

Hardware implementation of linear algebra operators for small floating point formats

Orégane Desrentes

- 1 Kalray
- 2 Floating point linear algebra
- 3 8 bits formats



Kalray



A complete offer or data-intensive applications

**A UNIFIED
VALUE PROPOSITION**



KALRAY
THE POWER OF MORE



PCIe Cards



Reference Solutions



MPPA® DPU
Manycore Processors



Data Management & Storage
Software Solutions



A complete offer for data-intensive applications

Two application domains:

Data Center acceleration

Computing



A complete offer or data-intensive applications

Two application domains:

Data Center acceleration

- Compression and decompression
- Encryption and decryption
- Erasure coding
- De-duplication

Computing



A complete offer or data-intensive applications

Two application domains:

Data Center acceleration

- Compression and decompression
- Encryption and decryption
- Erasure coding
- De-duplication

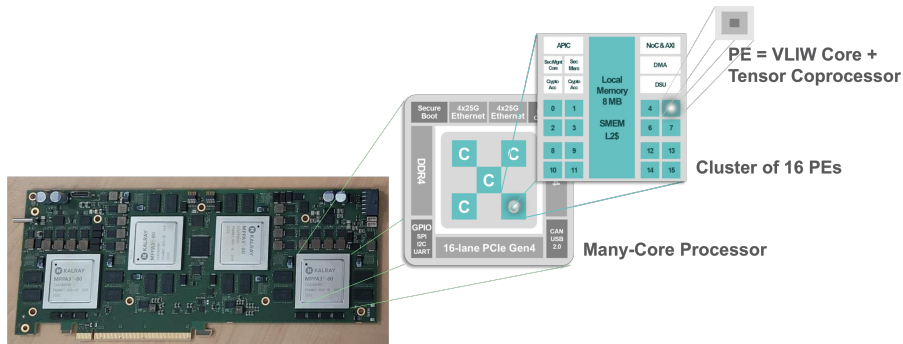
Computing

- Machine Learning
- Computer vision
- Pre/post processing
- Signal processing



Kalray MPPA® scalable many-core architecture

3rd-gen MPPA® processor: in TSMC 16nm technology, up to 1.2 GHz



Multiple Processors per Card



MPPA3 V2 Coolidge™ processing element (PE)

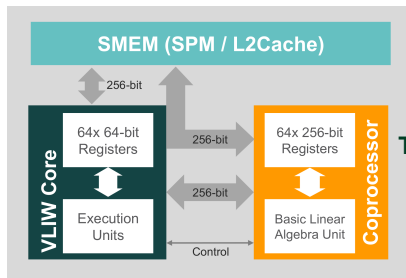
6-issue 64-bit VLIW core with a tightly-coupled tensor coprocessor

VLIW Core

- Scalar 32-bit and 64-bit INT & FP
- 8×8 -bit, 4×16 -bit, 2×32 -bit SIMD
- 128-bit 256-bit SIMD operations by bundling multiple instructions
- 256-bit load/store unit with masking

Tensor Coprocessor

- Matrix multiply-add on 4×4 tiles
- 512-bit multiply and add operands
- Matrix zip/unzip & transpose
- 256-bit load/store unit with masking
- Blocks of 256-bit registers used as circular buffer or as lookup table



MPPA3 V2 Coolidge™ processing element (PE)

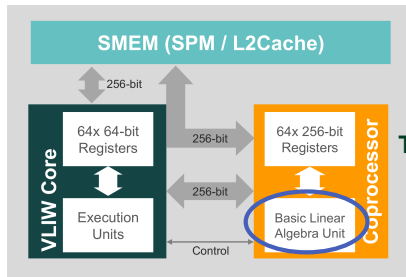
6-issue 64-bit VLIW core with a tightly-coupled tensor coprocessor

VLIW Core

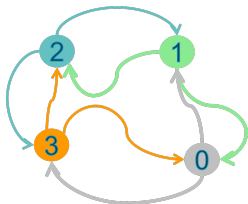
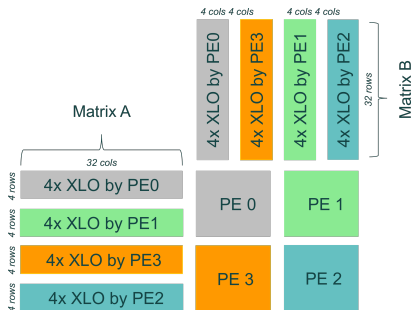
- Scalar 32-bit and 64-bit INT & FP
- 8×8 -bit, 4×16 -bit, 2×32 -bit SIMD
- 128-bit 256-bit SIMD operations by bundling multiple instructions
- 256-bit load/store unit with masking

Tensor Coprocessor

- Matrix multiply-add on 4×4 tiles
- 512-bit multiply and add operands
- Matrix zip/unzip & transpose
- 256-bit load/store unit with masking
- Blocks of 256-bit registers used as circular buffer or as lookup table



Coprocessor collective tensor operations



- PE operation: INT8.32
 $(4 \times 16) \cdot (16 \times 4) + = (4 \times 4)$
- Macro-scheme executed by 4 PEs
 - $8 \times$ 256-bit memory loads (XLO) per PE
 - $8 \times$ 256-bit data exchanges per PE
 - $8 \times$ matrix multiply-add operations per PE
- Matrix A and B are loaded by quarter by each PE which exchange one quarter with 2 different PEs
- Kernel for INT8.32:
 $(16 \times 32) \cdot (32 \times 16) + = (16 \times 16)$



Coprocessor Matrix multiply-accumulate operations

8-bit to 32-bit int matrix multiply-add: $(4 \times 16)_{\text{int}8} \cdot (16 \times 4)_{\text{int}8} + = (4 \times 4)_{\text{int}32}$

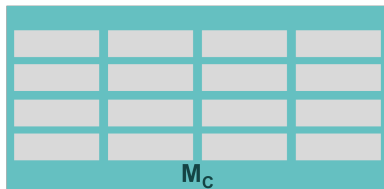
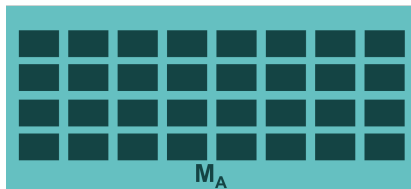
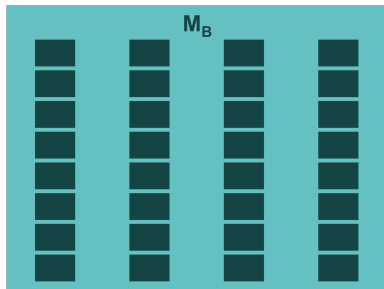
16-bit to 32-bit FP matrix multiply-add: $(4 \times 8)_{\text{FP}16} \cdot (8 \times 4)_{\text{FP}16} + = (4 \times 4)_{\text{FP}32}$

Signature

- 512-bit \times 512-bit $+ =$ 512-bit
- 256-bit register-pair multiplicands
- 256-bit register-pair accumulator

Performances

- 256 MADD eq. per cycle, 512 ops/c
- 128 FMA eq. per cycle, 256 flops/c
- 50 TOPS @1.2 GHz for 80 cores
- 25 TFLOPS @1.2 GHz for 80 core



Floating point linear algebra



Floating-Point numbers

- Computer representation for real numbers



Floating-Point numbers

- Computer representation for real numbers
- \mathbb{R} is infinite in range and precision



Floating-Point numbers

- Computer representation for real numbers
- \mathbb{R} is infinite in range and precision
- A computer is finite, in base 2



Floating-Point numbers

- Computer representation for real numbers
- \mathbb{R} is infinite in range and precision
- A computer is finite, in base 2

$$(-1)^S \times 2^E \times 1.F$$

S , E , F are stored in binary, in a finite format.



Floating-Point numbers

- Computer representation for real numbers
- \mathbb{R} is infinite in range and precision
- A computer is finite, in base 2

$$(-1)^S \times 2^E \times 1.F$$

S , E , F are stored in binary, in a finite format.

For floats: $w_S = 1$, $w_E = 8$, $w_F = 23$



Floating-Point special values

- $+0, -0$



Floating-Point special values

- $+0, -0$
 - $S = \text{sign}$
 - $E = 0$
 - $F = 0$



Floating-Point special values

- $+0, -0$
 - $S = \text{sign}$
 - $E = 0$
 - $F = 0$
- $+\infty, -\infty$



Floating-Point special values

- $+0, -0$
 - $S = \text{sign}$
 - $E = 0$
 - $F = 0$
- $+\infty, -\infty$
 - $S = \text{sign}$
 - $E = 2^{W_E} - 1$ (max number)
 - $F = 0$



Floating-Point special values

- $+0, -0$
 - $S = \text{sign}$
 - $E = 0$
 - $F = 0$
- $+\infty, -\infty$
 - $S = \text{sign}$
 - $E = 2^{W_E} - 1$ (max number)
 - $F = 0$
- NaN



Floating-Point special values

- $+0, -0$
 - $S = \text{sign}$
 - $E = 0$
 - $F = 0$
- $+\infty, -\infty$
 - $S = \text{sign}$
 - $E = 2^{w_E} - 1$ (max number)
 - $F = 0$
- NaN ($2^{w_F+1} - 2$)
 - $S = \text{any}$
 - $E = 2^{w_E} - 1$
 - $F \neq 0$



Floating-Point special values

- $+0, -0$
 - $S = \text{sign}$
 - $E = 0$
 - $F = 0$
- $+\infty, -\infty$
 - $S = \text{sign}$
 - $E = 2^{w_E} - 1$ (max number)
 - $F = 0$
- NaN ($2^{w_F+1} - 2$)
 - $S = \text{any}$
 - $E = 2^{w_E} - 1$
 - $F \neq 0$



Subnormals



Subnormals



Subnormals



Subnormals



Why we need them

Avoid numbers being flushed to 0 abruptly: gradual underflow.
Guarantees if $x \neq y$ then $x - y \neq 0$ and other useful properties.



Subnormals



Why we need them

Avoid numbers being flushed to 0 abruptly: gradual underflow.
 Guarantees if $x \neq y$ then $x - y \neq 0$ and other useful properties.

Why they are annoying

They have less precision than w_F .

