

## Zentralübung Rechnerstrukturen: Low-Power-Entwurf und Leistungsbewertung

### 2. Aufgabenblatt – Musterlösung

#### Low-Power-Entwurf

- Spannungsabsenkung:  $U : 0,8V \leftrightarrow 5V \Rightarrow U^2 : 0,64/25 = 0,0256$   
Frequenzerhöhung:  $f : 3000MHz/1MHz = 3000$   
Aus  $P \sim U^2 * f$  resultiert eine Zunahme der elektrischen Leistung um den Faktor  $0,0256 * 3000 = 76,8$   
Spannungsabsenkung:  $U : 0,6V \leftrightarrow 0,8V \Rightarrow U^2 : 0,36/0,64 = 0,5625$   
Frequenzerhöhung:  $f : 3,5GHz/3GHz \approx 1,167$   
 $P \sim U^2 * f = 0,5625 * 1,167 \approx 0,66$  bedeutet eine Abnahme der notwendigen elektrischen Leistung.
- Möchte man Prozessoren übertakten, so ist es nötig, dafür zu sorgen, dass die Taktflanken schneller steigen und so schneller gültige Signallevel vorliegen. Spannungserhöhung führt zu schnellerem Laden von  $C_{eff}$  und damit steileren Flanken.  
 $P_{switching} = C_{eff} * U^2 * f$   
Nachteil: Spannung fließt quadratisch in Formel ein.
- Aufgrund der immer weiter ansteigenden Integrationsdichte spielen mittlerweile die *Leckströme* eine erhebliche Rolle bei der Leistungsaufnahme. Wegen der Leckströme führt eine Verkleinerung der Strukturen nicht automatisch zu einer Reduzierung der Stromaufnahme, was eine höhere Taktung ermöglicht. Leckströme steigen mit höherer Temperatur.
- Ermitteln der Schaltwahrscheinlichkeit der gesamten Schaltung  $\mathbb{P}_{Schalt\ Gesamt}$  über die einzelnen Schaltwahrscheinlichkeiten  $\mathbb{P}_{Schalt\ Gatter}$  und diese über die Signalwahrscheinlichkeiten  $\mathbb{P}_{Gatter}(Ausgang = 0)$  bzw.  $\mathbb{P}_{Gatter}(Ausgang = 1)$ .

$$\begin{aligned}\mathbb{P}_{Schalt} &= \mathbb{P}(0 \rightarrow 1 \vee 1 \rightarrow 0) = \mathbb{P}(0 \rightarrow 1) + \mathbb{P}(1 \rightarrow 0) \\ &= \mathbb{P}(0) * \mathbb{P}_{neu}(1) + \mathbb{P}(1) * \mathbb{P}_{neu}(0) = \mathbb{P}(0) * \mathbb{P}(1) + \mathbb{P}(0) * \mathbb{P}(1) \\ &= 2 * \mathbb{P}(1) * \mathbb{P}(0) = 2 * \mathbb{P}(1) * (1 - \mathbb{P}(1))\end{aligned}$$

Betrachtung des ODER-Gatters:

- Signalwahrscheinlichkeit:

$$\begin{aligned}\mathbb{P}_{\text{Ausgang}}(1) &= 1 - \mathbb{P}_{\text{Ausgang}}(0) \\ &= 1 - \mathbb{P}(a = 0 \wedge b = 0) \\ &= 1 - \left(1 - \frac{1}{4}\right) * \left(1 - \frac{3}{4}\right) = 1 - \frac{3}{4} * \frac{1}{4} = \frac{13}{16}\end{aligned}$$

Alternative (direkte) Berechnung:

$$\begin{aligned}\mathbb{P}_{\text{Ausgang}}(1) &= \mathbb{P}(a = 1 \wedge b = 0) + \mathbb{P}(a = 0 \wedge b = 1) + \mathbb{P}(a = 1 \wedge b = 1) \\ &= \frac{1}{4} * \frac{1}{4} + \frac{3}{4} * \frac{3}{4} + \frac{1}{4} * \frac{3}{4} \\ &= \frac{1}{16} + \frac{9}{16} + \frac{3}{16} = \frac{13}{16}\end{aligned}$$

