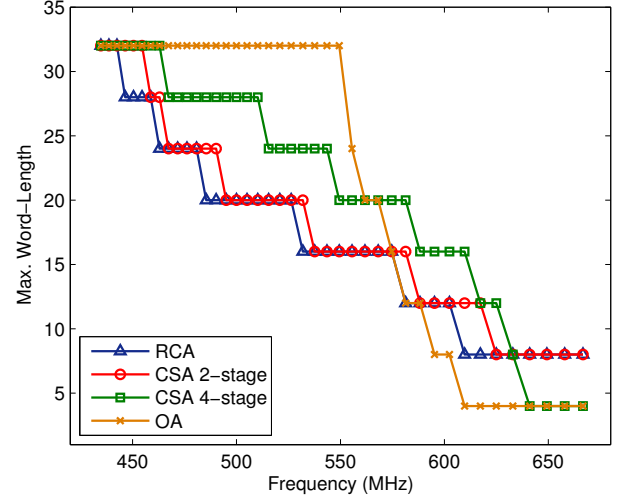


Fig. 3. The maximum word-length of different adder structures.

and OA is up to $3.8\times$ larger than RCA for some frequency values. Therefore, we incorporate silicon area as another design metric besides accuracy and performance, and evaluate the optimum adder structure under various design trade-offs. In this paper, we examine both the overclocking scenario and the conventional scenario for RCA, while only considering the conventional scenario for other adder structures.

In general, the evaluation could be of interest to a circuit designer in several ways. First, suppose the algorithm designer wish the circuit to run at a given frequency within a given area budget, while achieving the minimum output error. Decisions must then be made on which adder structure achieves this minimum, and whether the circuit should be overclocked.

As an example slice through the design space, we record the mean relative error (MRE) with respect to a range of operating frequencies for different design scenarios when the area budget is set to xxx. MRE can be calculated as given in (xxx). Notice that in Fig. xxx, the optimum adder structure is labeled.

Under this given area budget, both CSA and OA cannot be implemented with full precision, and this leads to large truncation errors

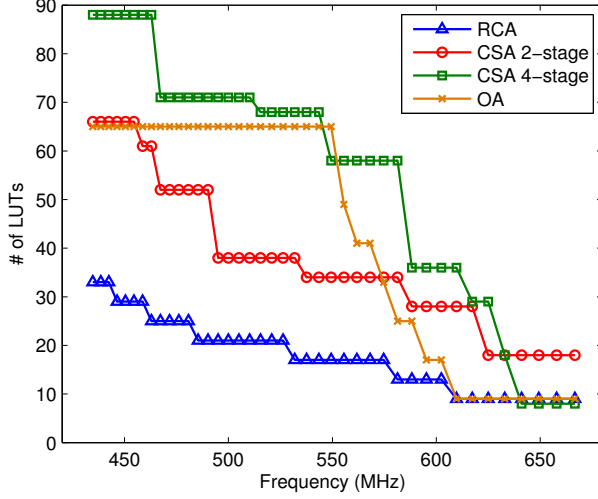


Fig. 4. The maximum word-length of different adder structures.

Algorithm 1 Digit Parallel Online Multiplication

```

1: Initialization:  $P_{[0]} = 0$ 
2: for  $j = 1, 2, \dots, N$  do
3:    $Xy_j \leftarrow X \cdot y_j$ 
4:    $W_{[j]} \leftarrow P_{[j-1]} + Xy_j$ 
5:    $z_j \leftarrow \text{sel}(W_{[j]})$ 
6:    $P_{[j+1]} \leftarrow r(W_{[j]} - Z_{[j]})$ 
7: end for
8:  $Z_{[N+1:2N]} \leftarrow \text{frac}(W_{[N]})$ 

```

IV. COMPARISON OF DIFFERENT MULTIPLIER STRUCTURES

A. Standard Binary Multiplier

B. Online Multiplier

1) *Algorithm:* Similar to online addition, online multiplication can also be implemented in a unrolled digit-parallel manner, of which the optimized algorithm is shown in Algorithm 1 [xxx]. In this algorithm, both inputs and outputs are N -digit numbers as given in (1). For a given iteration j , the product digit z_j is generated MSD-first through a selection function $\text{sel}()$. For any radix r and chosen digit set, there exists an appropriate selection method and a value of δ which ensure convergence. As the binary radix is used most commonly in computer arithmetic, we keep $r = 2$ throughout this paper with the corresponding redundant digit set $\{\bar{1}, 0, 1\}$. In this case $\text{sel}()$ is given by (6) [6].

$$\text{sel}(W_{[j]}) = \begin{cases} 1 & \text{if } W_{[j]} \geq \frac{1}{2} \\ 0 & \text{if } -\frac{1}{2} \leq W_{[j]} < \frac{1}{2} \\ \bar{1} & \text{if } W_{[j]} < -\frac{1}{2} \end{cases} \quad (6)$$

The parallel implementation of Algorithm 1 is illustrated in Fig. 5 where $N = 4$. The details of implementation are addressed in our previous work [xxx].

