Error Analysis of Matrix Multiplication with Narrow Range Floating-Point Arithmetic

Mantas Mikaitis

School of Computer Science, University of Leeds, Leeds, UK Joint work with Theo Mary, CNRS, Paris, France

The 30th Biennial Numerical Analysis Conference Glasgow, UK 26 June, 2025



8-bit floating point on the TOP500 (June 2025)



 \rightarrow 10-bit FP \rightarrow matrix arith. \rightarrow 19-bi \rightarrow 8-bit FP \rightarrow fast integer

Devices counted: P100, V100, A100, H100, MI210, MI250X, MI300X, Intel Data Center GPU, from https://www.top500.org.

 $\label{eq:stars} \begin{array}{l} \mbox{NVIDIA Blackwell throughputs (FLOPS)} \\ \mbox{fp8 } (9\times 10^{15}) \ \ \mbox{fp16 } (4.5\times 10^{15}) \ \ \mbox{fp64 } (0.04\times 10^{15}). \end{array}$

C.4 Value Table: P4, emin = -7, emax = 7

| 0x00 = 0.0000.000 = 0.0 | $0x40 = 0.1000.000 = +0b1.000 \times 2^{\circ}0 = 1.0$ | 0x80 = 1.0000.000 = NaN | $0xc0 = 1.1000.000 = -0b1.000 \times 2^{\circ}0 = -1.0$ |
|---|--|---|---|
| $0x01 = 0.0000.001 = +0b0.001 \times 2^{-7} = 0.0009765625$ | $0x41 = 0.1000.001 = +0b1.001 \times 2^{\circ}0 = 1.125$ | $0x81 = 1.0000.001 = -0b0.001 \times 2^{-7} = -0.0009765625$ | $0xc1 = 1.1000.001 = -0b1.001 \times 2^{-0} = -1.125$ |
| $0x02 = 0.0000.010 = +0b0.010 \times 2^{-7} = 0.001953125$ | $0x42 = 0.1000.010 = +0b1.010 \times 2^{\circ}0 = 1.25$ | $0x82 = 1.0000.010 = -0b0.010 \times 2^{-7} = -0.001953125$ | $0xc2 = 1.1000.010 = -0b1.010 \times 2^{-0} = -1.25$ |
| $0x03 = 0.0000.011 = +060.011 \times 2^{-7} = 0.0029296975$ | $0x43 = 0.1000.011 = +0b1.011 \times 270 = 1.375$ | $0x83 = 1.0000.011 = -0b0.011 \times 2^{-7} = -0.0029296875$ | $0xc3 = 1.1000.011 = -0b1.011 \times 270 = -1.375$ |
| $0 \times 04 = 0.0000 \ 100 = \pm 0 \pm 0.100 \times 2^{-7} = 0.00390625$ | $0x44 = 0.1000 100 = \pm 0b1 100 \times 200 = 1.5$ | $0 \times 84 = 1.0000, 100 = -0 \times 0.100 \times 2^{-7} = -0.00390625$ | $0xc4 = 1,1000,100 = -0b1,100 \times 270 = -1.5$ |
| 0+05 - 0.0000 101 - 1.050 101 x 27 7 - 0.0019929125 | 0=4E = 0.1000.101 = 1.0b1.101 × 070 = 1.60E | 0+95 - 1 0000 101 - 050 101 x 21 7 - 0 0019999125 | 0mmE = 1 1000 101 = 0b1 101 × 070 = 1 625 |
| | | | |
| $0x06 = 0.0000.110 = +000.110 \times 2 -7 = 0.000000378$ | 0146 = 0.1000.110 = +081.110 × 2 0 = 1.78 | 0x86 = 1.0000.110 = -080.110 x 2 -7 = -0.008859378 | 0xc6 = 1.1000.110 = -061.110 x 2 0 = -1.18 |
| $0x07 = 0.0000.111 = +060.111 \times 2^{-7} = 0.0068359375$ | $0x47 = 0.1000.111 = +051.111 \times 2^{\circ}0 = 1.875$ | $0x87 = 1.0000.111 = -060.111 \times 2^{-7} = -0.0068359375$ | $0xc7 = 1.1000.111 = -061.111 \times 2^{\circ}0 = -1.875$ |
| $0x08 = 0.0001.000 = +0b1.000 \times 2^{-7} = 0.0078125$ | $0x48 = 0.1001.000 = +0b1.000 \times 2^{-1} = 2.0$ | $0x88 = 1.0001.000 = -0b1.000 \times 2^{-7} = -0.0078125$ | $0xc8 = 1.1001.000 = -0b1.000 \times 2^{-1} = -2.0$ |
| $0x09 = 0.0001.001 = +0b1.001 \times 2^{-7} = 0.0087890625$ | $0x49 = 0.1001.001 = +0b1.001 \times 2^{-1} = 2.25$ | 0x89 = 1.0001.001 = -0b1.001×2 ⁻⁷ = -0.0087890625 | $0xc9 = 1.1001.001 = -0b1.001 \times 2^{-1} = -2.25$ |
| $0x0a = 0.0001.010 = +0b1.010 \times 2^{-7} = 0.009765625$ | $0x4a = 0.1001.010 = +0b1.010 \times 2^{-1} = 2.5$ | $0x8a = 1.0001.010 = -0b1.010 \times 2^{-7} = -0.009765625$ | $0xca = 1.1001.010 = -0b1.010 \times 2^{-1} = -2.5$ |
| $0x0b = 0.0001.011 = +0b1.011 \times 2^{-7} = 0.0107421875$ | $0x4b = 0.1001.011 = +0b1.011 \times 2^{\circ}1 = 2.75$ | $0x8b = 1.0001.011 = -0b1.011 \times 2^{\circ}-7 = -0.0107421875$ | $0xcb = 1.1001.011 = -0b1.011 \times 2^{-1} = -2.75$ |
| $0x0c = 0.0001.100 = +0b1.100 \times 2^{-7} = 0.01171875$ | $0x4c = 0.1001.100 = +0b1.100 \times 2^{\circ}1 = 3.0$ | $0x8c = 1.0001.100 = -0b1.100 \times 2^{\circ}-7 = -0.01171875$ | $0xcc = 1.1001.100 = -0b1.100 \times 2^{2} = -3.0$ |