They are encoded differently: $(-1)^S \times 2^E \times 0.F$

Where the first significant bit ? $0.F = 0.000000011001100$

⇒ They used to be treated in micro-code but now we do subnormal hardware



Short and biased history of floating point units

Short and biased history of floating point units

- In the 70s, adders and multipliers
 - $R = \circ(X + Y)$
 - $R = \circ(X \times Y)$



Short and biased history of floating point units

Short and biased history of floating point units

- In the 70s, adders and multipliers
- In 1985, IEEE-754 standard normalises the rounding $\circ (\dots)$



Short and biased history of floating point units

Short and biased history of floating point units

- In the 70s, adders and multipliers
- In 1985, IEEE-754 standard normalises the rounding $\circ(\dots)$
- In the 90s, FMA (fused multiply add): $R = \circ(X \times Y + Z)$
 - two operations in one instruction: *faster*
 - one single rounding: *more accurate*



Short and biased history of floating point units

Short and biased history of floating point units

- In the 70s, adders and multipliers
- In 1985, IEEE-754 standard normalises the rounding $\circ(\dots)$
- In the 90s, FMA (fused multiply add): $R = \circ(X \times Y + Z)$

- These days: $R \approx Z + \sum_{i=0}^{N-1} X_i \times Y_i$



Short and biased history of floating point units

Short and biased history of floating point units

- In the 70s, adders and multipliers
- In 1985, IEEE-754 standard normalises the rounding $\circ(\dots)$
- In the 90s, FMA (fused multiply add): $R = \circ(X \times Y + Z)$
- These days: $R = \circ\left(Z + \sum_{i=0}^{N-1} X_i \times Y_i\right)$



Short and biased history of floating point units

Short and biased history of floating point units

- In the 70s, adders and multipliers
- In 1985, IEEE-754 standard normalises the rounding $\circ(\dots)$
- In the 90s, FMA (fused multiply add): $R = \circ(X \times Y + Z)$
- These days: $R = \circ\left(Z + \sum_{i=0}^{N-1} X_i \times Y_i\right)$

Multiple applications

- Matrix multiplication (neural networks, graphical applications, scientific computing,...)
- Complex arithmetic (FFT,...)



Short and biased history of floating point units

Short and biased history of floating point units

- In the 70s, adders and multipliers
- In 1985, IEEE-754 standard normalises the rounding $\circ(\dots)$
- In the 90s, FMA (fused multiply add): $R = \circ(X \times Y + Z)$
- These days: $R = \circ\left(Z + \sum_{i=0}^{N-1} X_i \times Y_i\right)$

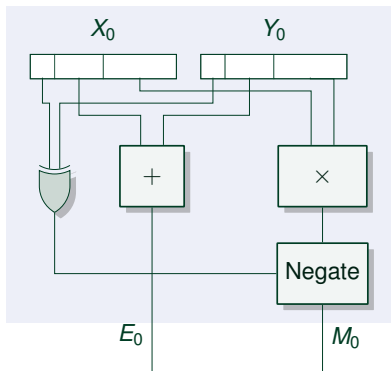
Multiple applications

- Matrix multiplication (neural networks, graphical applications, scientific computing,...)
- Complex arithmetic (FFT,...)

Objective: better than FMA chains

$$Z + \sum_{i=0}^{N-1} X_i \times Y_i \approx \circ(\dots \circ(\circ(Z + X_0 \times Y_0) + X_1 \times Y_1) \dots + X_N \times Y_N) \leftarrow$$

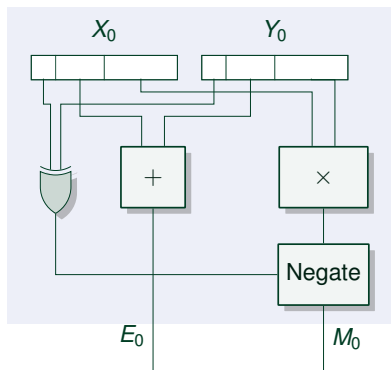
Product of floating point numbers



Easy, but the result is not a IEEE-754 floating point number:



Product of floating point numbers

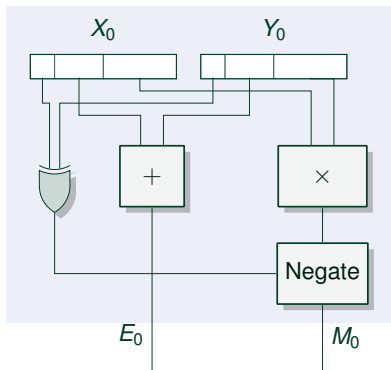


Easy, but the result is not a IEEE-754 floating point number:

- Not rounded
- Signed significands



Product of floating point numbers

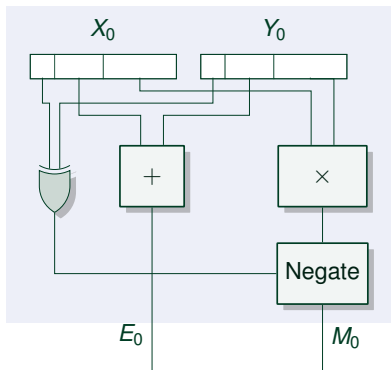


Easy, but the result is not a IEEE-754 floating point number:

- Not rounded
- Signed significands
- Significands are not normalised



Product of floating point numbers

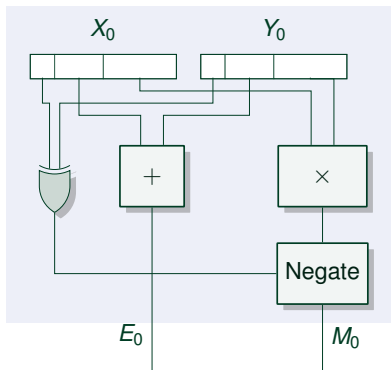


Easy, but the result is not a IEEE-754 floating point number:

- Not rounded
- Signed significands
- Significands are not normalised
 - Float is in $[1, 2[$ but product is in $[1, 4[$



Product of floating point numbers



Easy, but the result is not a IEEE-754 floating point number:

- Not rounded
- Signed significands
- Significands are not normalised
 - Float is in $[1, 2[$ but product is in $[1, 4[$
 - Product of normal and subnormal



Round a floating point number

Let's consider a decimal float with 5 digits precision.

Round by adding $\frac{1}{2}$ ulp (Unit in the Last Place) and truncating.

$$\pi \approx 3.1415927 \dots$$



Round a floating point number

Let's consider a decimal float with 5 digits precision.

Round by adding $\frac{1}{2}$ ulp (Unit in the Last Place) and truncating.

$\pi \approx 3.1415927 \dots$

$$\begin{array}{r} \mathbf{3.1415927} \\ + \quad 0.00005 \\ \hline \mathbf{3.1416427} \end{array}$$



Round a floating point number

Let's consider a decimal float with 5 digits precision.

Round by adding $\frac{1}{2}$ ulp (Unit in the Last Place) and truncating.

$$\pi \approx 3.1415927 \dots$$

3.1415927

+ 0.00005

3.1416



Round a floating point number

Let's consider a decimal float with 5 digits precision.
Round by adding $\frac{1}{2}$ ulp (Unit in the Last Place) and truncating.

$$\pi \approx 3.1415927 \dots$$

$$x = -15416500$$

$$\begin{array}{r} \mathbf{3.1415927} \\ + \quad 0.00005 \\ \hline \mathbf{3.1416} \end{array}$$



Round a floating point number

Let's consider a decimal float with 5 digits precision.
Round by adding $\frac{1}{2}$ ulp (Unit in the Last Place) and truncating.

$$\pi \approx 3.1415927 \dots$$

$$\begin{array}{r} \mathbf{3.1415927} \\ + \quad 0.00005 \\ \hline \mathbf{3.1416} \end{array}$$

$$x = -15416500$$

$$\begin{array}{r} \mathbf{1.5416500} \\ + \quad 0.00005 \\ \hline \mathbf{1.5417000} \end{array}$$



Round a floating point number

Let's consider a decimal float with 5 digits precision.
Round by adding $\frac{1}{2}$ ulp (Unit in the Last Place) and truncating.

$$\pi \approx 3.1415927 \dots$$

$$\begin{array}{r} \mathbf{3.1415927} \\ + \quad 0.00005 \\ \hline \mathbf{3.1416} \end{array}$$

$$x = -15416500$$

$$\begin{array}{r} \mathbf{1.5416500} \\ + \quad 0.00005 \\ \hline \mathbf{1.5417} \end{array}$$



Round a floating point number

Let's consider a decimal float with 5 digits precision.
Round by adding $\frac{1}{2}$ ulp (Unit in the Last Place) and truncating.

$$\pi \approx 3.1415927 \dots$$

$$\begin{array}{r} \mathbf{3.1415927} \\ + 0.00005 \\ \hline \mathbf{3.1416427} \end{array}$$

$$x = -15416500$$

$$\begin{array}{r} \mathbf{1.5416500} \\ + 0.00005 \\ \hline \mathbf{1.5417000} \end{array}$$

If exactly in the middle, round to an even float.

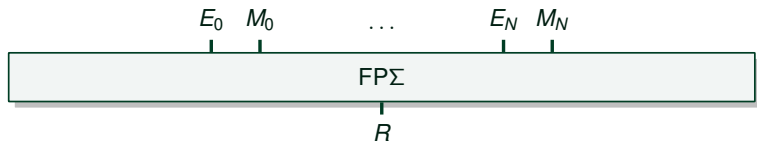
⇒ We need the following information to round:

$$\begin{array}{c} \mathbf{3.141592654} \\ \underbrace{\hspace{1.5cm}} \quad \underbrace{\hspace{1.5cm}} \\ \text{goes in } \circ(\pi) \quad \uparrow \quad \text{boolean information: is this 0? "sticky bit"} \end{array}$$

boolean information: is this digit over or under 5? "round bit"



Sum of floating point numbers



Rounding the sum of two floating point numbers

We sort $(E_0, M_0), (E_1, M_1)$ such that $E_0 \geq E_1$

+

			M_0			
			M_1			

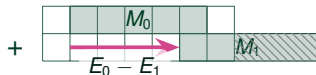
$$E_0 = E_1$$

Sum like integers



Rounding the sum of two floating point numbers

We sort $(E_0, M_0), (E_1, M_1)$ such that $E_0 \geq E_1$



$$E_0 = E_1$$

Sum like integers

$$E_0 \simeq E_1$$

Part of (shifted) M_1 is added to M_0 , and the rest is compressed in a "sticky bit"



Rounding the sum of two floating point numbers

We sort $(E_0, M_0), (E_1, M_1)$ such that $E_0 \geq E_1$



$$E_0 = E_1$$

Sum like integers



$$E_0 \simeq E_1$$

Part of (shifted) M_1 is added to M_0 , and the rest is compressed in a "sticky bit"



$$E_0 \gg E_1$$

M_1 is completely compressed in a "sticky bit"

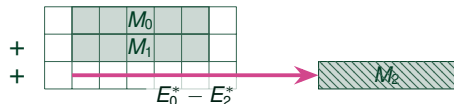
Rounding the sum of more than two floating point numbers - problem !!

Problem: cancellation

$$M_0 = -M_1 \text{ and } E_0 \gg E_2$$

$$M_0 + M_1 + M_2 = M_2$$

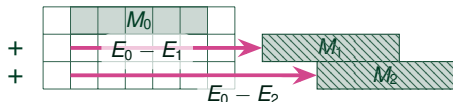
If M_2 has been totally compressed in a sticky bit, we cannot retrieve the result.



