#### A Novel Data Format for Approximate Arithmetic Computing

#### Mingze Gao, Qian Wang, Akshaya Nagendra, and Gang Qu

Electrical and Computer Engineering Department and Institute of Systems Research University of Maryland, College Park

### Introduction

- # What is Approximate Computing (AC)?
  - Approximate (error) vs. accurate
- # Why we need AC?
  - Power/energy efficiency
- # Why AC works?
  - Many of the applications are error-tolerable, e.g. Machine Learning, Image/Signal Processing
  - Disable partial computation
- # AC at different level
  - Arithmetic, Software, Compiler, Architecture, Memory, and Circuit

Sec Lab Dr. Gang Qu (gangqu@umd.edu)

#### Approximate Arithmetic

#### # Observation:

 Least Significant Bits (LSB) have much less contribution than Most Significant Bits (MSB) to the overall quality of the result.

#### # Approach:

Compute accurately on MSB

Apply approximation on LSB



#### Approximate Arithmetic

#### # Example: Compute S = A + B A = 0011 1010 0001 10002 B = 0101 1011 1011 10002

- Build a 16-bit approximate adder that[1]:
   8-bit accurate adder for high 8 bits
  - 8 OR gates for low 8 bits

shSec Lab

Error = 0.102% !

[1] H. R. Mahdiani, A. Ahmadi, S. M. Fakhraie, and C. Lucas, "Bio-inspired imprecise computational blocks for efficient VLSI implementation of soft-computing applications," Circuits Syst. I Regul. Pap. IEEE Trans., vol. 57, no. 4, pp. 850–862, 2010.



#### # We need a better approximate adder!

 $0000 \ 0000 \ 1111 \ 1001_2 = 249$ 

 $0000 \ 0000 \ 1011 \ 1000_2 = 184$ 

0000 0000 0111 1001<sub>2</sub> = 121 **Error = 18.4\%** 

 $A = 0000 \ 0000 \ 0101 \ 1000_2$  $B = 0000 \ 0000 \ 1011 \ 1000_2$ 

#### # However, what if the data is

#### Approximate Arithmetic

#### Approximate Arithmetic

- # What we have learned:
  - "static" approximate adder vs. "dynamic" data
- #Existing solutions:
  - Build additional discriminant circuit inside the approximate adder
- # Drawbacks:
  - Fail to deliver significant power savings
     Less accurate for multiplication



## Approximate Integer Format

# Contribution: a novel Approximate Integer Format (AIF) and the corresponding computation mechanisms # Desired properties of an ideal AIF From "static" to "dynamic" Cut-off the bitwidth of the operands Suitable for all arithmetic operations Provable error bound Applicable to fixed point arithmetic



-- Valid Block

- # AIF is based on the segmentation of operands.
  - An n-bit positive integer N is segmented into [n/k] blocks with k bits per block.
  - Example: n = 16, k = 4, there are 4 blocks.

A = 0000 1010 0001 10002

Definition 1: A valid block in a positive number is a block that has at least one '1' before or inside it.
A = 0000 1010 0001 00002

#### -- Sentinel Bits

Sentinel bits are used to truncate and round the less important bits to reduce bit-length of the operands
Definition 2: The ith sentinel bit st[i] of a number is defined as

 $st[i] = \begin{cases} 1, & block \ i \ is \ a \ valid \ block \\ 0, & block \ i \ is \ an \ invalid \ block \end{cases}$ 



#### --Precision Control

- Definition 3: The precision control 'pc' is the number of valid blocks in the number, from the leftmost one, that will be used in the computation.
- #Example:

 $1500_{10} = 0000 \ 0101 \ 1101 \ 1100_2 \ 800_{10} = 0000 \ 0011 \ 0010 \ 0000_2$ 

Both have 3 valid block, st = 0111
 If pc = 2, 2 blocks of each operand will be selected 150010 = 0000 0101 1101 11002 80010 = 0000 0011 0010 00002