| $0x04 = 0.0001101 = \pm 0b1101 \times 2^{-7} = 0.0126953125$ | $0x44 = 0.1001.101 = \pm 0b1.101 \times 2^{-1} = 3.25$ | $0x84 = 1.0001.101 = -0h1.101 \times 2^{-7} = -0.0126953125$ | $0xcd = 1\ 1001\ 101 = -0b1\ 101 \times 2^{-1} = -3\ 25$ |
| $0x0e = 0.0001110 = \pm 0b1110 \times 2^{-7} = 0.013671875$ | $0x4e = 0.1001.110 = \pm 0h1.110 \times 2^{2}1 = 3.5$ | 0x8e = 1 0001 110 = -0b1 110 × 2'-7 = -0 013671875 | $0xce = 1.1001.110 = -0b1.110 \times 271 = -3.5$ |
| $0x0f = 0.0001111 = \pm 0b1111 \times 2^{-7} = 0.0146484375$ | $0x4f = 0.1001.111 = \pm 0b1.111 \times 271 = 3.75$ | $0xBf = 1.0001.111 = -0b1.111 \times 2^{-7} = -0.0146484375$ | $0xcf = 1.1001.111 = -0b1.111 \times 271 = -3.75$ |
| | | | |
| 0x10 = 0.0010.000 = +001.000 x 2 *6 = 0.018628 | 0150 = 0.1010.000 = +051.000 × 2 2 = 4.0 | 0x30 = 1.0010.000 = -001.000 x 2 -6 = -0.010628 | 0xd0 = 1.1010.000 = -001.000 x 2 2 = -4.0 |
| $0x11 = 0.0010.001 = +001.001 \times 2 - 6 = 0.017578125$ | $0x51 = 0.1010.001 = +051.001 \times 2.2 = 4.5$ | $0x91 = 1.0010.001 = -061.001 \times 2 - 6 = -0.017578125$ | $0xd1 = 1.1010.001 = -061.001 \times 2.2 = -4.5$ |
| $0x12 = 0.0010.010 = +061.010 \times 2^{-6} = 0.01963125$ | $0x52 = 0.1010.010 = +051.010 \times 2^{-2} = 5.0$ | $0x92 = 1.0010.010 = -061.010 \times 2^{-6} = -0.01953125$ | $0xd2 = 1.1010.010 = -061.010 \times 2.2 = -5.0$ |
| $0x13 = 0.0010.011 = +061.011 \times 2^{-6} = 0.021484375$ | $0x53 = 0.1010.011 = +051.011 \times 2^{-2} = 5.5$ | $0x93 = 1.0010.011 = -061.011 \times 2^{-6} = -0.021484375$ | $0xd3 = 1.1010.011 = -061.011 \times 2.2 = -5.5$ |
| $0x14 = 0.0010.100 = +0b1.100 \times 2^{-6} = 0.0234375$ | $0x54 = 0.1010.100 = +0b1.100 \times 2^{\circ}2 = 6.0$ | $0x94 = 1.0010.100 = -0b1.100 \times 2^{\circ}-6 = -0.0234375$ | $0xd4 = 1.1010.100 = -0b1.100 \times 2^{2} = -6.0$ |
| 0x15 = 0.0010.101 = +0b1.101×2 ⁻⁶ = 0.025390625 | $0x55 = 0.1010.101 = +0b1.101 \times 2^{\circ}2 = 6.5$ | 0x95 = 1.0010.101 = -0b1.101×2 ⁻⁶ = -0.025390625 | 0xd5 = 1.1010.101 = -0b1.101×2 ² = -6.5 |
| $0x16 = 0.0010.110 = +0b1.110 \times 2^{-6} = 0.02734375$ | $0x56 = 0.1010.110 = +0b1.110 \times 2^{-2} = 7.0$ | $0x96 = 1.0010.110 = -0b1.110 \times 2^{-6} = -0.02734375$ | $0xd6 = 1.1010.110 = -0b1.110 \times 272 = -7.0$ |
| $0x17 = 0.0010.111 = +0b1.111 \times 2^{-6} = 0.029296875$ | $0x57 = 0.1010.111 = +0b1.111 \times 2^{-2} = 7.5$ | $0x97 = 1.0010.111 = -0b1.111 \times 2^{-6} = -0.029296875$ | $0xd7 = 1.1010.111 = -0b1.111 \times 272 = -7.5$ |
| $0x18 = 0.0011.000 = +0b1.000 \times 2^{-5} = 0.03125$ | $0x58 = 0.1011.000 = +0b1.000 \times 2^{-3} = 8.0$ | $0x98 = 1.0011.000 = -0b1.000 \times 2^{-5} = -0.03125$ | $0xd8 = 1.1011.000 = -0b1.000 \times 2^{-3} = -8.0$ |
| $0x19 = 0.0011.001 = \pm 0b1.001 \times 2^{\circ} = 0.03515625$ | $0x59 = 0.1011.001 = \pm 0b1.001 \times 2^{\circ}3 = 9.0$ | $0x99 = 1.0011.001 = -0b1.001 \times 2^{\circ} - 5 = -0.03515625$ | $0xd9 = 1.1011.001 = -0b1.001 \times 2^{2}3 = -9.0$ |
| $0x1a = 0.0011.010 = \pm 0b1.010 \times 2^{\circ} = 0.0390625$ | $0x5x = 0.1011.010 = \pm 0b1.010 \times 2^{\circ}3 = 10.0$ | $0x9a = 1.0011.010 = -0b1.010 \times 2^{\circ}-5 = -0.0390625$ | $0 \text{ rds} = 1.1011.010 = -0 \text{ b} 1.010 \times 2^{-3} = -10.0$ |
| 0x1b = 0.0011.011 = +0b1.011 × 27-5 = 0.04296875 | $0xEb = 0.1011.011 = \pm 0b1.011 \times 273 = 11.0$ | $0x9h = 1.0011.011 = -0h1.011 \times 2^{-5} = -0.04296975$ | $0xdb = 1.1011.011 = -0b1.011 \times 27 = -11.0$ |
| | | | |
| 0x1c = 0.0011.100 = +001.100 x 2 -8 = 0.046678 | 0150 = 0.1011.100 = +081.100 × 2 3 = 12.0 | 0x9C = 1.0011.100 = -001.100 x 2 -8 = -0.040878 | 0xdc = 1.1011.100 = -061.100 x 2 3 = -12.0 |
