Verlässliche Systemsoftware

Übungen zur Vorlesung

Festkommaarithmetik

Phillip Raffeck, Simon Schuster, Peter Ulbrich

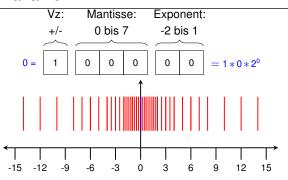
Technische Universität Dortmund Lehrstuhl für Informatik 12 (Arbeitsgruppe Systemsoftware) https://sys.cs.tu-dortmund.de

Wintersemester 2020



Raffeck, Schuster, Ulbrich VSS (WS20) 1

Fließkommazahlen

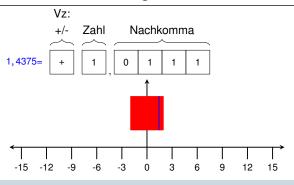


IEEE 754

- Noch komplexer:
 - normalisierte/denormalisierte Darstellung
 - Rundung, Fehlersemantik, ...
 - NaN, ∞, ...
- https://ieeexplore.ieee.org/document/4610935/



Festkommazahlen: Grundlagen



C-Standard und Zahlendarstellung

- Zahlendarstellung im Standard nicht festgelegt:
 - Einerkomplement
 - Vorzeichen und Magnitude
 - Zweierkomplement
- Heute meist Zweierkomplement ~ kein dediziertes Vorzeichenbit



Festkommaarithmethik - Motivation

```
func:
                                          {r7, lr}
                                    push
                                    sub
                                          sp, #8
                                    add
                                          r7, sp, #0
                                    ldr
                                          r3, [pc, #28]; float a
                                    str
                                          r3, [r7, #4]
                                    1dr
                                          r3. [pc. #28] : float b
float func(void){
                                          r3, [r7, #0]
                                    str
  volatile float a = 23.42:
                                    1dr
                                          r2. [r7. #4]
  volatile float b = 12.34:
                                          r3. [r7, #0]
                                    ldr
  return a * b:
                                    adds
                                          r0, r2, #0;
                                                       Param 1
                                    adds r1, r3, #0; Param 2
                                    bl
                                          3a6c <__aeabi_fmul>
                                    adds r3, r0, #0
                                    adds r0, r3, #0
Setup
                                    mov
                                          sp, r7
Plattform: ARM Cortex-M0+
                                    add
                                          sp. #8
                                          {r7, pc}
                                    pop
Compiler: arm-gcc
```

- Funktion __aeabi_fmul : 300 Zeilen Assembler
- Keine Fließkommaeinheit (engl. floating-point unit, FPU) vorhanden
- Emulation der Fließkommaarithmetik in Software



Festkommaarithmetik – Q-Notation

- Mikrocontroller ohne Fließkommaeinheit
- Kein EAN für Fließkommazahlen
 - Festkommaarithmethik mit Ganzzahlen
- Zahlenformat häufig in Q-Notation [1] angegeben
- $Qm.n \sim$ Festkommazahl mit
 - *m* Bit vor dem Komma, *n* nach dem Komma, ein Vorzeichenbit
 - Wertebereich: $[-2^m, 2^m 2^{-n}]$
 - Auflösung: 2⁻ⁿ
- Implementierung für Übungsaufgabe vorgegeben

Implementierung als Integer

→ passendes Q-Format ist anwendungsspezifisch



Q-Notation – Beziehung zu Fließkommazahlen

von Fließkomma nach Qm.n

- 1 Multiplikation mit 2ⁿ
- 2 Runden auf die nächste Ganzzahl

von Qm.n nach Fließkomma

- Umwandlung in Fließkommazahl ~ cast
- 2 Multiplikation mit 2^{-n}



Operationen – Addition/Subtraktion

Addition und Subtraktion wie bei Ganzzahlen

Addition

Subtraktion

Operationen – Multiplikation/Division

Braucht Zwischenergebnis von doppelter Bitbreite


```
a \cdot b
\stackrel{Q.n}{=} (a \cdot 10^n) \cdot (b \cdot 10^n)
= (a \cdot b) \cdot 10^{2n}
\neq (a \cdot b) \cdot 10^n
```

$$\frac{a}{b} \stackrel{Q.n}{=} \frac{a \cdot 10^n}{b \cdot 10^n}$$
$$= \frac{a}{b} \neq \frac{a}{b} \cdot 10^n$$

- Siehe Implementierung in fixedpoint.c
- Vorsicht: Rundungsfehler durch Transformationen



Literatur



Erick L. Oberstar.

Fixed-point representation & fractional math.

Technical report, Oberstar Consulting, August 2007.