-Rounding

When we use sentinel bits to truncate the valid blocks, rounding is needed

- # Two rounding techniques:
  - Classic rounding
    - For multiplication and division
  - Efficient rounding
    - For addition and subtraction



# Classic Rounding

# Definition 4: The classic rounding of a number N at the ith LSB means adding the ith bit to the (i+1)th bit setting ith bit and bits to its right to zero  $# Example: N = 263_{10} = 0000 0001 0000 0111_2$ From the 3<sup>rd</sup> least significant bit • N = 0000 0001 0000 10002 From the 4th bit **N = 0000 0001 0000 0000**2



# Efficient Rounding

Definition 5: The efficient rounding of A+B at the ith bit is

Atrunc + Btrunc + Cinround

Atrunc and Btrunc are obtained by truncating the i least significant bits from A and B
 Cinround = (Ai&Bi), AND ith bits of A and B



## Efficient Rounding

To compute S = A + BA = 0011 1010 1001 10002 B = 0000 1011 1011 10002

S' = Atrunc + Btrunc + Cinround 0011 1010 + 1011 + 1

0100 0110

Truncate A and B Atrunc = 0011 1010 0000 00002 Btrunc = 0000 1011 0000 00002

-- Example

Compute round-off carry in Cinround = A7 & B7 = 1

S using efficient rounding: 0100 0110 0000 00002 = 1792010 Accurate S: 0100 0110 0101 00002 = 1800010



#### Approximate Integer Format

Given a 4-block operand A= b<sub>3</sub>b<sub>2</sub>b<sub>1</sub>b<sub>0</sub>.
 Only five possible values of A's sentinel bits st<sub>a</sub>: 0000, 0001, 0011, 0111, 1111.

For the first four cases, the data A will be stored in following format:

$$\mathbf{st}_{\mathbf{a}}$$
  $\mathbf{b}_{\mathbf{2}}$   $\mathbf{b}_{\mathbf{1}}$   $\mathbf{b}_{\mathbf{0}}$ 

For the last case of 1111, A will be stored

$$st_a b_3 b_2 b_1$$



as:

# AIF Arithmetic



- **Compute the sentinel bits of the result**  $S: st_s = st_A | st_B$
- Truncate ith to (i-pc+1)th blocks of A and B to obtain A' and B', respectively
  - Suppose the leftmost '1' in st<sub>s</sub> is in st[i], and we plan to pick pc valid blocks
- # Compute S' = A' + B' and Cout
- # Update  $st_s$  by  $st_s[i+1] = C_{out}$
- Reformulate S in AIF using st<sub>s</sub> and S'. Padding O's if necessary



## AIF Arithmetic -- Addition Example

Original data A = 0011 1010 0001 10002 B = 0000 1011 1011 10002

Compute S'  $0011 \ 1010 + 0$   $+ \ 1011 + 1$  $= \ 0100 \ 0110$ 

nSec Lab



		$\mathbf{\nabla}$	200	22.
C	om	pu'	te :	STs
÷	· · · ·		111	1
1.1 -		N. 13.	111	1
1.3	1	14 .1	111	1. 2
0		3.1		1.1
	10	1.17	111	1
			┶┶┶╺	



Reformulate S in AIF: 1111 0100 0110 00002 = 1792010 Accurate S: 0100 0101 1101 00002 = 1787210

# AIF Arithmetic --Multiplication

- Round the leftmost pc valid blocks of A and B into A' and B'
  Compute S' = A' \* B'
  Compute sentinel bits st<sub>S</sub> using st<sub>A</sub> and st<sub>B</sub> and carry out
  Shift and reformulate S in AIF using S'
  - and st<sub>s</sub>. Padding O's if necessary



#### Error analysis

# Let rounding error of A and B are era and erB, respectively.
# Error of AIF based addition:
■ Eradd = 2\*max(era, erB)
# Error of AIF based multiplication:
■ era + erB + era\*erB
■ era<<1, erB<<1, Ermul ≈ era + erB</li>