| $0x1d = 0.0011.101 = +001.101 \times 2 - 5 = 0.05078125$ | $0x5d = 0.1011.101 = +051.101 \times 2.3 = 13.0$ | $0x9d = 1.0011.101 = -061.101 \times 2 - 5 = -0.05078125$ | $0xdd = 1.1011.101 = -061.101 \times 2.3 = -13.0$ |
| $0x1e = 0.0011.110 = +0b1.110 \times 2^{2} - 5 = 0.0546875$ | $0x5e = 0.1011.110 = +0b1.110 \times 2^{-3} = 14.0$ | $0x9e = 1.0011.110 = -0b1.110 \times 2^{-5} = -0.0546875$ | $0xde = 1.1011.110 = -0b1.110 \times 2^{-3} = -14.0$ |
| $0x1f = 0.0011.111 = +0b1.111 \times 2^{-5} = 0.05859375$ | $0\pi 5f = 0.1011.111 = +0b1.111 \times 2^{-3} = 15.0$ | $0x9f = 1.0011.111 = -0b1.111 \times 2^{-5} = -0.05859375$ | $0xdf = 1.1011.111 = -0b1.111 \times 2^{-3} = -15.0$ |
| $0x20 = 0.0100.000 = +0b1.000 \times 2^{-4} = 0.0625$ | $0x60 = 0.1100.000 = +0b1.000 \times 2'4 = 16.0$ | $0xa0 = 1.0100.000 = -0b1.000 \times 2^{\circ} - 4 = -0.0625$ | $0xe0 = 1.1100.000 = -0b1.000 \times 2^{2} = -16.0$ |
| $0x21 = 0.0100.001 = +0b1.001 \times 2^{-4} = 0.0703125$ | $0x61 = 0.1100.001 = +0b1.001 \times 2'4 = 18.0$ | $0xa1 = 1.0100.001 = -0b1.001 \times 2^{\circ}-4 = -0.0703125$ | $0xe1 = 1.1100.001 = -0b1.001 \times 2^{2} = -18.0$ |
| $0x22 = 0.0100.010 = +0b1.010 \times 2^{-4} = 0.078125$ | $0x62 = 0.1100.010 = +0b1.010 \times 2^{-4} = 20.0$ | $0xa2 = 1.0100.010 = -0b1.010 \times 2^{-4} = -0.078125$ | $0xe2 = 1.1100.010 = -0b1.010 \times 2^{-4} = -20.0$ |
| 0x23 = 0.0100.011 = +0b1.011×2-4 = 0.0859375 | $0x63 = 0.1100.011 = +0b1.011 \times 2^{-4} = 22.0$ | $0xa3 = 1.0100.011 = -0b1.011 \times 2^{-4} = -0.0859375$ | $0xe3 = 1.1100.011 = -0b1.011 \times 2^{-4} = -22.0$ |
| $0x24 = 0.0100.100 = +0b1.100 \times 2^{-4} = 0.09375$ | $0x64 = 0.1100.100 = +0b1.100 \times 274 = 24.0$ | $0xa4 = 1.0100.100 = -0b1.100 \times 2^{-4} = -0.09075$ | $0xe4 = 1.1100.100 = -0b1.100 \times 274 = -24.0$ |
| $0x25 = 0.0100.101 = +0b1.101 \times 2^{-4} = 0.1015625$ | $0x65 = 0.1100.101 = +0b1.101 \times 274 = 26.0$ | $0xa5 = 1.0100.101 = -0b1.101 \times 2^{-4} = -0.1015625$ | $0xe5 = 1.1100.101 = -0b1.101 \times 2.4 = -26.0$ |
| $0x26 = 0.0100.110 = \pm 0b1.110 \times 2^{\circ} - 4 = 0.109375$ | $0 = 66 = 0.1100.110 = \pm 0 1.110 \times 2.4 = 28.0$ | $0_{xa6} = 1.0100.110 = -0b1.110 \times 2^{\circ} A = -0.109375$ | $0xe6 = 1.1100.110 = -0b1.110 \times 224 = -28.0$ |
| 0x27 - 0.0100 111 - +0b1 111 x 21 4 - 0.1171975 | $0x67 = 0.1100.111 = +001.111 \times 201 = 20.0$ | 0xe7 - 1.0100 111 0b1 111 x 22 40.1171975 | 0xe7 - 1 1100 111 0b1 111 × 27 |
| | | | |
| $0x28 = 0.0101.000 = +001.000 \times 2 - 3 = 0.125$ | $0168 = 0.1101.000 = +001.000 \times 2.8 = 32.0$ | 0x88 = 1.0101.000 = -001.000 x 2 -3 = -0.125 | 0xe8 = 1.1101.000 = -051.000 x 2 5 = -32.0 |
| $0x29 = 0.0101.001 = +061.001 \times 2 - 3 = 0.140625$ | $0x69 = 0.1101.001 = +061.001 \times 2.5 = 36.0$ | $0xa9 = 1.0101.001 = -061.001 \times 2 - 3 = -0.140626$ | $0xe9 = 1.1101.001 = -061.001 \times 2.6 = -36.0$ |
| $0x2a = 0.0101.010 = +061.010 \times 2 - 3 = 0.16626$ | $0x6a = 0.1101.010 = +061.010 \times 2.5 = 40.0$ | $0xaa = 1.0101.010 = -001.010 \times 2 - 3 = -0.15625$ | $0xea = 1.1101.010 = -061.010 \times 26 = -40.0$ |
| $0x2b = 0.0101.011 = +0b1.011 \times 2^{-3} = 0.171875$ | $0x6b = 0.1101.011 = +0b1.011 \times 2^{-5} = 44.0$ | $0xab = 1.0101.011 = -0b1.011 \times 2^{-3} = -0.171875$ | $0 \text{xeb} = 1.1101.011 = -0b1.011 \times 2^{-5} = -44.0$ |
| $0x2c = 0.0101.100 = +0b1.100 \times 2^{-3} = 0.1875$ | $0x6c = 0.1101.100 = +0b1.100 \times 2^{-5} = 48.0$ | $0xac = 1.0101.100 = -0b1.100 \times 2^{-3} = -0.1875$ | $0 \text{xec} = 1.1101.100 = -0b1.100 \times 2\% = -48.0$ |
| $0x2d = 0.0101.101 = +0b1.101 \times 2^{-3} = 0.203125$ | $0\pi 6d = 0.1101.101 = +0b1.101 \times 2^{\circ}5 = 52.0$ | $0xad = 1.0101.101 = -0b1.101 \times 2^{-3} = -0.203125$ | $0xed = 1.1101.101 = -0b1.101 \times 2^{-5} = -52.0$ |