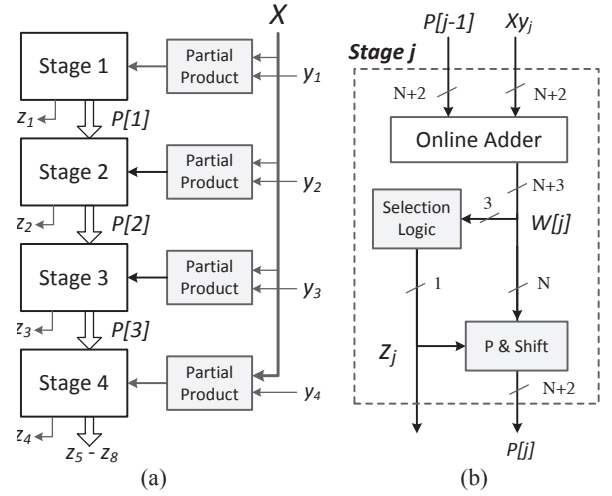


Fig. 5. (a) Structure of a 4-digit online multiplier. (b) Structure of one stage. The word-length of all signals are labeled in digits, N is the word-length of input signals.

2) Structure Optimization for Half Precision Results:

In actual applications, normally the multiplier is connected with other arithmetic operators. If the multiplication result is utilized for subsequent operations, and a consistent word-length is used throughout the system, then only the most significant half of the product is required. In a conventional multiplier with standard binary arithmetic, this is achieved by either truncating or rounding the least significant half of the products. However, both the computation time and the structure remains unchanged, because results are generated from LSDs.

In comparison, OM offers the flexibility to simplify the structure corresponding to the required precision. This is possible because in an OM, the product digits are generated initially from the MSD, and there is no carry propagation from the LSD to the MSD with the employment of the redundant number system. This will potentially lead to a more area efficient design. For example, the modified structure diagram of a 4-digit OM is illustrated in Fig. 6(a). The word-length of signals within each stage can be correspondingly reduced, as shown in Fig. 6(b). Notice that instead of keeping identical word-length throughout the stages (in Fig. 5), in the modified structure the signal word-lengths depend on the stage number j and hence fewer bits of signals are needed for LSD stages.

V. CONCLUSION

ACKNOWLEDGMENT

REFERENCES

- [1] R. Galli and A. F. Tenca, "Design and evaluation of online arithmetic for signal processing applications on FPGAs," in *Proc. SPIE Advanced Signal Processing Algorithms, Architectures and Implementations*, Aug 2001, pp. 134–144.
- [2] M. Dimmler, A. Tisserand, U. Holmbeg, and R. Longchamp, "On-line arithmetic for real-time control of microsystems," *IEEE/ASME Trans. on Mechatronics*, vol. 4, no. 2, pp. 213–217, Jun 1999.
- [3] M. D. Ercegovac and T. Lang, *Digital arithmetic*. Morgan Kaufmann, 2003.
- [4] T. Kim, W. Jao, and S. Tjiang, "Circuit optimization using carry-save-adder cells," *IEEE Trans. Computer-Aided Design of Integrated Circuits and Systems*, vol. 17, no. 10, pp. 974–984, 1998.

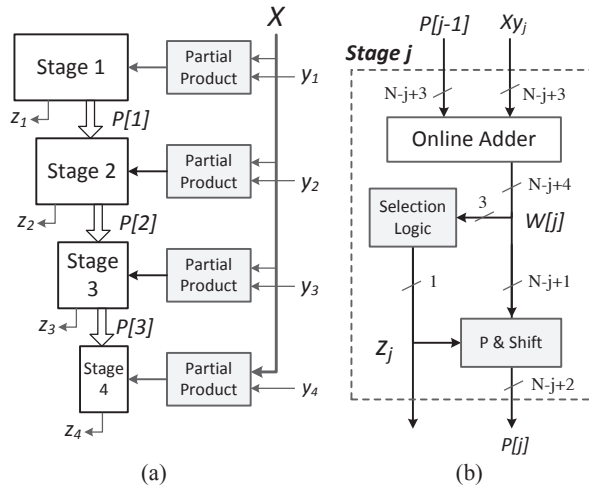


Fig. 6. (a) A 4-digit online multiplier with half precision results. (b) Structure of stage j , with the word-length of all signals labeled.

- [5] A. Avizienis, "Signed-digit number representations for fast parallel arithmetic," *IRE Trans. on Electronic Computers*, vol. EC-10, no. 3, pp. 389–400, 1961.
- [6] K. S. Trivedi and M. Ercegovac, "On-line algorithms for division and multiplication," *IEEE Trans. on Computers*, vol. 26, pp. 161–167, 1975.