# Negative AIF

- # Deal with negative integer
- Cannot use previous equation to compute st
   Solution:
  - Re-define the valid block
  - A valid block in a negative number is a block that has at least one '0' before or inside it.
  - Replace logic OR with logic AND when computing st
  - Arithmetic operations remains the same



# High level programming language

- # Introduce appropriate instructions and data type.
- #Incorporate it in compiler and ISA
- # Example: define apxint as the AIF data

Dr. Gang Qu (gangqu@umd.edu)



type

# Compute in Caution

- Condition criterion, e.g. if, while condition
- Data value that the result is very sensitive to
- Functions that have periodical property, e.g. sin, cos, and modulo operation



### Experiment Results: HW Cost

Overhead Comparison of Arithmetic Units and The Sentinel Bits Computing Circuits

1		cells	area	power(nW)
	8-bit checker	10	23.93	61475.68
**	16-bit checker	17	39.89	132145.68
-	8-bit adder	91	212.12	752614.84
	16-bit adder	230	322.33	2234652.24
1	32-bit adder	498	1116.93	4818821.94
	8-bit multiplier	377	1037.62	2829745.1
-	16-bit multiplier	1406	4208.68	10815807.06
	32-bit multiplier	4916	15126.01	34033690.3

#### Fibonacci Sequence: Accuracy

#First 40 elements in Fibonacci Sequence

- A number is sum of its previous two: 1, 1, 2, 3, 5, 8, 13, 21, 34, ...
- Test the error propagation

#	pc=2	pc=3	pc=4	#	pc=2	pc=3	pc=4	#	pc=2	pc=3	pc=4
1~13	0	0	0	22	-0.02732	3.49E-05	0	31	-0.02291	-0.00035	2.30E-06
14	-0.00984	0	0	23	-0.02692	0	0	32	-0.02267	-0.00059	-8.51E-06
15	-0.01317	0	0	24	-0.02707	1.33E-05	0	33	-0.02276	-0.0005	-4.38E-06
16	-0.01691	0	0	25	-0.02912	0.000404	8.24E-06	34	-0.02273	-0.00053	-5.96E-06
17	-0.01858	0	0	26	-0.03225	0.000336	1.02E-05	35	-0.02274	-0.00052	-5.36E-06
18	-0.01985	0	0	27	-0.0375	0.000362	9.44E-06	36	-0.02274	-0.00052	-5.59E-06
19	-0.02173	-0.00044	0	28	-0.03947	0.000352	9.72E-06	37	-0.0328	-0.0001	1.05E-06
20	-0.02905	0.000183	0	29	-0.04118	0.000356	9.61E-06	38	-0.03621	-0.00046	-5.53E-06
21	-0.02625	-5.65E-05	0	30	-0.04205	0.000354	9.66E-06	39	-0.04003	-0.00064	-3.02E-06
								40	-0.04174	-0.00077	-3.98E-06



# **Real Life Application**

Which one is the original image?

**PSNR = 41.76** 





#### PSNR = 89.64

--IDC

PSNR =116.29





original image



Dr. Gang Qu (gangqu@umd.edu)

25

### Real Life Applications

and the second second

-Quality

				1	
AIF	IDCT	KNN	FFT	Kmeans	SVM
modules					
32_8_2	41.7579	92.9%	0.0081	1.125%	60.94%
32_8_3	64.6398	93.2%	2.79E-04	0.125%	85.11%
32_8_4	89.8324	93.1%	1.61E-05	0	85.98%
32_8_5	112.0872	93.2%	3.02E-06	0	85.59%
32_8_6	116.2944	93.2%	7.11E-08	0	86.42%
baseline	116.2949	93.2%	0	0	86.42%
Error		Classification		miss-	Classification
Metric	PSNR	accuracy	ARES	clustered%	accuracy

and the second second second second second



#### Final Result --Power Savings Normalized power consumptions 1.2 1 0.8 0.6 0.4 0.2 n idct knn fft kmeans sym 32 8 2 **32\_8\_3 32\_8\_4 32\_8\_5 32 8 6** accurate shSec Lab Dr. Gang Qu (gangqu@umd.edu)

# Thank you!