| 0x2e = 0.0101.110 = +0b1.110×2 ⁻³ = 0.21875 | 0x6e = 0.1101.110 = +0b1.110×2'5 = 56.0 | 0xae = 1.0101.110 = -0b1.110×2 ⁻³ = -0.21875 | 0xee = 1.1101.110 = -0b1.110×25 = -56.0 |
| 0x2f = 0.0101.111 = +0b1.111×2~3 = 0.234375 | $0x6f = 0.1101.111 = +0b1.111 \times 2^{\circ}5 = 60.0$ | 0xaf = 1.0101.111 = -0b1.111×2 ⁻³ = -0.234375 | 0xef = 1.1101.111 = -0b1.111×25 = -60.0 |
| $0x30 = 0.0110.000 = +0b1.000 \times 2^{-2} = 0.25$ | $0x70 = 0.1110.000 = +0b1.000 \times 2^{-6} = 64.0$ | $0xb0 = 1.0110.000 = -0b1.000 \times 2^{-2} = -0.25$ | $0xf0 = 1.1110.000 = -0b1.000 \times 2^{-6} = -64.0$ |
| $0x31 = 0.0110.001 = \pm 0b1.001 \times 2^{-2} = 0.28125$ | $0x71 = 0.1110.001 = \pm 0b1.001 \times 276 = 72.0$ | $0yh1 = 1.0110.001 = -0h1.001 \times 2^{-2} = -0.28125$ | $0xf1 = 1.1110.001 = -0h1.001 \times 276 = -72.0$ |
| $0x32 = 0.0110.010 = +0b1.010 \times 2^{-2} = 0.3125$ | $0\pi72 = 0.1110.010 = +0b1.010 \times 276 = 80.0$ | $0xb2 = 1.0110.010 = -0b1.010 \times 2^{-2} = -0.3125$ | $0xf2 = 1.1110.010 = -0b1.010 \times 276 = -80.0$ |
| 0x33 - 0.0110.011 - +051.011 x 27-2 - 0.34375 | 0x73 - 0.1110.011 - +051.011 × 275 - 88.0 | 0xb3 = 1.0110.011 = -0b1.011 × 27-2 = -0.34375 | $0xf3 = 1.1110.011 = -0b1.011 \times 276 = -88.0$ |
| 0+94 - 0.0110.100 - 1.051.100 x 22 0 - 0.975 | 0+74 - 0 1110 100 - 1 001 100 - 000 | 0xbd = 1.0110.001 = 0b1.0010(1 × 2 = 0.000)0 | 0x10 = 1.1110.001 = 001.001/x10 = 00.0 |
| | | | |
| 0x00 = 0.0110.101 = +001.101 x 2 -2 = 0.40626 | 0.75 - 0.1110.101 = +001.101 × 26 = 104.0 | $0.100 = 1.01101101 = -001.101 \times 2 - 2 = -0.40625$ | $3410 - 1.1110.101 = -001.101 \times 28 = -104.0$ |
| 0x36 = 0.0110.110 = +001.110 × 2~2 = 0.4375 | $0176 = 0.1110.110 = +061.110 \times 2.6 = 112.0$ | 0x06 = 1.0110.110 = -001.110×2-2 = -0.43/5 | oxi6 = 1.1110.110 = -061.110×276 = -112.0 |
| $0x37 = 0.0110.111 = +001.111 \times 2^{2} - 2 = 0.46875$ | $0x/7 = 0.1110.111 = +0b1.111 \times 2^{\circ}b = 120.0$ | uxbr = 1.0110.111 = -0b1.111×2*-2 = -0.46875 | $0 \times r r = 1.1110.111 = -061.111 \times 2.6 = -120.0$ |
| $0x38 = 0.0111.000 = +0b1.000 \times 2^{-1} = 0.5$ | $0x78 = 0.1111.000 = +0b1.000 \times 277 = 128.0$ | $0xb8 = 1.0111.000 = -0b1.000 \times 2^{-1} = -0.5$ | $0xf8 = 1.1111.000 = -0b1.000 \times 27 = -128.0$ |
| $0x39 = 0.0111.001 = +0b1.001 \times 2^{\circ} - 1 = 0.5625$ | $0\pi79 = 0.1111.001 = +0b1.001 \times 277 = 144.0$ | $0xb9 = 1.0111.001 = -0b1.001 \times 2^{-1} = -0.5625$ | $0xf9 = 1.1111.001 = -0b1.001 \times 27 = -144.0$ |
| $0x3a = 0.0111.010 = +0b1.010 \times 2^{-1} = 0.625$ | $0x7a = 0.1111.010 = +0b1.010 \times 277 = 160.0$ | 0xba = 1.0111.010 = -0b1.010×2 ⁻ -1 = -0.625 | $0xfa = 1.1111.010 = -0b1.010 \times 27 = -160.0$ |
| $0x3b = 0.0111.011 = +0b1.011 \times 2^{-1} = 0.6875$ | $0\pi7b = 0.1111.011 = +0b1.011 \times 277 = 176.0$ | 0xbb = 1.0111.011 = -0b1.011×2°-1 = -0.6875 | $0xfb = 1.1111.011 = -0b1.011 \times 27 = -176.0$ |
| 0x3c = 0.0111.100 = +0b1.100×2 ⁻¹ = 0.75 | 0x7c = 0.1111.100 = +0b1.100×277 = 192.0 | 0xbc = 1.0111.100 = -0b1.100×2 ⁻ -1 = -0.75 | 0xfc = 1.1111.100 = -0b1.100×27 = -192.0 |
| 0x3d = 0.0111.101 = +0b1.101×2~1 = 0.8125 | 0x7d = 0.1111.101 = +0b1.101×27 = 208.0 | 0xbd = 1.0111.101 = -0b1.101×2 ⁻¹ = -0.8125 | 0xfd = 1.1111.101 = -0b1.101×27 = -208.0 |
| 0x3e = 0.0111.110 = +0b1.110×2-1 = 0.875 | 0x7e = 0.1111.110 = +0b1.110×27 = 224.0 | 0xbe = 1.0111.110 = -0b1.110×2 ⁻¹ = -0.875 | 0xfe = 1.1111.110 = -0b1.110×27 = -224.0 |
| $0x3f = 0.0111.111 = +0b1.111 \times 2^{-1} = 0.9375$ | 0x7f = 0.1111.111 = +Inf | 0xbf = 1.0111.111 = -0b1.111×2'-1 = -0.9375 | 0xff = 1.1111.111 = -Inf |
| | | 0.0010 | |