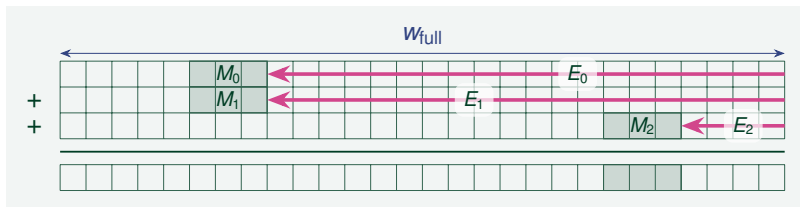
Problem: multi-sticky

$$E_0 \gg E_1 \text{ and } E_0 \gg E_2$$

If M_1 et M_2 were compressed, we cannot round the result M_0



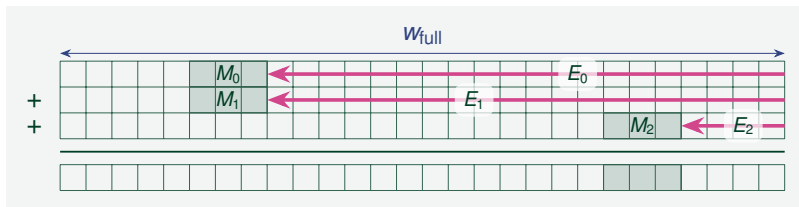
An easy but expensive solution: the Kulisch accumulator¹



- Method:
 - Convert to fixpoint
 - Sum like integer
 - Use Leading Zero Count to get final exponent and significand

¹U. W. Kulisch, "Advanced Arithmetic for the Digital Computer: Design of Arithmetic Units", 2002

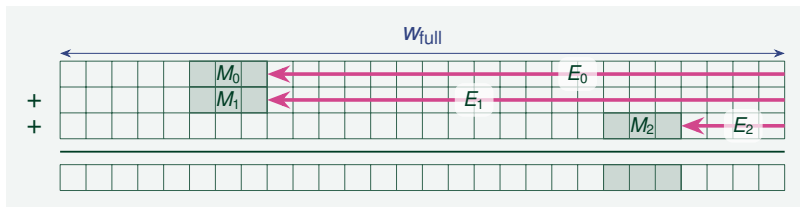
An easy but expensive solution: the Kulisch accumulator¹



- Method:
 - Convert to fixpoint
 - Sum like integer
 - Use Leading Zero Count to get final exponent and significand
- $w_{full} = 2^{w_e} + w_f - 1$

¹U. W. Kulisch, "Advanced Arithmetic for the Digital Computer: Design of Arithmetic Units", 2002

An easy but expensive solution: the Kulisch accumulator¹

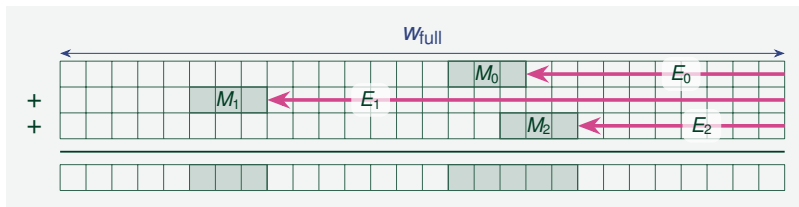


- Method:
 - Convert to fixpoint
 - Sum like integer
 - Use Leading Zero Count to get final exponent and significand
- $w_{full} = 2^{w_e} + w_f - 1$
- Good when w_e is small (or for a format with little range like Int or Posit when $es \leq 2$)

¹U. W. Kulisch, "Advanced Arithmetic for the Digital Computer: Design of Arithmetic Units", 2002



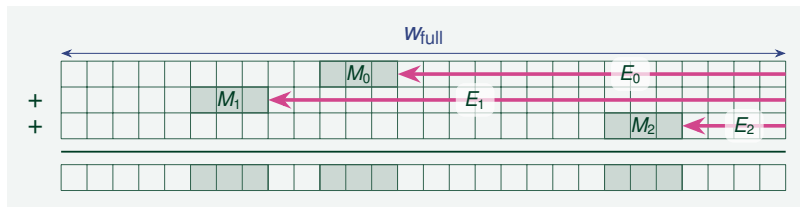
An easy but expensive solution: the Kulisch accumulator¹



- Method:
 - Convert to fixpoint
 - Sum like integer
 - Use Leading Zero Count to get final exponent and significand
- $w_{full} = 2^{w_e} + w_f - 1$
- Good when w_e is small (or for a format with little range like Int or Posit when $es \leq 2$)

¹U. W. Kulisch, "Advanced Arithmetic for the Digital Computer: Design of Arithmetic Units", 2002

An easy but expensive solution: the Kulisch accumulator¹



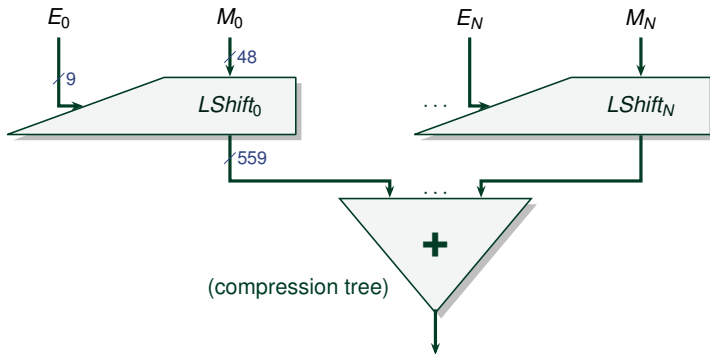
- Method:
 - Convert to fixpoint
 - Sum like integer
 - Use Leading Zero Count to get final exponent and significand
- $w_{full} = 2^{w_e} + w_f - 1$
- Good when w_e is small (or for a format with little range like Int or Posit when $es \leq 2$)

¹U. W. Kulisch, "Advanced Arithmetic for the Digital Computer: Design of Arithmetic Units", 2002

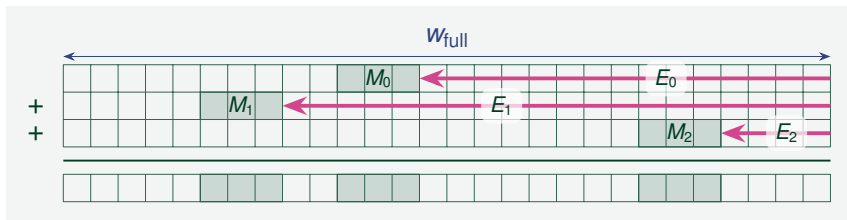


Kulisch accumulator: Architecture

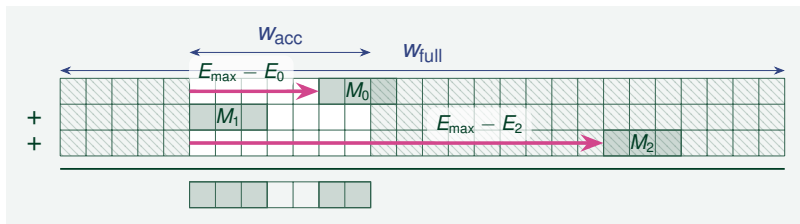
Size in bits for FP32 dot-product



Kulisch accumulator: Expensive and empty



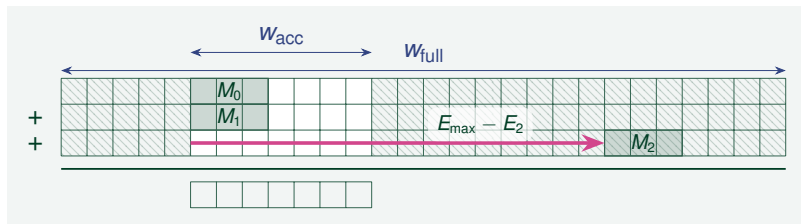
A less expensive version: truncated Kulisch



- Method similar to floating point sum:
 - Choose $w_{\text{acc}} < w_{\text{full}}$ (arbitrarily)
 - Align all numbers on the biggest one, throw away any bits that don't fill in w_{acc} (no sticky, we don't care)
 - Sum like integer, round to float



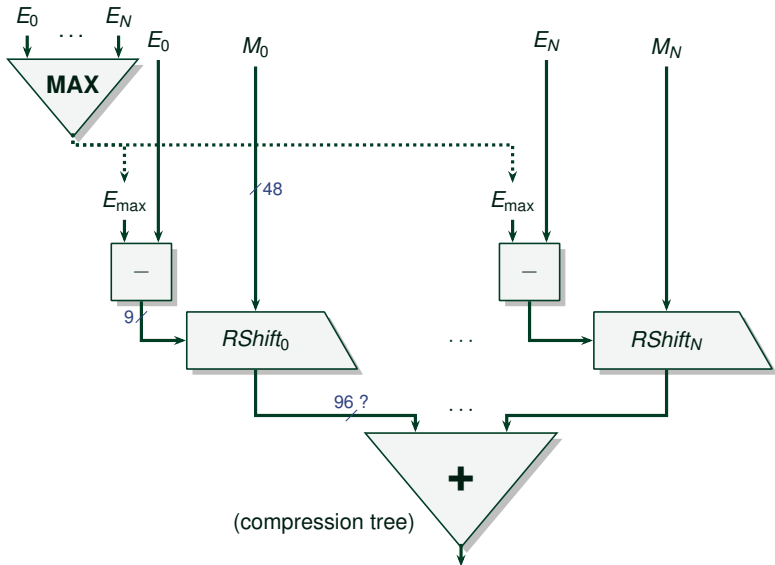
A less expensive version: truncated Kulisch



- Method similar to floating point sum:
 - Choose $w_{\text{acc}} < w_{\text{full}}$ (arbitrarily)
 - Align all numbers on the biggest one, throw away any bits that don't fill in w_{acc} (no sticky, we don't care)
 - Sum like integer, round to float
- Inexact computation

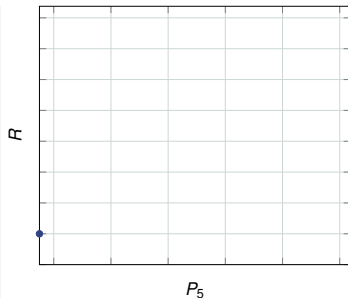
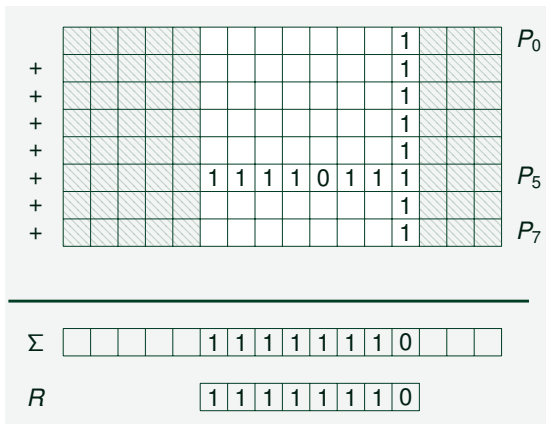


Truncated Kulisch accumulator: Architecture



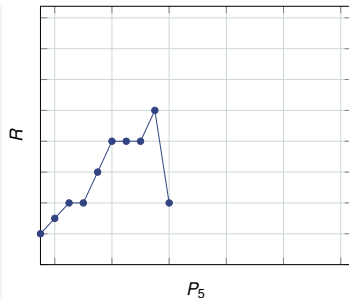
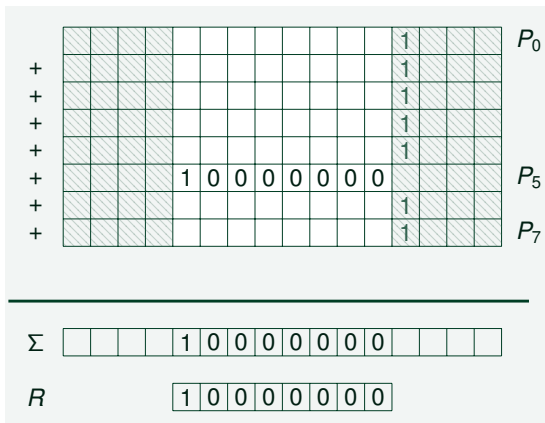
Non monotonicity of truncated Kulisch

Recently presented in detail by Mantas Mikaitis.



