

Abgabe von Bruno Stendal, Martin Baer, Lukas Gewinner und Christian Schäfer

4. Aufgabenblatt zum Kurs

TI2: Rechnerarchitektur

von Bernadette Keßler

bis Freitag, den .2022, 10:15 Uhr.

Zahlendarstellung und Rechnen

Beantworten Sie folgende Fragen:

1. Bei Fließkommazahlen ist die Position des Kommas nicht wie bei Festkommazahlen fixiert, das heißt, dass man keine Kompromisse zwischen Wertebereich und Genauigkeit eingehen muss. Auf diese Weise können sowohl sehr große Zahlen, die jedoch nur wenige Nachkommastellen besitzen, als auch sehr genaue Zahlen, die viele Nachkommastellen jedoch wenige Vorkommastellen besitzen, gespeichert werden.
2. Eine Fließkommazahl besteht aus drei Abschnitten, die Mantisse, Charakteristik und S heißen. Dabei steht S für das Vorzeichenbit. Die Fließkommazahl wird nach der Exponentialschreibweise mit $z = \text{Mantisse} \cdot \text{Charakteristik}$ dargestellt. Die Charakteristik besteht intern aus der Basis 2 mit einem Exponenten, welchem je nach Größe Bits reserviert werden.
3. Der Vorteil dem Exponenten mehr Bits zuzuteilen ist, dass der Wertebereich der darstellbaren Zahlen größer wird und der Vorteil der Mantisse mehr Bits zuzuteilen ist, dass man genauere Zahlen darstellen kann.
4. Ein Überlauf entsteht, wenn man eine größere Zahl als die größte darstellbare Zahl darstellen will. Das heißt, es sind mehr Bits gefordert, als für den Exponenten zugeteilt wurden. Beim Unterlauf ist es genau andersherum, wenn man eine kleinere Zahl als die kleinste darstellbare Zahl darstellen will. Anders gesagt entsteht eine Lücke zwischen der kleinsten darstellbaren positiven und größten darstellbaren negativen Zahl rund um die null.
5. Bei IEEE-754 werden bestimmte Standards für die Darstellung von Fließkommazahlen gesetzt. Beispielsweise werden bei 32 Bits die Bits in das „Most significant bit“ für das Vorzeichen, 8 Bits für die Charakteristik/Exponenten und 23 Bits für die Mantisse aufgeteilt. Negative Exponenten werden mittels Offset=127 ermöglicht. Des Weiteren wird die Normalisierung verwendet und durch das hidden Bit etwas Speicher eingespart, aber dann eine extra Darstellung der „0“ gefordert. Overflows über maxreal und minreal werden als +- unendlich dargestellt und eine unzulässige Darstellung NaN (not a number).
6. Betragsmäßig größte bzw. kleinste darstellbare Zahl im IEEE-754 32bit Standard:

Größte darstellbare Zahl: $(2 - 2^{-23}) \cdot 2^{127} \approx 3,410^{38}$

Kleinste darstellbare Zahl

normalisiert: $2^{-126} \approx 1,210^{-38}$

denormalisiert: $2^{-23} \cdot 2^{-126} \approx 1,410^{-45}$

RA-Übung 4

IEEE P 754: $-592.183940 + 0.91213$

Umrechnung in Binär wurde auf letztem Zettel gemacht, hier nicht mehr ausführlich.

$592.183940 \stackrel{\wedge}{=} 1001010000.00101111000101101...$
 $\swarrow \cdot 2^{-9}$, dann noch 23 Bit significant

1,00101000000101111000101 | 101...

↑ ist implizit 23 guard round → aufrunden

also 1. 0 0 0 0 0 0 0 0 0 0 0 0 1 10

characteristic: müssen mit 2^9 multiplizieren, zusätzlich 127 Offset, also 136

$136 = 10001000_2$, Zahl ist negativ, also $\text{sign} = 1$

Es ergibt sich:

$1 \mid 10001000 \mid 00101000000101111000110$
 S Char. signif.

$$0.91213 \hat{=} 0.\underset{\substack{\text{Bin.} \\ 2^1}}{111010011000000101011010} \underbrace{000001}_{\text{abrunden}} \dots$$

character: 2^{-1} , offset 127 \rightarrow 126 $\hat{=}$ 01111110, sign = 0

↳ 0 | 0111110 | 110100110000010101010

Addition (bzw. Subtraktion) $2^9, 2^{-1} \rightarrow 9 - (-1) = 10 \rightarrow$ Rechts-Shift von 0,9.1...