M. Mikaitis (Leeds)

Narrow range floating point

| Format | precision | min pos. | max pos. | и |
|-------------------|-----------|-----------------|---------------------------------------|------------------|
| binary64 (double) | 53 | 2^{-1022} | $\sim 1.798 	imes 10^{308}$ | 2 ⁻⁵³ |
| binary32 (single) | 24 | 2^{-126} | $\sim 3.403 	imes 10^{38}$ | 2^{-24} |
| tf32 (19-bit) | 11 | 2^{-126} | $\sim 3.401 	imes 10^{38}$ | 2^{-11} |
| bfloat16 | 8 | 2^{-126} | \sim 3.389 $	imes$ 10 ³⁸ | 2 ⁻⁸ |
| binary16 | 11 | 2^{-14} | 65504 | 2^{-11} |
| fp8-E4M3 | 4 | 2 ⁻⁶ | 448 | 2^{-4} |
| fp8-E5M2 | 3 | 2^{-14} | 57344 | 2^{-3} |
| fp6-E2M3 | 4 | 2 ⁰ | 7.5 | 2^{-4} |
| fp6-E3M2 | 3 | 2^{-2} | 28 | 2^{-3} |
| fp4-E2M1 | 2 | 2 ⁰ | 6 | 2^{-2} |

Mixed-precision matrix multipliers

Formats with narrow ranges are available in matrix multiply operation.