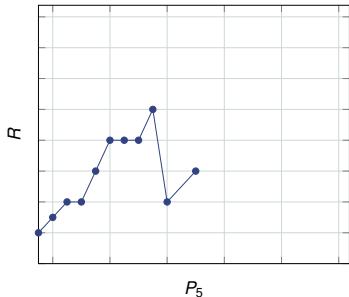
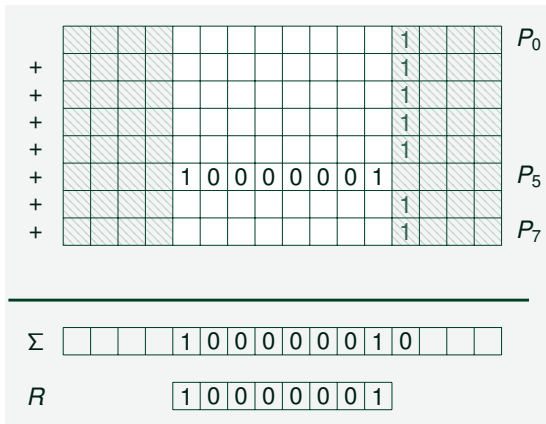
Non monotonicity of truncated Kulisch

Recently presented in detail by Mantas Mikaitis.



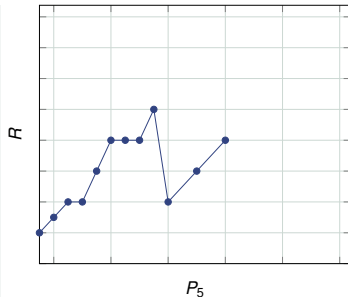
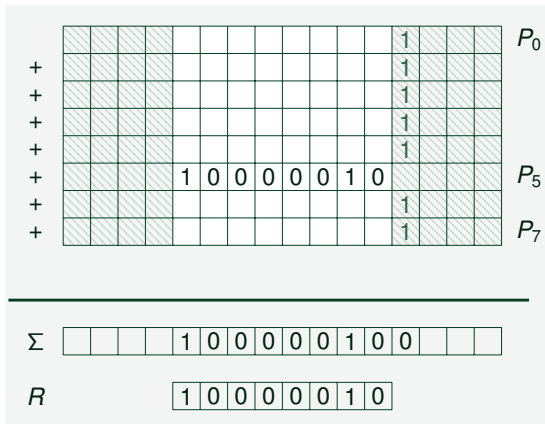
Non monotonicity of truncated Kulisch

Recently presented in detail by Mantas Mikaitis.



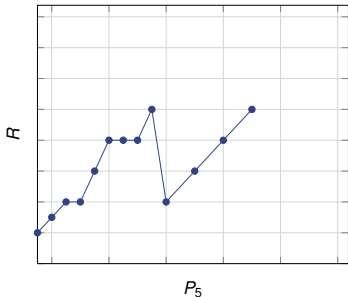
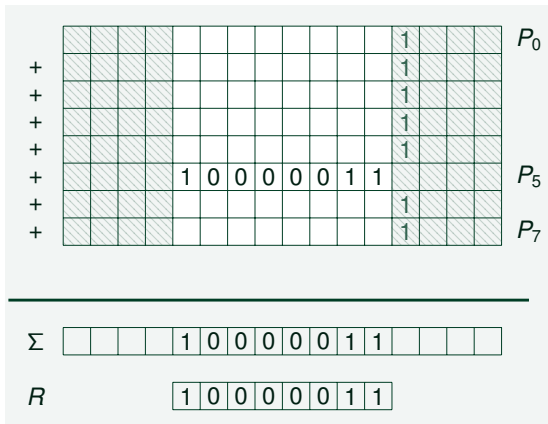
Non monotonicity of truncated Kulisch

Recently presented in detail by Mantas Mikaitis.



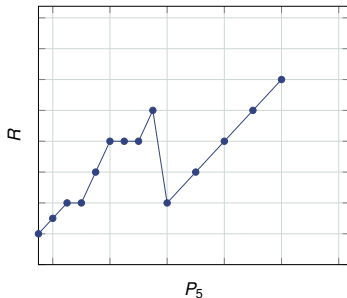
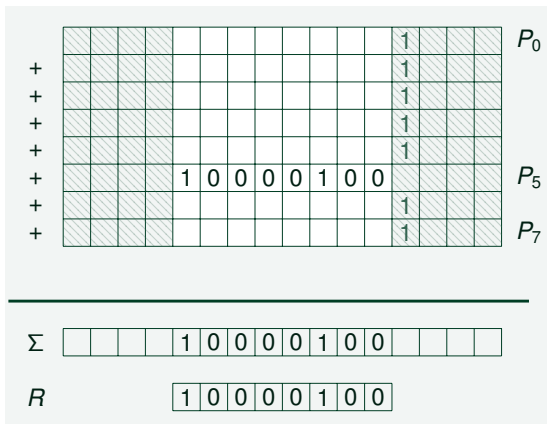
Non monotonicity of truncated Kulisch

Recently presented in detail by Mantas Mikaitis.



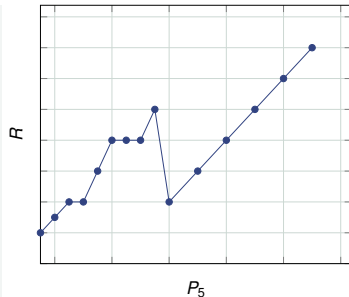
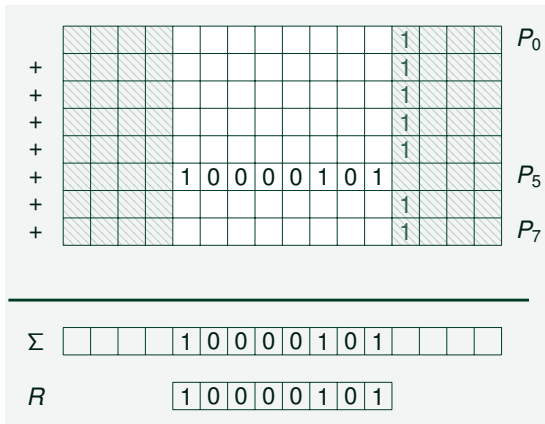
Non monotonicity of truncated Kulisch

Recently presented in detail by Mantas Mikaitis.



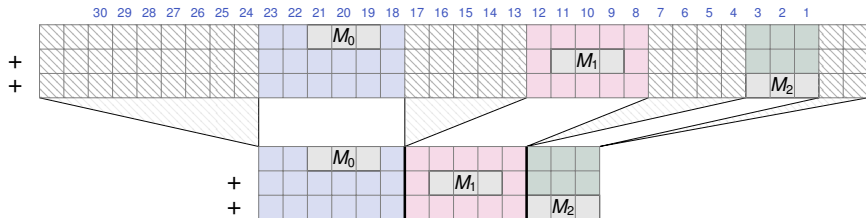
Non monotonicity of truncated Kulisch


Recently presented in detail by Mantas Mikaitis.



Compressed Kulisch (not truncated)¹

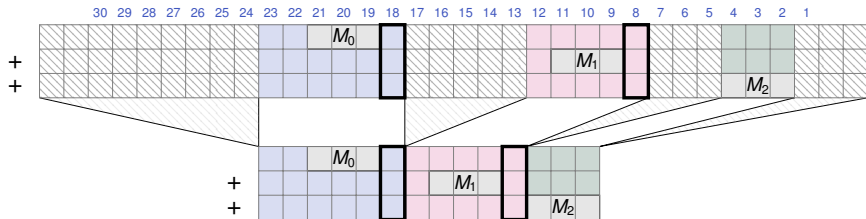
Correctly rounded sum of products



¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

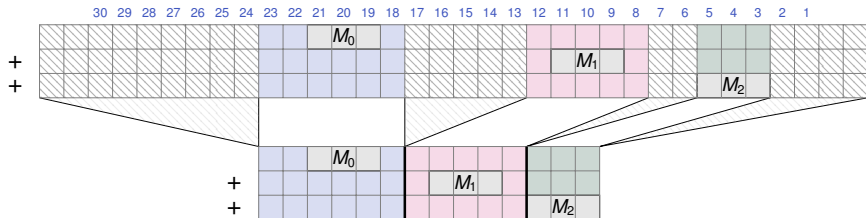
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators"

Compressed Kulisch (not truncated)¹

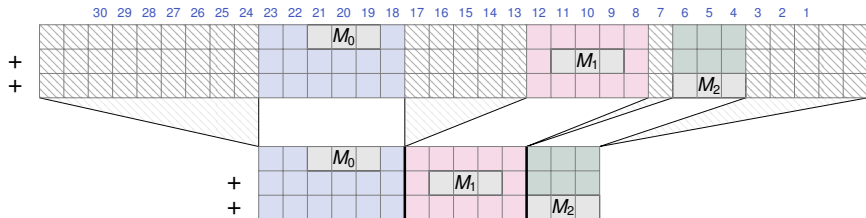
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

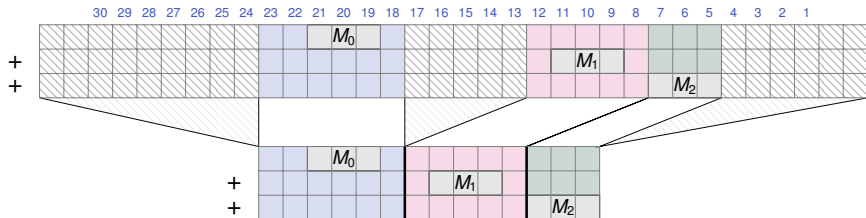
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

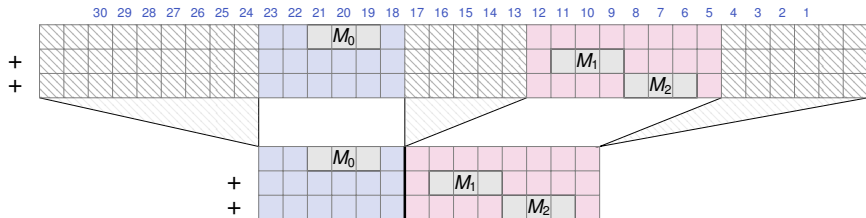
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