- Schaltwahrscheinlichkeit:

$$\begin{aligned}\mathbb{P}_{\text{Schalt}} &= 2 * \mathbb{P}(1) * (1 - \mathbb{P}(1)) \\ &= 2 * \frac{13}{16} * \left(1 - \frac{13}{16}\right) = \frac{2 * 13 * 3}{16 * 16} = \frac{39}{128}\end{aligned}$$

5. Verwenden der Summe der Schaltwahrscheinlichkeiten als Metrik um beide Varianten zu vergleichen.

Variante 1:

- Beide linken Gatter:  $\mathbb{P}_{\text{links}}(1) = \frac{1}{2} * \frac{1}{2} = \frac{1}{4}$ ,  $\mathbb{P}_{\text{Schalt-links}} = 2 * \frac{1}{4} * \frac{3}{4} = \frac{3}{8}$
- Rechtes Gatter: Signalwahrscheinlichkeiten für Eingänge des rechten Gatters = Ausgangssignalwahrscheinlichkeiten der linken Gatter  
 $\mathbb{P}_{\text{rechts}}(1) = \frac{1}{4} * \frac{1}{4} = \frac{1}{16}$ ,  $\mathbb{P}_{\text{Schalt-rechts}} = 2 * \frac{1}{16} * \frac{15}{16} = \frac{15}{128}$
- $\text{Summe}_{\text{Schaltw'keiten}} = \frac{3}{8} + \frac{3}{8} + \frac{15}{128} = \frac{111}{128}$

Variante 2:

- Linkes Gatter:  $\mathbb{P}_{\text{links}}(1) = \frac{1}{2} * \frac{1}{2} = \frac{1}{4}$ ,  $\mathbb{P}_{\text{Schalt-links}} = 2 * \frac{1}{4} * \frac{3}{4} = \frac{3}{8}$
- Mittleres Gatter:  $\mathbb{P}_{\text{mitte}}(1) = \frac{1}{2} * \frac{1}{4} = \frac{1}{8}$ ,  $\mathbb{P}_{\text{Schalt-mitte}} = 2 * \frac{1}{8} * \frac{7}{8} = \frac{7}{32}$
- Rechtes Gatter:  $\mathbb{P}_{\text{rechts}}(1) = \frac{1}{2} * \frac{1}{8} = \frac{1}{16}$ ,  $\mathbb{P}_{\text{Schalt-rechts}} = 2 * \frac{1}{16} * \frac{15}{16} = \frac{15}{128}$
- $\text{Summe}_{\text{Schaltw'keiten}} = \frac{3}{8} + \frac{7}{32} + \frac{15}{128} = \frac{91}{128}$

Damit kann aus der höheren Summe der Schaltwahrscheinlichkeiten in Variante 1 ein höherer Leistungsverbrauch resultieren, da es wahrscheinlicher ist, dass ein beliebiges Gatter schaltet und somit zusätzliche Leistung aufnimmt. Die Schaltstruktur hat jedoch eine geringere Durchlaufzeit (2 Ebenen) im Gegensatz zu Variante 2 (3 Ebenen).

## Leistungsbewertung

- Die Zykluszeit hängt von der Organisation und der Technologie ab.  
Die Anzahl der Instruktionen ist bedingt durch die Befehlssatzarchitektur und die Güte des Compilers.  
Die Zyklen pro Instruktion werden durch die Organisation und die Befehlssatzarchitektur beeinflusst.

$$2. \quad f = \frac{i \cdot CPI}{t}, \quad MIPS = \frac{f}{CPI \cdot 10^6}$$

$$f_A = \frac{3,5 \cdot 10^6 \cdot \frac{7}{5}}{2 \cdot 10^{-3} s} = 2450 \text{ MHz} \quad MIPS_A = \frac{2,45 \cdot 10^9}{\frac{7}{5} \cdot 10^6 s} = 1750 \text{ MIPS}$$

$$f_B = \frac{1,5 \cdot 10^6 \cdot \frac{3}{2}}{2 \cdot 10^{-3} s} = 1125 \text{ MHz}, \quad MIPS_B = \frac{1,125 \cdot 10^9}{\frac{3}{2} \cdot 10^6 s} = 750 \text{ MIPS}$$