Hardware matrix multipliers in mixed precision

• Example above is 4×4 , but dimensions differ across architectures.

• Internal dot product precision, rounding, and subnormal support.

While the input formats have narrow ranges, the output is less constrained on both the precision and range.

M. Mikaitis (Leeds)

Narrow range floating point

Part 1: Basic single-word algorithm

Goal: Given *A* and *B*, matrices in, for example, binary64, multiply them accurately using mixed-precision MMAs.

- Scale input matrices A and B.
- ② Round input matrices to the input format.
- Multiply scaled and rounded A and B in the accumulation format.
- Scale the output matrix.

$$C = \Lambda^{-1} \Big(\mathrm{fl}(\Lambda A) \mathrm{fl}(BM) \Big) M^{-1}$$

- A and M are nonsingular diagonal matrices with diagonal coefficients λ_i and μ_i respectively.
- Scale coefficients λ_i and μ_i are powers of two.

Let θ be the maximum value we can afford in the scaled A and B.

Scaling by powers of two means the maximum entry per row of A or column of B is in $(\theta/2, \theta]$.

We should maximise θ to reduce number of underflows, but at the same time remove possibility of overflow.

Choose:

$$\theta = \min(f_{\max}, \sqrt{F_{\max}/n}).$$

which avoids overflow in the input and in the accumulation of n products.

- Take $A \in \mathbb{R}^{4 \times 4}$ and $B \in \mathbb{R}^{4 \times 4}$.
- Set fp8-E4M3 as the *input format* with $f_{max} = 448$.
- Set binary16 as the accumulation format with $F_{\rm max} = 65504$.
- No subnormal floating-point numbers.
- This gives $\min(448, \sqrt{65504/4}) = \min(448, 127.9687) \approx 127 = \theta$.

Scaling factors

In this case before rounding matrices to the *input format* we need to scale them such that 127 is the maximum value that appears.

- 127 is lower than $f_{\rm max} = 448$ no *input format* overflows.
- $127 \times 127 = 16129$ and if we accumulate four such products we get $64616 < F_{max} = 65504$. No accumulation format overflows.

Take

$$A = \begin{bmatrix} 500 & 1 & 1 & 2^{-6} \\ 128 & 128 & 128 & 128 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}, B = \begin{bmatrix} 1 & 128 & 1 & 1 \\ 1 & 128 & 1 & 1 \\ 1 & 128 & 1 & 1 \\ 1 & 128 & 1 & 1 \end{bmatrix}$$

We have

| | 502.015625 | 64258 | 502.015625 | 502.015625 |
|------|------------|-------|------------|------------|
| ٨D | 512 | 65536 | 512 | 512 |
| AD = | 4 | 512 | 4 | 4 |
| | 4 | 512 | 4 | 4 |

Overflows in the above example if no scaling is applied

(Input) $500 > f_{\max} = 448$ and (output) $65536 > F_{\max} = 65504$.

•

$$C = \Lambda^{-1} (\operatorname{fl}(\Lambda A) \operatorname{fl}(BM)) M^{-1}, \quad \theta = 127$$

Step 1: Scale A and B.

$$\Lambda A = \begin{bmatrix} 2^{-2} & 0 & 0 & 0 \\ 0 & 2^{-1} & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 500 & 1 & 1 & 2^{-6} \\ 128 & 128 & 128 & 128 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 125 & 2^{-2} & 2^{-2} & 2^{-8} \\ 64 & 64 & 64 & 64 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$
$$BM = \begin{bmatrix} 1 & 128 & 1 & 1 \\ 1 & 128 & 1 & 1 \\ 1 & 128 & 1 & 1 \\ 1 & 128 & 1 & 1 \\ 1 & 128 & 1 & 1 \\ 1 & 128 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 2^{-1} & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \end{bmatrix}$$

How the scale coefficients are calculated

For example, take the first row of A. The largest value is 500 and we need to get it below $\theta = 127$. $\lambda_1 = 2^{\lfloor \log_2(127/500) \rfloor} = 2^{-2}$.

M. Mikaitis (Leeds)

$$C = \Lambda^{-1} \Big(\mathrm{fl}(\Lambda A) \mathrm{fl}(BM) \Big) M^{-1}$$

Step 2: Round to the *input format* fp8-E4M3 ($f_{\min} = 2^{-6}$).

$$\mathsf{fl}(\Lambda A) = \mathsf{fl}\left(\begin{bmatrix} 125 & 2^{-2} & 2^{-2} & 2^{-8} \\ 64 & 64 & 64 & 64 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}\right) = \begin{bmatrix} 125 & 2^{-2} & 2^{-2} & \mathbf{0} \\ 64 & 64 & 64 & 64 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$
$$\mathsf{fl}(BM) = \mathsf{fl}\left(\begin{bmatrix} 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \end{bmatrix}\right) = \begin{bmatrix} 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \end{bmatrix}$$

Underflow in the above example

Notice that since subnormals are off, numbers $\leq f_{\min}/2$ will round to zero, causing underflow. This happened to $\Lambda A(1,4) = 2^{-8}$, which resulted from scaling the first row of A, where originally $A(1,4) = 2^{-6}$.