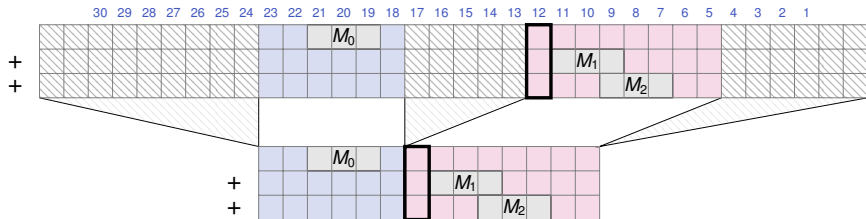
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

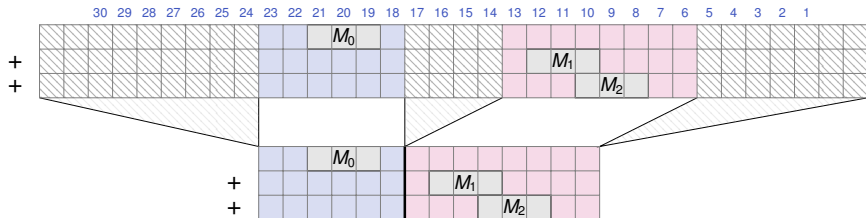
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

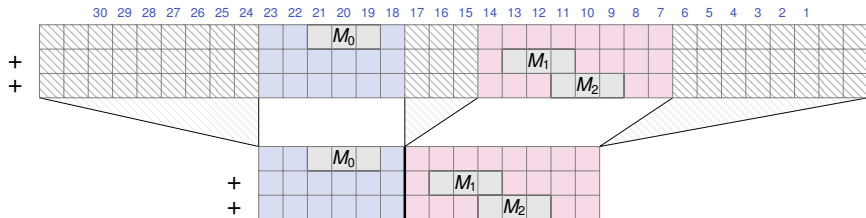
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

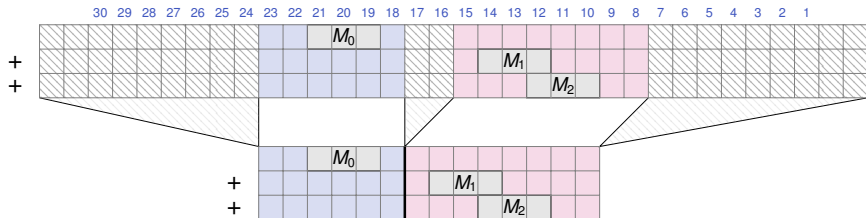
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

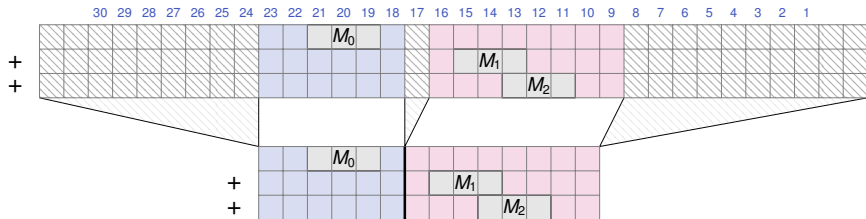
Correctly rounded sum of products



¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

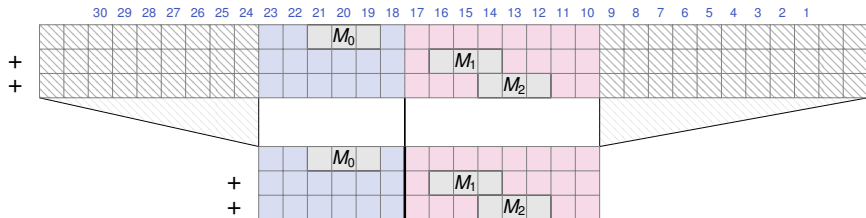
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators"

Compressed Kulisch (not truncated)¹

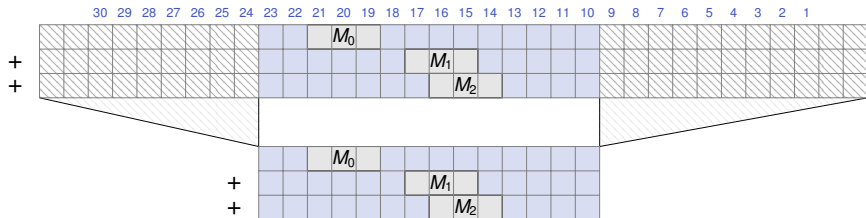
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

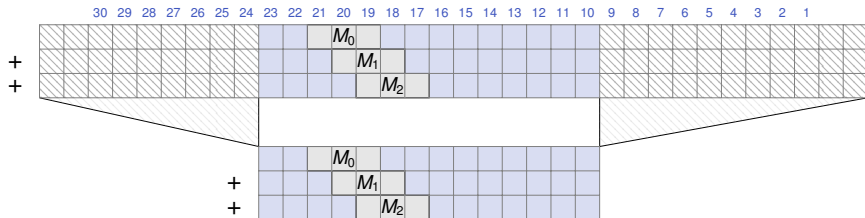
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

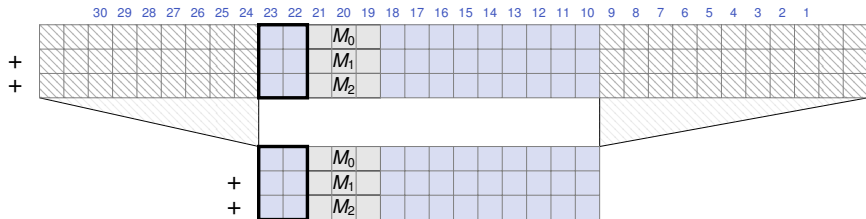
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

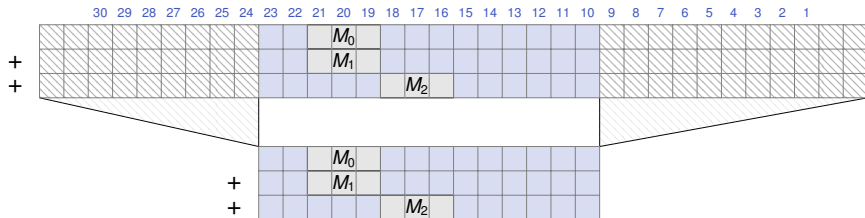
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

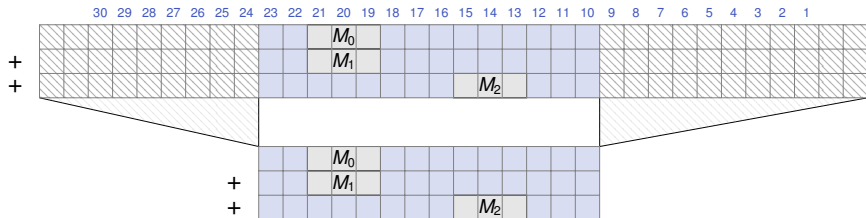
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

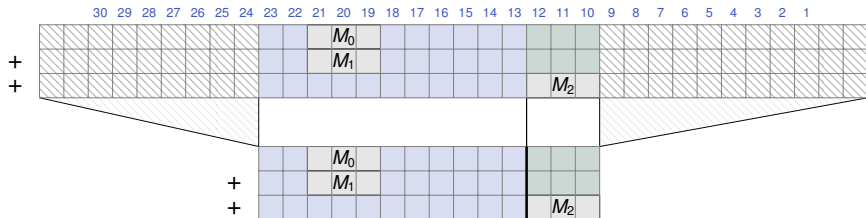
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

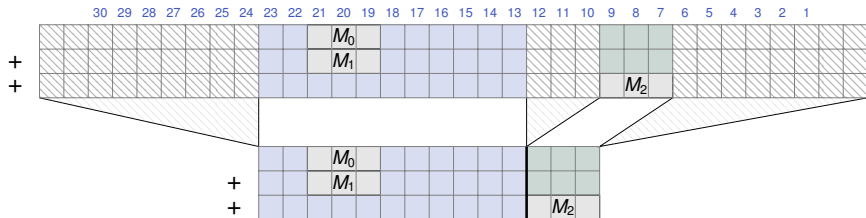
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

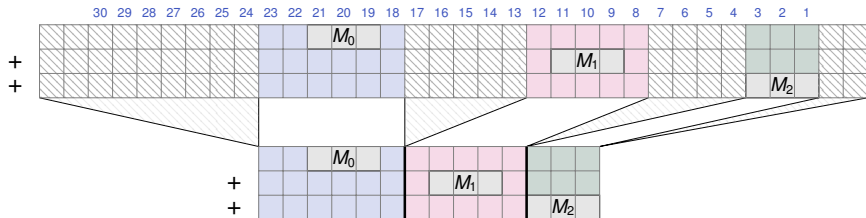
Correctly rounded sum of products




¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Compressed Kulisch (not truncated)¹

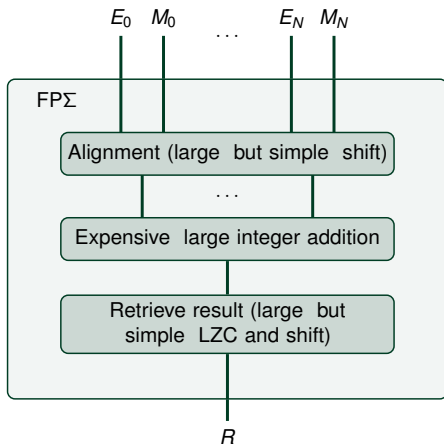
Correctly rounded sum of products



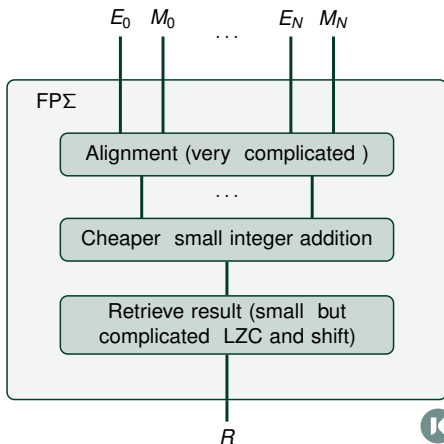
¹O. Desrentes, B. Dupont de Dinechin, F. de Dinechin, "Exact Fused Dot Product Add Operators" 

