

Technische Grundlagen der Informatik – Kapitel 8

Prof. Dr. Andreas Koch
Fachbereich Informatik
TU Darmstadt



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Kapitel 8: Themen

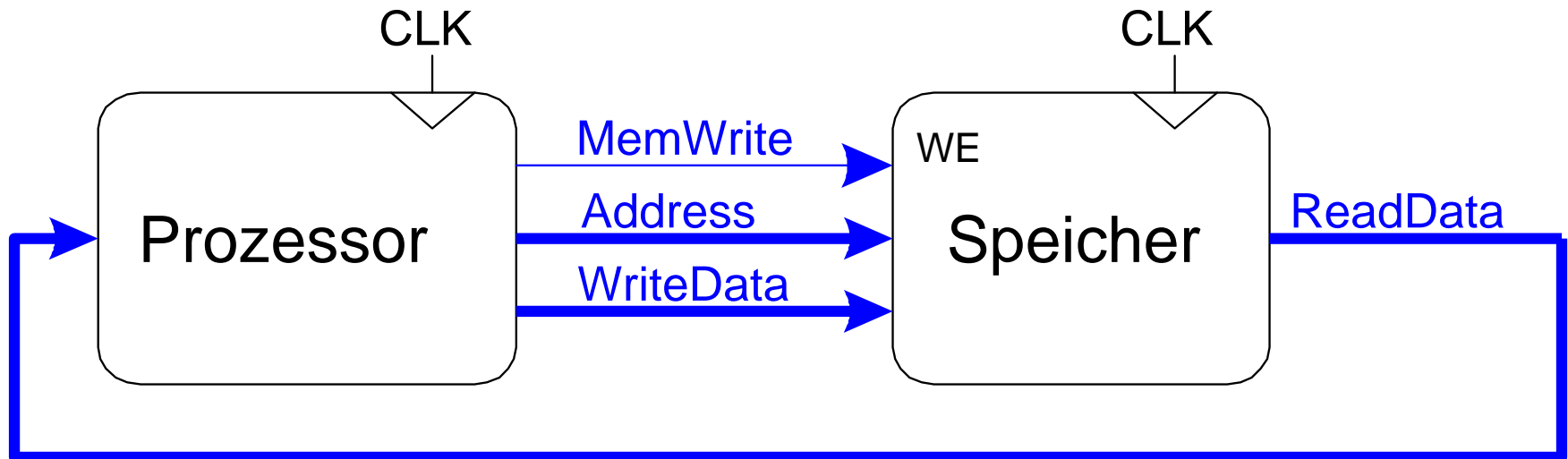


TECHNISCHE
UNIVERSITÄT
DARMSTADT

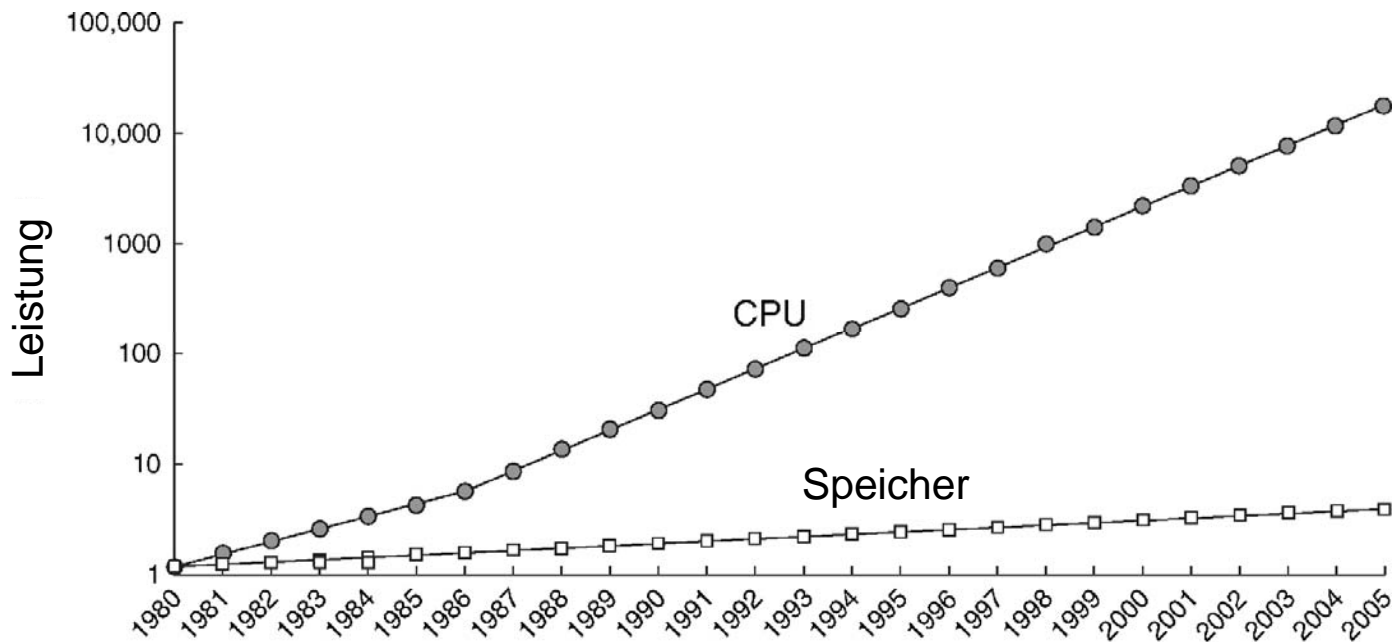
- **Einleitung**
- **Leistungsvergleich von Speichersystemen**
- **Caches**
- **Virtueller Speicher**
- **Speichereinblendung von Ein-/Ausgabegeräten**
- **Zusammenfassung**