M. Mikaitis (Leeds)

$$C = \Lambda^{-1} \Big(\mathrm{fl}(\Lambda A) \mathrm{fl}(BM) \Big) M^{-1}$$

Step 3: Perform matrix multiply in the accumulation format binary16 (T = 11, $F_{max} = 65504$).

$$\begin{bmatrix} 125 & 2^{-2} & 2^{-2} & 0\\ 64 & 64 & 64 & 64\\ 1 & 1 & 1 & 1\\ 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 64 & 1 & 1\\ 1 & 64 & 1 & 1\\ 1 & 64 & 1 & 1\\ 1 & 64 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 125.5 & 8032 & 125.5 & 125.5\\ 256 & 16384 & 256 & 256\\ 4 & 256 & 4 & 4\\ 4 & 256 & 4 & 4 \end{bmatrix}$$

$$C = \Lambda^{-1} \Big(\mathrm{fl}(\Lambda A) \mathrm{fl}(BM) \Big) M^{-1}$$

Step 4: Undo the scaling.

| [4 | 0 | 0 | 0] | [125. | .5 8 | 3032 | 125.5 | 125.5 | 5] | Γ1 | 0 | 0 | 0] | |
|----|---|---|----|---------------|------|-------|-------|-------|----|----|---|---|----|---|
| 0 | 2 | 0 | 0 | 256 | 51 | 6384 | 256 | 256 | | 0 | 2 | 0 | 0 | |
| 0 | 0 | 1 | 0 | 4 | | 256 | 4 | 4 | | 0 | 0 | 1 | 0 | = |
| 0 | 0 | 0 | 1 | 4 | | 256 | 4 | 4 | | L0 | 0 | 0 | 1 | |
| | | | | Г | 502 | 64256 | 5 502 | 502] | | | | | | |
| | | | | | 512 | 65536 | 5 512 | 512 | | | | | | |
| | | | | | 4 | 512 | 4 | 4 | | | | | | |
| | | | | | 4 | 512 | 4 | 4 | | | | | | |

$$C = \Lambda^{-1} \Big(\mathrm{fl}(\Lambda A) \mathrm{fl}(BM) \Big) M^{-1}$$

Comparison. Our result computed with mixed-precision MMA:

| | 502 | 64256 | 502 | 502 |
|------------------|-----|-------|-----|-----|
| $\Lambda P \sim$ | 512 | 65536 | 512 | 512 |
| $AD \approx$ | 4 | 512 | 4 | 4 |
| | 4 | 512 | 4 | 4 |

And the exact result

| | 502.015625 | 64258 | 502.015625 | 502.015625 |
|------|------------|-------|------------|------------|
| ٨D | 512 | 65536 | 512 | 512 |
| AD = | 4 | 512 | 4 | 4 |
| | 4 | 512 | 4 | 4 |

Part 2: Error analysis

Model 1

The following model describes a mixed-precision MMA operation to compute C = AB, assuming round-to-nearest ties-to-even is used. We have two FP formats:

Input format with precision t, unit roundoff u = 2^{-t}, exponent in [e_{min}, e_{max}], range of normalized values ±[f_{min}, f_{max}]. The maximum distance between any number in [-f_{min}, f_{min}] and the nearest FP number is

 $g_{\min} = egin{cases} f_{\min}/2 & ext{if subnormals are not available} \ uf_{\min} & ext{with subnormals (gradual underflow)} \end{cases}$

Accumulation format with T ≥ t, U = 2^{-T}, exponent in
[E_{min}, E_{max}] ⊇ [e_{min}, e_{max}], and range of norm. numbers ±[F_{min}, F_{max}]. The
maximum distance between any number in [-F_{min}, F_{min}] and the nearest FP
number is

$$G_{\min} = \begin{cases} F_{\min}/2 & \text{if subnormals are not available} \\ UF_{\min} & \text{with subnormals (gradual underflow)} \end{cases}$$

Rounding error model based on [Demmel, 1984]

Take $x, y \in \mathbb{R}$. Assuming no overflows occur, the rounding operator to the *input format* is described as

$$\mathsf{fl}(x) = x(1+\delta) + \eta, \quad |\delta| \le u, \quad |\eta| \le g_{\min}, \quad \eta \delta = 0,$$

and arithmetic operations in the accumulation format as

$$\mathsf{FL}(x \operatorname{op} y) = (x \operatorname{op} y)(1 + \delta) + \eta, \quad |\delta| \leq U, \quad |\eta| \leq G_{\min}, \quad \eta \delta = 0,$$

with op $\in \{+, -, \times, \div\}$.

Here $\eta \delta = 0$ accounts for only one type of error, rounding or overflow.

Single-word algorithm:

$$\|\widehat{C} - AB\|_{\infty} \lesssim \left(2u + nU + 4n^2\theta^{-1}g_{\min} + 4n^2\theta^{-2}G_{\min}\right) \|A\|_{\infty} \|B\|_{\infty}.$$

Part 3: Multi-word algorithm

Step 2: Round to the *input format*, in **double-word representation**.

We will round each ΛA and BM to two fp8-E4M3 matrices instead of one.