Architecture comparison

Kulisch

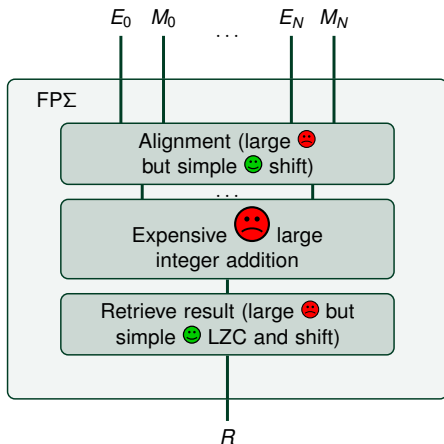


Compressed

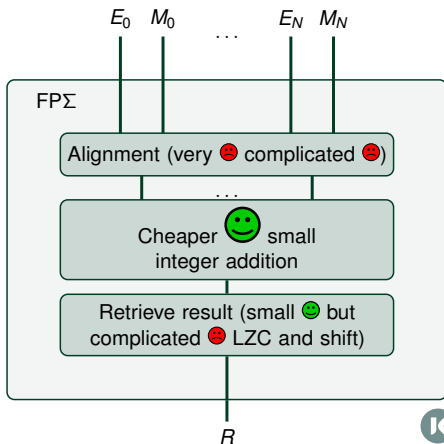


Architecture comparison

Kulisch

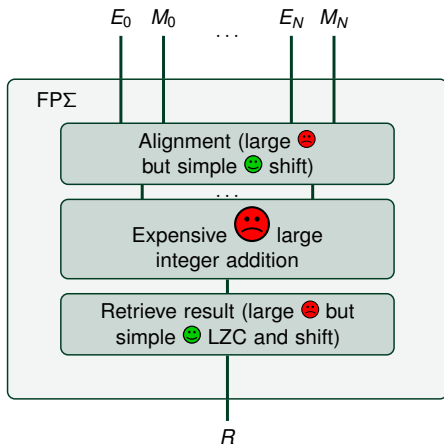


Compressed

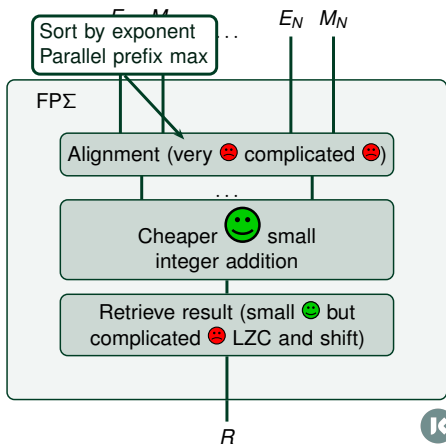


Architecture comparison

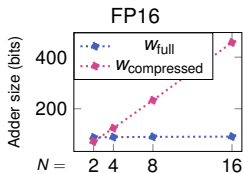
Kulisch



Compressed



Results

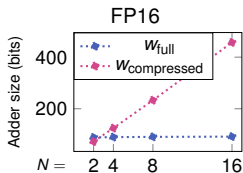


$$W_{full} \sim 2^{W_E}$$

$$W_{compressed} \sim N \times W_F$$



Results



$$W_{full} \sim 2^{w_E}$$

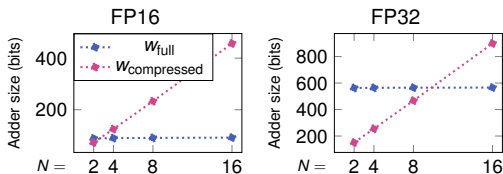
$$W_{compressed} \sim N \times w_F$$

3 situations:

- The format has a lot of precision compared to the range, $w_{compressed} > w_{full}$ even for small N



Results



$$w_{full} \sim 2^{w_E}$$

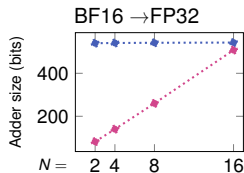
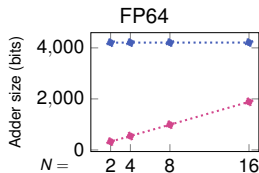
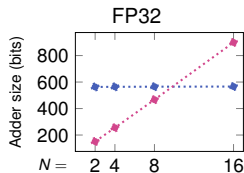
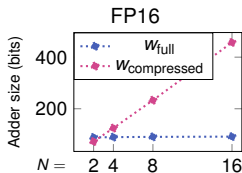
$$w_{compressed} \sim N \times w_F$$

3 situations:

- The format has a lot of precision compared to the range, $w_{compressed} > w_{full}$ even for small N
- The format is more balanced, $w_{compressed} < w_{full}$ for small N



Results



$$w_{full} \sim 2^{w_E}$$

$$w_{compressed} \sim N \times w_F$$

3 situations:

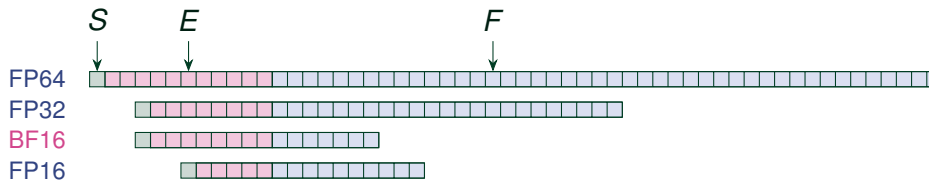
- The format has a lot of precision compared to the range, $w_{compressed} > w_{full}$ even for small N
- The format is more balanced, $w_{compressed} < w_{full}$ for small N
- The format has a lot of range compared to the precision $w_{compressed} < w_{full}$ until larger N



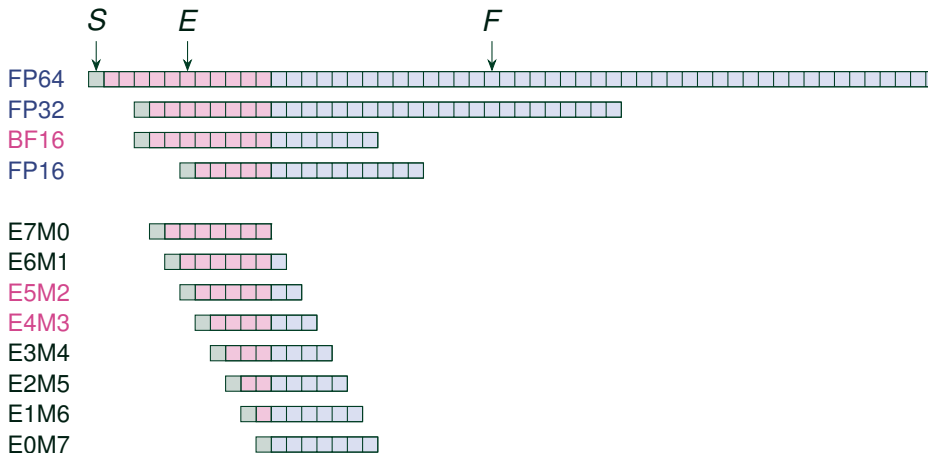
8 bits formats



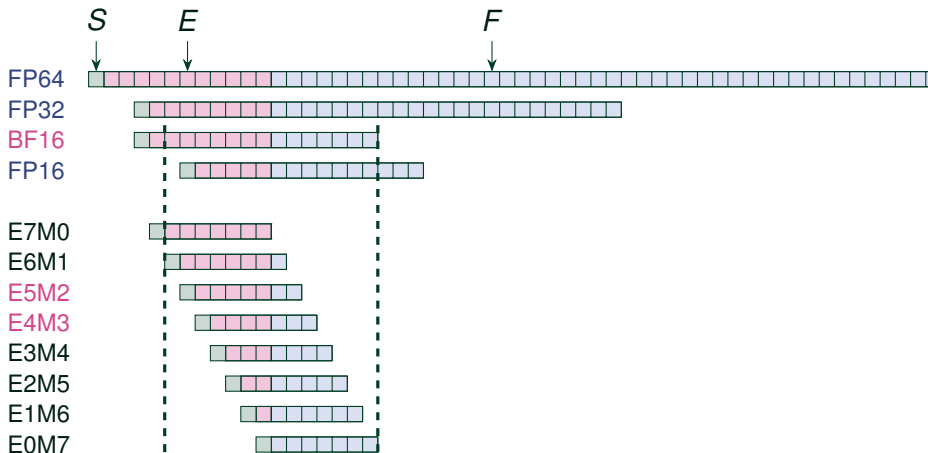
Formats



Formats



Formats



Special values?

8 bits is not a lot of bits (256 different values)