Handwritten notes:

implizite 1 →

→ sticky gesetzt

G.R.S

SideNote: Bin der Meinung in der VL ist was falsch, ich brauche ja bei der Papier-Rechnung 1+23 Bits wegen der impliziten 1? und dann eben erst guard, round, sticky

→ sticky gesetzt

G.R.S

→ abgerunden, ist schon

$\text{sign} = 1$ (da Zahlen in Rechnung getauscht)

↳ 1 | 10001000 | 00100111101011100000101
 2^9 (wie oben)

$$4 \hat{=} -1001001111,01011100000101_2 \stackrel{=}{=} 591,35968017...$$

$$3981.1729 \times -2.91762$$

$$3981.1729 \hat{=} \underbrace{111110001101001011000100}_{\cdot 2^{11}} | 0011 \dots \text{abrunden}$$

$$\text{character: } 11 + 127 = 138, \text{ sign} = 0$$

$$0 | 10001010 | 11110001101001011000100$$

$$(-)2.91762 \hat{=} \underbrace{101110101011101001001001}_{\cdot 2^{-1}} | 00111 \text{ abrunden}$$

$$\text{character: } 1 + 127 = 128, \text{ sign} = 1$$

$$1 | 10000000 | 01110101011101001001001$$

Sorry, aber die Multiplikation der Mantissen werde ich mal nicht aufschreiben...

$$\text{Ergebnis ist } 1011010101111000110010 | 011 \dots \text{abrunden}$$

um eine Stelle verschieben \Rightarrow in character: +1

$$\text{sign} = 1 \text{ xor } 0 \rightarrow 1 \quad \text{character: } \begin{array}{r} 10001010 \\ + 10000000 \\ - 01111111 \end{array} \} = 00000001$$

10001011 aber +1 wegen Normalisierung oben, also

$$10001100$$

$$1 | 10001100 | 0110101011111000110010$$

$$(-) = 140 \rightarrow 2^{13}, \text{ also } 1011010101111, 1000110010$$

$$\hat{=} -11615,548828125.$$

Floating Point Rechner

Implementieren Sie eine Floating Point Unit in Software. Zwei gegebene 32-bit IEEE-754 Zahlen (operand1 und operand2) sollen addiert werden können. Das Ergebnis soll ebenfalls im 32-bit IEEE-754 Floating Point Format zurückgegeben werden. Verwenden Sie keine Floating Point Register / Befehle außer MOVD.

```
1 ;Bruno Stendal, Martin Baer, Lukas Gewinner, Christian Schaefer
2         global calc_add
3
4 calc_add:
5         mov rax, 0;
6         movd edx, xmm0;
7         shr edx, 31;
8         mov ah, dl;
9
10        movd edx, xmm0;
11        shl edx, 1;
12        shr edx, 24;
13        mov al, dl;
14
15        movd edx, xmm0;
16        shl edx, 9;
17        shr edx, 9;
18        mov esi, edx;
19        add esi, 8388608;
20
21        mov rcx, 0;
22        mov rdx, 0;
23
24        movd edx, xmm1;
25        shr edx, 31;
26        mov ch, dl;
27
28        movd edx, xmm1;
29        shl edx, 1;
30        shr edx, 24;
31        mov cl, dl;
32
33        movd edx, xmm1;
34        shl edx, 9;
35        shr edx, 9;
36        mov edi, edx;
37        add edi, 8388608;
38
39        ;ah sign a
40        ;al expo a
```

```

41         ;esi mat a
42
43         cmp ah, 1;
44         jne .matisse;
45         xor esi, 0xff;
46
47         ;ch sign b
48         ;cl expo b
49         ;edi mat b
50
51         cmp ch, 1;
52         jne .matisse;
53         xor edi, 0xff;
54
55         mov rdx, 0;
56
57     .matisse:
58         cmp al, cl;
59         jl .shiftA;
60         jg .shiftB;
61         je .noShift;
62
63     .shiftA:
64
65         mov rdx, 0;
66         mov dl, cl;
67         mov cl, 0;
68
69         mov cl, dl;
70         sub cl, al;
71
72         mov al, cl;
73         mov edx, 1;
74
75         shr esi, cl;
76
77         mov al, dl;
78         mov cl, dl;
79
80         jmp .noShift;
81
82     .shiftB:
83
84         mov rdx, 0;
85         mov dl, cl;
86         mov cl, 0;
87

```

```

88     mov cl, al;
89     sub cl, dl;
90
91     shr edi, cl;
92
93     mov cl, dl;
94     jmp .noShift;
95
96 .noShift:
97
98     mov rdx, 0;
99     mov edx, esi;
100    add edx, edi;
101
102
103    mov rbx, 0;
104    mov ebx, edx;
105    shr ebx, 24;
106
107    cmp ebx, 0;
108    jne .expoPlus;
109    je .end;
110
111 .expoPlus:
112     cmp ebx, 0;
113     jne .loop;
114     je .end;
115 .loop:
116     add al, 1;
117     shr rbx, 1;
118     jmp .expoPlus;
119
120 .end:
121     shl edx, 9;
122     shr edx, 9;
123
124     cmp ah, 1;
125     jne .endBuild
126     xor edx, 0xff;
127     cmp ah, 1;
128     je .endBuild
129
130     cmp ch, 1;
131     jne .endBuild
132     xor edx, 0xff;
133
134 .endBuild:

```

```
135
136     mov rbx, 0;
137     mov bh, ah;
138     mov bl, al;
139     shl rbx, 23;
140     add ebx, edx;
141     movd xmm0, ebx;
142     ret
```