Es ist Prozessor B zu wählen, weil

- ohne Berechnung: Gleich schnell in der Abarbeitung bei wesentlich weniger Instruktionen (1,5 vs. 3,5 Mio Instruktionen)
- halbe Taktfrequenz ( $P \sim U^2 \cdot f$ , Fertigung)

### 3. Benchmark-Berechnung

- Anzahl Instruktionen:

$$i = \sum i_{typ} = (300 + 75 + 150 + 25) \cdot 10^3 = 550.000$$

- Taktzyklen:

$$c = \sum i_{typ} \cdot c_{typ} = (300 \cdot 1 + 75 \cdot 2 + 150 \cdot 3 + 25 \cdot 4) \cdot 10^3 = 1.000.000$$

- Zykluszeit bei 4GHz Taktfrequenz:

$$t = \frac{1}{f} = \frac{1}{4 \text{ GHz}} = 0,25 \cdot 10^{-9} s = 0,25 \text{ ns}$$

- Ausführungszeit:

$$t_{exec} = c \cdot t_{cyc} = 1000 \cdot 10^3 \cdot 0,25 \cdot 10^{-9} = 250 \cdot 10^{-6} s = 250 \mu s$$

- CPI:

$$CPI = \frac{c}{i} = \frac{1000 \cdot 10^3}{550 \cdot 10^3} = \frac{100}{55} = \frac{20}{11} \approx 1,82$$

- MIPS:

$$MIPS = \frac{i}{t \cdot 10^6} = \frac{550.000}{250} = 2200$$

- MFLOPS: wie MIPS, wobei Anzahl der Befehle und Ausführungszeit nur für ließkommaberechnung

$$MFLOPS = \frac{75.000}{(75.000 \cdot 2) \cdot (0,25 \cdot 10^{-9}) \cdot 10^6} = \frac{1}{0,5 \cdot 10^{-3}} = 2000$$

4. (vergl. Hennessy and Patterson, Computer Architecture A Quantitative Approach, 4. Auflage, S. 43-44.)

Es ändern sich nur die Zyklen pro Instruktion, Taktrate und Anzahl der Instruktionen bleiben gleich.

Der unoptimierte CPI-Wert errechnet sich nach:

$$CPI_{base} = \sum_{i=1}^n CPI_i * Anteil_i = (4 * 25\%) + (1,33 * 75\%) = 2,0$$

Die Zyklen pro Instruktion mit neuem FPSQR:

$CPI_a$  kann durch Abziehen der gesparten Zyklen erfolgen:

$$\begin{aligned} CPI_a &= CPI_{base} - 0,02 * (CPI_{oldFPSQR} - CPI_{newFPSQR}) \\ &= 2,0 - 0,02 * (20 - 2) = 1,64 \end{aligned}$$

Die Alternative mit dem neuen  $CPI_{FP}$ -Wert errechnet sich analog zum  $CPI_{base}$ :

$$CPI_b = (1,33 * 75\%) + (2,5 * 25\%) = 1,625$$

Aufgrund des geringeren CPI-Werts bietet sich die Alternative b mit den verbesserten Zyklen pro Gleitkommaoperation an.

Berechnung des Gewinns (Speedup) durch die Verwendung der Alternative (b) gegenüber dem vorherigen System (base):

$$\begin{aligned} Speedup_{(b)} &= \frac{CPU\ time_{base}}{CPU\ time_b} \\ &= \frac{i * Taktrate * CPI_{base}}{i * Taktrate * CPI_b} \\ &= \frac{CPI_{base}}{CPI_b} \end{aligned}$$

Eingesetzt ergibt sich:

$$Speedup_b = \frac{2,00}{1,62} \approx 1,23$$

→ Alternative (b) ist 1,23-mal schneller als das bisherige System.

5. a) Der Laufzeitunterschied zwischen dem Base und Peak-Setup ist in den erlaubten Optimierungen zu suchen. Während Base nur konservative Standardoptimierungen erlaubt und gleiche Compileroptionen für alle Benchmarks vorschreibt, erlaubt Peak das aggressive Optimieren für die individuelle Architektur.

Für 483.xalancbmk fällt auf, dass die Laufzeitunterschiede vernachlässigbar sind, dies lässt zwei Schlüsse zu:

- entweder waren die durchgeführten Optimierungen nicht wirkungsvoll,
- oder weitere Optimierungen wurden nicht angestrebt.