- Rechenleistung hängt ab von:
 - Prozessorleistung
 - Leistung des Speichersystems

Speicherschnittstelle



- Annahme bisher in der Vorlesung: Speicherzugriffe dauern 1 Takt
- Ist aber seit den 1980'er Jahren **nicht mehr** wahr



Herausforderungen beim Entwurf von Speichersystemen



TECHNISCHE
UNIVERSITÄT
DARMSTADT

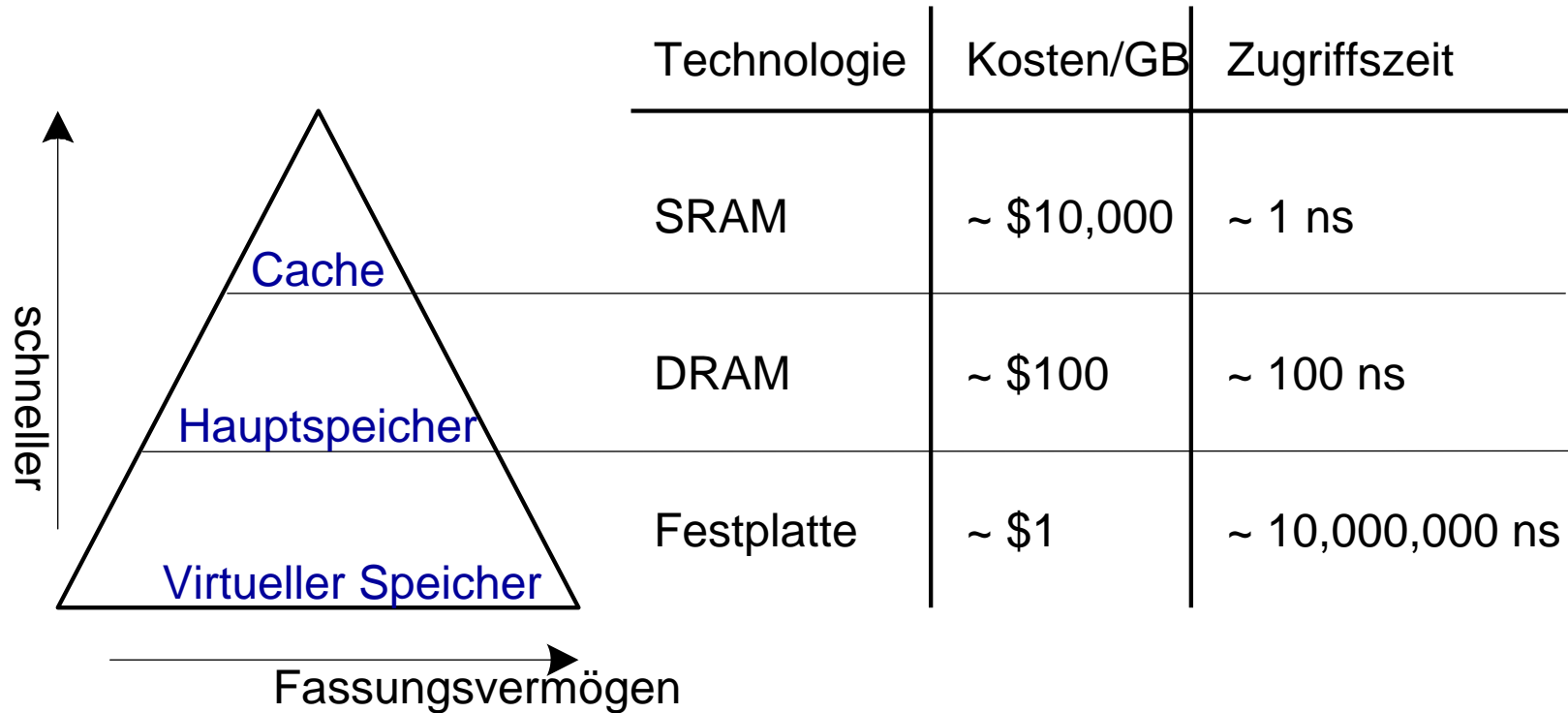
- Speichersystem soll so schnell sein wie Prozessor
 - Zumindest dem Anschein nach ...
- **Idealer** Speicher:
 - Schnell
 - Billig
 - Hohes Fassungsvermögen