C.4 Value Table: P4, P = 4, emax = 7

0d00 = 0.0000000 + 0.0	0d40 = 0.1000000 + 0b1000020 = 1.0	0d80 = 1.0000000 + Inf
0d01 = 0.0000001 + 0b0100020 = 1.125	0d41 = 0.1000001 + 0b1000021 = 1.125	0d81 = 1.0000001 + 0b1000020 = -0.009765625
0d02 = 0.0000010 + 0b0010020 = 1.25	0d42 = 0.1000010 + 0b1000022 = 1.25	0d82 = 1.0000010 + 0b1000021 = -0.001953125
0d03 = 0.0000011 + 0b0110020 = 1.375	0d43 = 0.1000011 + 0b1000023 = 1.375	0d83 = 1.0000011 + 0b1000022 = -0.003928125
0d04 = 0.0001000 + 0b0000020 = 1.5	0d44 = 0.1000100 + 0b1000024 = 1.5	0d84 = 1.0000100 + 0b1000023 = -0.00785625
0d05 = 0.0001001 + 0b0100020 = 1.625	0d45 = 0.1000101 + 0b1000025 = 1.625	0d85 = 1.0000101 + 0b1000024 = -0.011784375
0d06 = 0.0001010 + 0b0000020 = 1.75	0d46 = 0.1000110 + 0b1000026 = 1.75	0d86 = 1.0000110 + 0b1000025 = -0.0157125
0d07 = 0.0001011 + 0b0100020 = 1.875	0d47 = 0.1000111 + 0b1000027 = 1.875	0d87 = 1.0000111 + 0b1000026 = -0.019640625
0d08 = 0.0002000 + 0b0000020 = 2.0	0d48 = 0.1001000 + 0b1000028 = 2.0	0d88 = 1.0001000 + 0b1000027 = -0.02356875
0d09 = 0.0002001 + 0b0100020 = 2.125	0d49 = 0.1001001 + 0b1000029 = 2.125	0d89 = 1.0001001 + 0b1000028 = -0.027496875
0d0a = 0.0002010 + 0b0000020 = 2.25	0d4a = 0.1001010 + 0b1000030 = 2.25	0d8a = 1.0001010 + 0b1000029 = -0.031425
0d0b = 0.0002011 + 0b0100020 = 2.375	0d4b = 0.1001011 + 0b1000031 = 2.375	0d8b = 1.0001011 + 0b1000030 = -0.035353125
0d0c = 0.0003000 + 0b0000020 = 2.5	0d4c = 0.1001100 + 0b1000032 = 2.5	0d8c = 1.0001100 + 0b1000031 = -0.03928125
0d0d = 0.0003001 + 0b0100020 = 2.625	0d4d = 0.1001101 + 0b1000033 = 2.625	0d8d = 1.0001101 + 0b1000032 = -0.043209375
0d0e = 0.0003010 + 0b0000020 = 2.75	0d4e = 0.1001110 + 0b1000034 = 2.75	0d8e = 1.0001110 + 0b1000033 = -0.0471375
0d0f = 0.0003011 + 0b0100020 = 2.875	0d4f = 0.1001111 + 0b1000035 = 2.875	0d8f = 1.0001111 + 0b1000034 = -0.051065625
0d10 = 0.0004000 + 0b0000020 = 3.0	0d50 = 0.1010000 + 0b1000036 = 3.0	0d90 = 1.0010000 + 0b1000035 = -0.05499375
0d11 = 0.0004001 + 0b0100020 = 3.125	0d51 = 0.1010001 + 0b1000037 = 3.125	0d91 = 1.0010001 + 0b1000036 = -0.058921875
0d12 = 0.0004010 + 0b0000020 = 3.25	0d52 = 0.1010010 + 0b1000038 = 3.25	0d92 = 1.0010010 + 0b1000037 = -0.06285
0d13 = 0.0004011 + 0b0100020 = 3.375	0d53 = 0.1010011 + 0b1000039 = 3.375	0d93 = 1.0010011 + 0b1000038 = -0.066778125
0d14 = 0.0005000 + 0b0000020 = 3.5	0d54 = 0.1010100 + 0b1000040 = 3.5	0d94 = 1.0010100 + 0b1000039 = -0.07070625
0d15 = 0.0005001 + 0b0100020 = 3.625	0d55 = 0.1010101 + 0b1000041 = 3.625	0d95 = 1.0010101 + 0b1000040 = -0.074634375
0d16 = 0.0005010 + 0b0000020 = 3.75	0d56 = 0.1010110 + 0b1000042 = 3.75	0d96 = 1.0010110 + 0b1000041 = -0.0785625
0d17 = 0.0005011 + 0b0100020 = 3.875	0d57 = 0.1010111 + 0b1000043 = 3.875	0d97 = 1.0010111 + 0b1000042 = -0.082490625
0d18 = 0.0006000 + 0b0000020 = 4.0	0d58 = 0.1011000 + 0b1000044 = 4.0	0d98 = 1.0011000 + 0b1000043 = -0.08641875
0d19 = 0.0006001 + 0b0100020 = 4.125	0d59 = 0.1011001 + 0b1000045 = 4.125	0d99 = 1.0011001 + 0b1000044 = -0.090346875
0d1a = 0.0006010 + 0b0000020 = 4.25	0d5a = 0.1011010 + 0b1000046 = 4.25	0d9a = 1.0011010 + 0b1000045 = -0.094275
0d1b = 0.0006011 + 0b0100020 = 4.375	0d5b = 0.1011011 + 0b1000047 = 4.375	0d9b = 1.0011011 + 0b1000046 = -0.098203125
0d1c = 0.0007000 + 0b0000020 = 4.5	0d5c = 0.1011100 + 0b1000048 = 4.5	0d9c = 1.0011100 + 0b1000047 = -0.10213125
0d1d = 0.0007001 + 0b0100020 = 4.625	0d5d = 0.1011101 + 0b1000049 = 4.625	0d9d = 1.0011101 + 0b1000048 = -0.106059375
0d1e = 0.0007010 + 0b0000020 = 4.75	0d5e = 0.1011110 + 0b1000050 = 4.75	0d9e = 1.0011110 + 0b1000049 = -0.110000000
0d1f = 0.0007011 + 0b0100020 = 4.875	0d5f = 0.1011111 + 0b1000051 = 4.875	0d9f = 1.0011111 + 0b1000050 = -0.113930625
0d20 = 0.0008000 + 0b0000020 = 5.0	0d60 = 0.1012000 + 0b1000052 = 5.0	0da0 = 1.0012000 + 0b1000051 = -0.11786125
0d21 = 0.0008001 + 0b0100020 = 5.125	0d61 = 0.1012001 + 0b1000053 = 5.125	0da1 = 1.0012001 + 0b1000052 = -0.121791875
0d22 = 0.0008010 + 0b0000020 = 5.25	0d62 = 0.1012010 + 0b1000054 = 5.25	0da2 = 1.0012010 + 0b1000053 = -0.1257225
0d23 = 0.0008011 + 0b0100020 = 5.375	0d63 = 0.1012011 + 0b1000055 = 5.375	0da3 = 1.0012011 + 0b1000054 = -0.129653125
0d24 = 0.0009000 + 0b0000020 = 5.5	0d64 = 0.1012100 + 0b1000056 = 5.5	0da4 = 1.0012100 + 0b1000055 = -0.13358375
0d25 = 0.0009001 + 0b0100020 = 5.625	0d65 = 0.1012101 + 0b1000057 = 5.625	0da5 = 1.0012101 + 0b1000056 = -0.137514375
0d26 = 0.0009010 + 0b0000020 = 5.75	0d66 = 0.1012110 + 0b1000058 = 5.75	0da6 = 1.0012110 + 0b1000057 = -0.141445000
0d27 = 0.0009011 + 0b0100020 = 5.875	0d67 = 0.1012111 + 0b1000059 = 5.875	0da7 = 1.0012111 + 0b1000058 = -0.145375625
0d28 = 0.0010000 + 0b0000020 = 6.0	0d68 = 0.1013000 + 0b1000060 = 6.0	0da8 = 1.0013000 + 0b1000059 = -0.14930625
0d29 = 0.0010001 + 0b0100020 = 6.125	0d69 = 0.1013001 + 0b1000061 = 6.125	0da9 = 1.0013001 + 0b1000060 = -0.153236875
0d2a = 0.0010010 + 0b0000020 = 6.25	0d6a = 0.1013010 + 0b1000062 = 6.25	0daa = 1.0013010 + 0b1000061 = -0.1571675
0d2b = 0.0010011 + 0b0100020 = 6.375	0d6b = 0.1013011 + 0b1000063 = 6.375	0dab = 1.0013011 + 0b1000062 = -0.161098125
0d2c = 0.0011000 + 0b0000020 = 6.5	0d6c = 0.1013100 + 0b1000064 = 6.5	0dac = 1.0013100 + 0b1000063 = -0.16502875
0d2d = 0.0011001 + 0b0100020 = 6.625	0d6d = 0.1013101 + 0b1000065 = 6.625	0dad = 1.0013101 + 0b1000064 = -0.168959375
0d2e = 0.0011010 + 0b0000020 = 6.75	0d6e = 0.1013110 + 0b1000066 = 6.75	0dae = 1.0013110 + 0b1000065 = -0.172890000
0d2f = 0.0011011 + 0b0100020 = 6.875	0d6f = 0.1013111 + 0b1000067 = 6.875	0daf = 1.0013111 + 0b1000066 = -0.176820625
0d30 = 0.0012000 + 0b0000020 = 7.0	0d70 = 0.1100000 + 0b1000068 = 7.0	0db0 = 1.0014000 + 0b1000067 = -0.18075125
0d31 = 0.0012001 + 0b0100020 = 7.125	0d71 = 0.1100001 + 0b1000069 = 7.125	0db1 = 1.0014001 + 0b1000068 = -0.184681875
0d32 = 0.0012010 + 0b0000020 = 7.25	0d72 = 0.1100010 + 0b1000070 = 7.25	0db2 = 1.0014010 + 0b1000069 = -0.1886125
0d33 = 0.0012011 + 0b0100020 = 7.375	0d73 = 0.1100011 + 0b1000071 = 7.375	0db3 = 1.0014011 + 0b1000070 = -0.192543125
0d34 = 0.0013000 + 0b0000020 = 7.5	0d74 = 0.1100100 + 0b1000072 = 7.5	0db4 = 1.0014100 + 0b1000071 = -0.19647375
0d35 = 0.0013001 + 0b0100020 = 7.625	0d75 = 0.1100101 + 0b1000073 = 7.625	0db5 = 1.0014101 + 0b1000072 = -0.200404375
0d36 = 0.0013010 + 0b0000020 = 7.75	0d76 = 0.1100110 + 0b1000074 = 7.75	0db6 = 1.0014110 + 0b1000073 = -0.204335000
0d37 = 0.0013011 + 0b0100020 = 7.875	0d77 = 0.1100111 + 0b1000075 = 7.875	0db7 = 1.0014111 + 0b1000074 = -0.208265625
0d38 = 0.0014000 + 0b0000020 = 8.0	0d78 = 0.1101000 + 0b1000076 = 8.0	0db8 = 1.0015000 + 0b1000075 = -0.21219625
0d39 = 0.0014001 + 0b0100020 = 8.125	0d79 = 0.1101001 + 0b1000077 = 8.125	0db9 = 1.0015001 + 0b1000076 = -0.216126875
0d3a = 0.0014010 + 0b0000020 = 8.25	0d7a = 0.1101010 + 0b1000078 = 8.25	0dba = 1.0015010 + 0b1000077 = -0.2200575
0d3b = 0.0014011 + 0b0100020 = 8.375	0d7b = 0.1101011 + 0b1000079 = 8.375	0dbb = 1.0015011 + 0b1000078 = -0.223988125
0d3c = 0.0015000 + 0b0000020 = 8.5	0d7c = 0.1101100 + 0b1000080 = 8.5	0dbc = 1.0015100 + 0b1000079 = -0.22791875
0d3d = 0.0015001 + 0b0100020 = 8.625	0d7d = 0.1101101 + 0b1000081 = 8.625	0dbd = 1.0015101 + 0b1000080 = -0.231849375
0d3e = 0.0015010 + 0b0000020 = 8.75	0d7e = 0.1101110 + 0b1000082 = 8.75	0dbe = 1.0015110 + 0b1000081 = -0.235780000
0d3f = 0.0015011 + 0b0100020 = 8.875	0d7f = 0.1101111 + 0b1000083 = 8.875	0dbf = 1.0015111 + 0b1000082 = -0.239710625



Special values ?

8 bits is not a lot of bits (256 different values)

NaN

Having $2^{w_F+1} - 2$ NaN values is a waste of encoding space.
Should we keep at least one ? **Yes**



Special values ?

8 bits is not a lot of bits (256 different values)

NaN

Having $2^{w_F+1} - 2$ NaN values is a waste of encoding space.
Should we keep at least one ? **Yes**

∞

Having 2 infinities is a waste of encoding space.
Should we keep them anyway ? **Maybe**



Special values ?

8 bits is not a lot of bits (256 different values)

NaN

Having $2^{w_F+1} - 2$ NaN values is a waste of encoding space.
Should we keep at least one ? **Yes**

∞

Having 2 infinities is a waste of encoding space.
Should we keep them anyway ? **Maybe**

-0

Having 2 zeros is a waste of encoding space.
Should we keep -0 ? **Maybe**



Use of encoding space (E4M3 example)

Graphcore¹:

+0				NaN			
sn				-sn			
sn				-sn			
sn				-sn			
sn				-sn			
sn				-sn			
sn				-sn			
sn				-sn			
sn				-sn			

¹B. Noune, P. Jones, D. Justus, D. Masters, and C. Luschi, "8-bit numerical formats for deep neural networks," 2022

Use of encoding space (E4M3 example)

IEEE WG P3109 standard (with saturation)³:

+0				NaN			
sn				-sn			
sn				-sn			
sn				-sn			
sn				-sn			
sn				-sn			
sn				-sn			
sn				-sn			

³"IEEE Working Group P3109 Interim Report on 8-bit Binary Floating-point Formats"



Bias ?

For a real exponent E , it is encoded in the format like a positive number, by adding a bias.