Compute the first word (first of the two matrices):

$$A^{(0)} = \mathsf{fl}(\Lambda A) = \mathsf{fl}\left(\begin{bmatrix} 125 & 2^{-2} & 2^{-2} & 2^{-8} \\ 64 & 64 & 64 & 64 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}\right) = \begin{bmatrix} 125 & 2^{-2} & 2^{-2} & \mathbf{0} \\ 64 & 64 & 64 & 64 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$
$$B^{(0)} = \mathsf{fl}(BM) = \mathsf{fl}\left(\begin{bmatrix} 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \end{bmatrix}\right) = \begin{bmatrix} 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \\ 1 & 64 & 1 & 1 \end{bmatrix}$$

Step 2: Round to the *input format* fp8-E4M3, in **double-word representation**.

Compute the second word (rounding/underflow error in the first step):

Extra scaling

Notice the division by $u^1 = 2^{-4}$ before rounding, which is done to reduce underflows in the input format. In general, the multi-word split is $A^{(i)} = fI\left(\left(\Lambda A - \sum_{k=0}^{i-1} u^k A^{(k)}\right)/u^i\right).$

Step 3: Perform matrix products and add them in the *accumulation format* binary16.

p-word case

After splitting ΛA and BM into $A^{(0)}, \ldots, A^{(p-1)}$ and $B^{(0)}, \ldots, B^{(p-1)}$, approximate matrix multiply by p(p+1)/2 products

$$C \approx \Lambda^{-1} \bigg(\sum_{i+j < \rho} u^{i+j} A^{(i)} B^{(j)} \bigg) M^{-1}.$$

In our double-word case

$$A^{(0)}B^{(0)} + uA^{(1)}B^{(0)} =$$

$$C \approx \Lambda^{-1} \bigg(\sum_{i+j < \rho} u^{i+j} A^{(i)} B^{(j)} \bigg) M^{-1}.$$

Step 4: Undo the scaling.

| 4 0 0 | 0 2 0 | 0 0 1 | 0 0 0 | 125.50390625 256 4 | 8032.25 16384 256 | 5 125.50390 256 4 | 625 125.503 25 4 | 390625 56 | 1 0 0 | 0 2 0 | 0 0 1 | 0 0 0 | _ |
|-------------|-------------|-------------|-------------|--------------------------|-------------------------|-------------------------|------------------------|--------------|-------------|-------------|-------------|-------------|---|
| 0 | 0 | 0 | 1 | 4 | 256 | 4 | 4 | t] | 0 | 0 | 0 | 1 | |
| | | | | 502.015625 | 64258 | 502.015625 | 502.015625 |] | | | | | |
| | | | | 512 | 65536 | 512 | 512 | | | | | | |
| | | | | 4 | 512 | 4 | 4 | - AD. | | | | | |
| | | | | 4 | 512 | 4 | 4 | | | | | | |

Single-word algorithm:

$$\|\widehat{C}-AB\|_{\infty} \lesssim \left(2u+nU+4n^2\theta^{-1}g_{\min}+4n^2\theta^{-2}G_{\min}\right)\|A\|_{\infty}\|B\|_{\infty}.$$

p-word algorithm:

$$\begin{split} \|\widehat{C} - AB\|_{\infty} \lesssim \Big((p+1)u^p + 4nu^{p-1}\theta^{-1}g_{\min} \\ &+ (n+p^2)U + 2p(p+1)n^2\theta^{-2}G_{\min}\Big)\|A\|_{\infty}\|B\|_{\infty}. \end{split}$$

Part 4: Numerical experiments

We generate $A \in \mathbb{R}^{10 \times n}$ and $B \in \mathbb{R}^{n \times 10}$ and vary n.

Elements in $[-10^{10}, -10^{-10}] \cup [10^{-10}, 10^{10}].$

Measure the accuracy with $\frac{\|\widehat{C}-C\|_{\infty}}{\|A\|_{\infty}\|B\|_{\infty}}$ where *C* is computed in binary64.

We check with subnormals on/off to detect any improvements due to *gradual underflow*.

We also plot the variants of MMA without any range (exponent) limitations.

Numerical experiments I



Numerical experiments II



Numerical experiments III



Summary

- Underflows in narrow-range FP formats not a problem, provided three types of scaling are used.
- Shown algorithms require minimal bit-level manipulations.
- Simple algorithms to obtain binary32 accuracy.
- Other algorithms exist combining MMA and binary64 addition.

SIAM SISC paper

T. Mary, and M. Mikaitis. Error Analysis of Matrix Multiplication with Narrow Range Floating-Point Arithmetic. **Preprint. Accepted for SIAM SISC**. Mar. 2025. https://bit.ly/42dqpkn.

Slides at http://mmikaitis.github.io/talks

References I

J. Demmel

Underflow and the reliability of numerical software SIAM J. Sci. Comput., 5:4. Dec. 1984.

- P. Blanchard, N. J. Higham, F. Lopez, T. Mary, and S. Pranesh Mixed precision block fused multiply-add: Error analysis and application to GPU tensor cores SIAM J. Sci. Comput., 42:3. Jan. 2020.
- M. Fasi, N. J. Higham, F. Lopez, T. Mary, and M. Mikaitis Matrix Multiplication in Multiword Arithmetic: Error Analysis and Application to GPU Tensor Cores SIAM J. Sci. Comput., 45:1. Feb. 2023.

M. Fasi and M. Mikaitis CPFloat: A C library for Simulating Low-Precision Arithmetic ACM Trans. Math. Software, 49. 2023