Ein Blick in die Sektion Peak Optimization Flags der Webseite verrät, dass außer `basepeak=yes` (welches nur die Anzahl der laufenden Kopien des Programms auf dem System beeinflusst) keine Compileroptimierungen angestrebt wurden. Dies steht im Gegensatz zu allen weiteren Programmen, die für den Peak-Lauf mit aggressiven Compileroptimierungen übersetzt wurden.

- b) Es gilt:  $SPEC_{ratio} = \frac{Referenzzeit_x}{Laufzeit_x \text{ auf Testsystem}}$  für einen Benchmark  $x$ .

Somit ergibt sich durch Umstellen und Einsetzen der  $SPEC_{ratio}$  und der Laufzeit aus der Tabelle:

$$Referenzzeit_{462.libquantum} = 613 * 33,8 \text{ s} = 20719,4 \text{ s}$$

- c) Die Suche ergibt *Ultra Enterprise 2* von Sun Microsystems. Hierauf weist die errechnete Referenzlaufzeit des ausgewählten Benchmarks hin, die annähernd übereinstimmen:

Benchmark	Referenzzeit <sub>errechnet</sub>	Laufzeit <sub>Ultra Enterprise 2</sub>
462.libquantum	20719,4	20704

Auch der  $SPEC_{int\_base}^{2006} = 1.00$  spricht dafür.

6. a) Bedienzeiten:  $X_i = t_{\text{Zugriff}} + t_{\text{Übertragung}}$

$$X_1 = 12 \text{ ms} + \frac{100 \text{ kB}}{6000 \text{ kB/s}} = 28,67 \text{ ms}$$

$$X_2 = 10 \text{ ms} + \frac{100 \text{ kB}}{7500 \text{ kB/s}} = 23,33 \text{ ms}$$

$$X_3 = 8 \text{ ms} + \frac{100 \text{ kB}}{8000 \text{ kB/s}} = 20,5 \text{ ms}$$

- b) Maximaler Durchsatz:  $D_{i\max} = \frac{1}{X_i}$

$$D_{1\max} = \frac{1}{28,67 \text{ ms}} = 34,88 \frac{1}{\text{s}}$$

$$D_{2\max} = \frac{1}{23,33 \text{ ms}} = 42,86 \frac{1}{\text{s}}$$

$$D_{3\max} = \frac{1}{20,5 \text{ ms}} = 48,78 \frac{1}{\text{s}}$$

Nur Platten mit  $D_{\max} > A$  können eingesetzt werden, da sonst die Festplatte nicht genügend Zeit hat, um alle Aufträge rechtzeitig zu bedienen. Aufgrund von  $A = 40/s$  sind somit nur die Platten 2 und 3 einsetzbar.

- c) Auslastung:  $U_i = D/D_{i\max} = D * X_i$ , hier  $D = A$

$$U_2 = D * X_2 = 40 \frac{1}{\text{s}} * 23,33 \text{ ms} = 0,93, \text{ d.h. } 93 \% \text{ Auslastung}$$

$$U_3 = D * X_3 = 40 \frac{1}{\text{s}} * 20,5 \text{ ms} = 0,82, \text{ d.h. } 82 \% \text{ Auslastung}$$

- d) Gesetz von Little:  $Q = W * D$

Q: Anzahl von Aufträgen in der Warteschlange

W: Wartezeit

D: Durchsatz

d.h.  $W_i = \frac{Q_i}{D}$ , wobei abermals gilt  $D = A$

$$W_2 = \frac{Q_2}{D} = \frac{3}{40/s} = 75 \text{ ms}$$

$$W_3 = \frac{Q_3}{D} = \frac{2}{40/s} = 50 \text{ ms}$$

Reaktionszeit des Gesamtsystems aus Warteschlange und Festplatte:

- $\text{Reaktionszeit}_i = \text{Wartezeit}_i + \text{Bedienzeit}_i$

- einsetzen ergibt:

$$\text{Reaktionszeit}_2 = 75 \text{ ms} + 23,33 \text{ ms} = 98,33 \text{ ms}$$

$$\text{Reaktionszeit}_3 = 50 \text{ ms} + 20,5 \text{ ms} = 70,5 \text{ ms}$$

Damit ist das System mit Platte 3 vorzuziehen, da es schneller reagiert.