	FP16 (IEEE-754)
w_E	5
Bias	$2^{w_E-1} - 1 = 15$
Max exp	$2^{w_E-1} - 1 = 15$
Min exp	$-2^{w_E-1} + 2 = -14$



Bias ?

For a real exponent E , it is encoded in the format like a positive number, by adding a bias.

	FP16 (IEEE-754)	E5M2 (until recently)
w_E	5	5
Bias	$2^{w_E-1} - 1 = 15$	$2^{w_E-1} - 1 = 15$
Max exp	$2^{w_E-1} - 1 = 15$	$2^{w_E-1} = 16$
Min exp	$-2^{w_E-1} + 2 = -14$	$-2^{w_E-1} + 2 = -14$



Bias ?

For a real exponent E , it is encoded in the format like a positive number, by adding a bias.

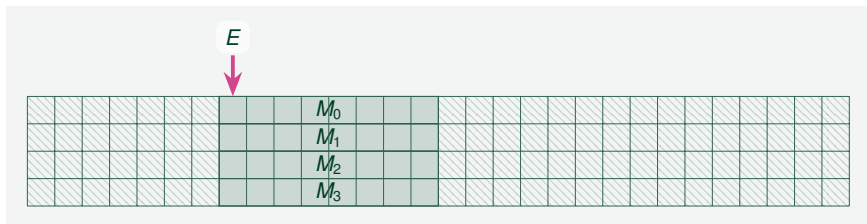
	FP16 (IEEE-754)	E5M2 (until recently)	E5M2 (IEEE WG P3109)
w_E	5	5	5
Bias	$2^{w_E-1} - 1 = 15$	$2^{w_E-1} - 1 = 15$	$2^{w_E-1} = 16$
Max exp	$2^{w_E-1} - 1 = 15$	$2^{w_E-1} = 16$	$2^{w_E} - 1 = 15$
Min exp	$-2^{w_E-1} + 2 = -14$	$-2^{w_E-1} + 2 = -14$	$-2^{w_E-1} + 1 = -15$



Block floating point

History

A vector of fixpoint numbers that share an exponent.



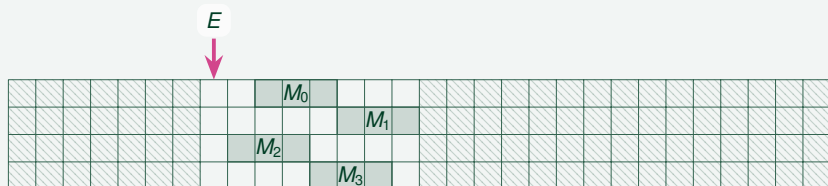
Block floating point

History

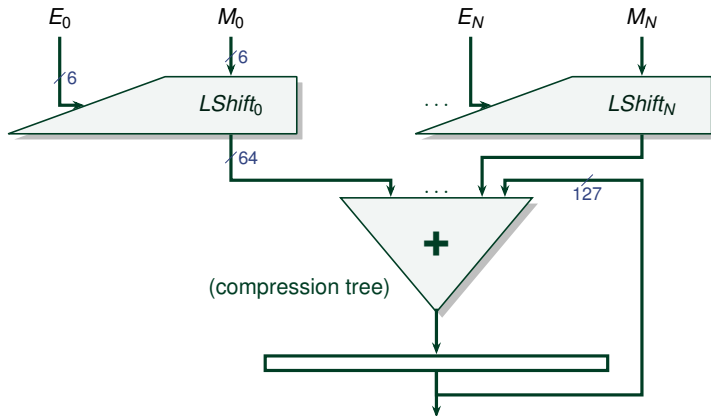
A vector of fixpoint numbers that share an exponent.

Recent variant

A vector of small floating point numbers that share a bias modifier/scaling factor



Kulisch architecture for exact 8 bits dot product

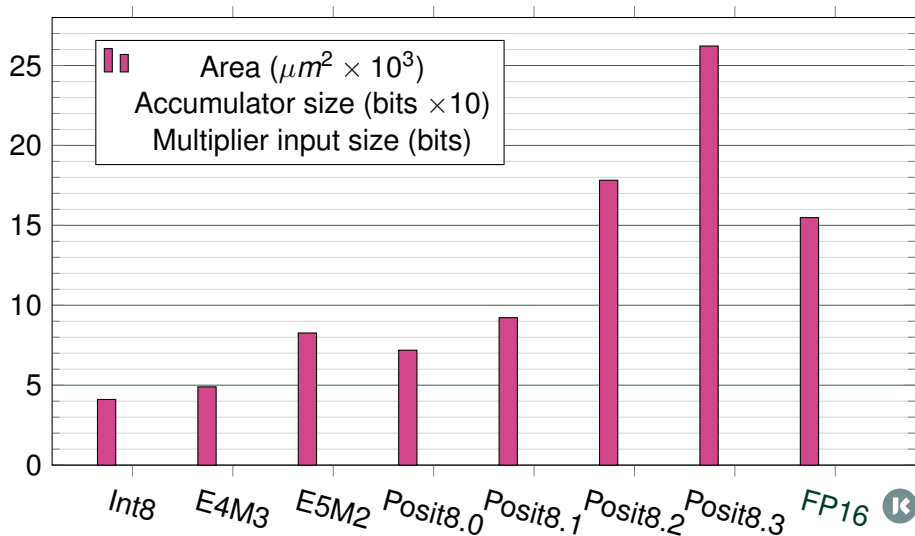


Kulisch architecture for exact 8 bits dot product

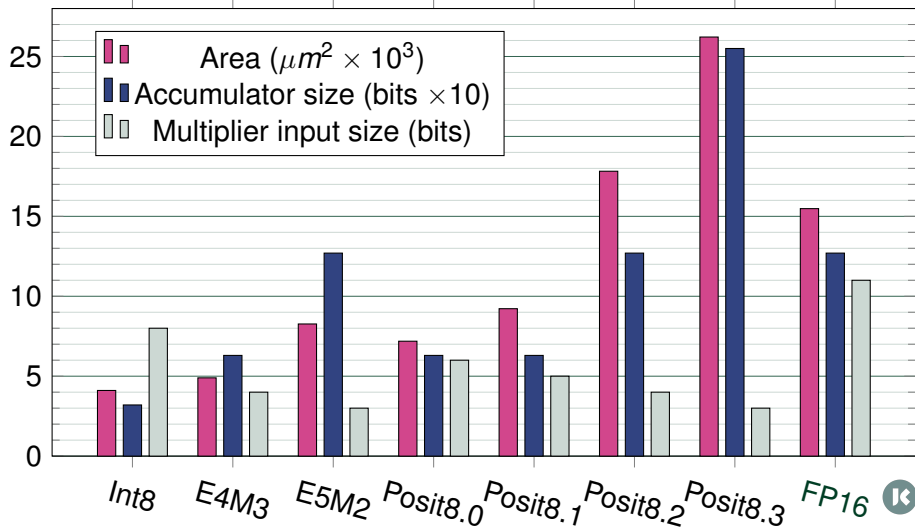
Format	Product				Accumulator	
	size	LSB	MSB	w (in bits)	MSB	w_{acc} (in bits)
INT8	8×8	0	15	16	31	32
E4M3	4×4	-18	16	36	44	63+1
E5M2	3×3	-32	30	64	94	127+1
Posit8.0	6×6	-6	6	26	50	63+1
Posit8.1	5×5	-24	24	50	38	63+1
Posit8.2	4×4	-48	48	98	78	127+1
Posit8.3	3×3	-96	96	194	158	255+1
FP16	11×11	-48	30	80	78	127+1



Comparison with other 8 bits formats: posits and integers

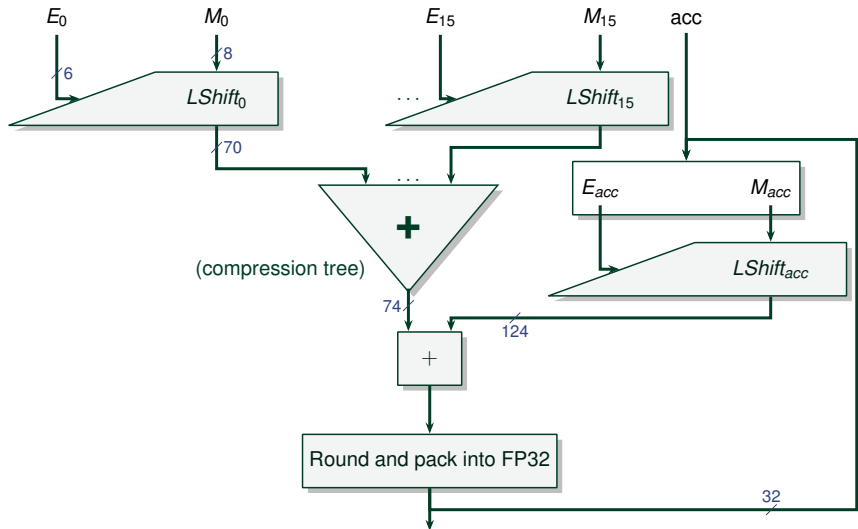


Comparison with other 8 bits formats: posits and integers



How it is implemented in the MPPA

Instead of accumulating in large fixpoint, use a FP32



Conclusion

Back to the 70s



Conclusion

Back to the 70s

- NVIDIA¹ does some sort of truncated Kulisch with their E5M2 and E4M3



^aB. Hickmann and D. Bradford, "Experimental analysis of matrix multiplication functional units", 2019



Conclusion

Back to the 70s

- NVIDIA¹ does some sort of truncated Kulisch with their E5M2 and E4M3
- Intel² is similar



^bB. Hickmann, J. Chen, M. Rotzin, A. Yang, M. Urbanski, S. Avancha, "Intel Nervana Neural Network Processor-T (NNP-T) Fused Floating Point Many-Term Dot Product", 2020



Conclusion

Back to the 70s

- NVIDIA¹ does some sort of truncated Kulisch with their E5M2 and E4M3
- Intel² is similar
- Arm³ is now doing some exact Kulisch with rounding by block of 4 products, on E5M2 and E4M3 with scaling factor



¹D. Lutz, A. Saini, M. Kroes, T. Elmer, H. Valsaraju, "Fused FP8 4-Way Dot Product with Scaling and FP32 Accumulation", 2024



Conclusion

Back to the 70s

- NVIDIA¹ does some sort of truncated Kulisch with their E5M2 and E4M3
- Intel² is similar
- Arm³ is now doing some exact Kulisch with rounding by block of 4 products, on E5M2 and E4M3 with scaling factor
- Kalray is continuing the exact Kulisch with rounding by block of 8? products, with E5M2 and E4M3 (but which ones ?)



Conclusion

Back to the 70s

- NVIDIA¹ does some sort of truncated Kulisch with their E5M2 and E4M3
- Intel² is similar
- Arm³ is now doing some exact Kulisch with rounding by block of 4 products, on E5M2 and E4M3 with scaling factor
- Kalray is continuing the exact Kulisch with rounding by block of 8? products, with E5M2 and E4M3 (but which ones ?)
- IEEE WG P3109 formats are a description of everything that can be done