Praktisch: Nur zwei von drei Eigenschaften realisierbar!

- Verwende **Hierarchie** von Speichern



Speicherhierarchie: Beispiele für Ebenen



- Nutze Lokalität zur Beschleunigung von Speicherzugriffen aus
- **Zeitliche Lokalität:**
 - Ein gerade benutztes Datum wird wahrscheinlich bald wieder gebraucht
 - **Ausnutzung:** gerade benutzte Daten in höheren Ebenen der Speicherhierarchie halten
- **Räumliche Lokalität:**
 - Um ein benutztes Datum herum liegende Daten werden wahrscheinlich auch bald gebraucht
 - **Ausnutzung:** Beim Zugriff auf ein Datum auch benachbarte Daten in höhere Ebenen der Speicherhierarchie bringen

- **Treffer (*hit*):** Datum wird auf dieser Ebene der Speicherhierarchie gefunden
- **Verfehlt (*miss*):** nicht gefunden (suche auf tieferer Ebene)

Hit-Rate $= \# \text{ hits} / \# \text{ Speicherzugriffe}$
 $= 1 - \text{Miss Rate}$

Miss-Rate (*MR*) $= \# \text{ misses} / \# \text{ Speicherzugriffe}$
 $= 1 - \text{Hit Rate}$

- **Durchschnittliche Speicherzugriffszeit (*average memory access time, AMAT*):** Durchschnittliche Zeit die der Prozessor braucht, um auf ein Datum zuzugreifen

$$\text{AMAT} = t_{\text{cache}} + MR_{\text{cache}} (t_{MM} + MR_{MM} t_{VM})$$

MM = *main memory*, Hauptspeicher

VM = *virtual memory*, Virtueller Speicher

Beispiel 1: Leistung eines Speichersystems



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Ein Programm hat 2.000 Load- und Store-Befehle
 - 1.250 der Daten werden im Cache vorgefunden
 - Der Rest kommt aus anderen Ebenen der Speicherhierarchie
-
- **Was sind die Hit- und Miss-Rates des Caches?**

Beispiel 1: Leistung eines Speichersystems

- Ein Programm hat 2000 Load- und Store-Befehle
- 1250 der Daten werden im Cache vorgefunden
- Der Rest kommt aus anderen Ebenen der Speicherhierarchie
- Was sind die Hit- und Miss-Rates des Caches?

$$\text{Hit Rate} = 1250/2000 = \mathbf{0,625}$$

$$\text{Miss Rate} = 750/2000 = \mathbf{0,375} = 1 - \text{Hit Rate}$$

Beispiel 2: Leistung eines Speichersystems



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Annahme: Prozessor hat zwei Hierarchieebenen:
 - Cache
 - Hauptspeicher
- $t_{\text{cache}} = 1$ Takt
- $t_{MM} = 100$ Takt
- **Wie lang ist die AMAT des Programmes aus Beispiel 1?**

Beispiel 2: Leistung eines Speichersystems



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Annahme: Prozessor hat zwei Hierarchieebenen:
 - Cache
 - Hauptspeicher
- $t_{\text{cache}} = 1$ Takt
- $t_{MM} = 100$ Takte
- Wie lang ist die AMAT des Programmes aus Beispiel 1?

$$\begin{aligned}\text{AMAT} &= t_{\text{cache}} + MR_{\text{cache}} t_{MM} \\ &= [1 + 0,375 * 100] \text{ Takte} \\ &= \mathbf{38,5 \text{ Takte}}\end{aligned}$$

Gene Amdahl, 1922 -

- **Amdahls Gesetz:** Die Optimierung eines Subsystems bringt nur dann etwas, wenn das Subsystem tatsächlich großen Einfluss auf die Gesamtrechenleistung hat
- Gründete drei Firmen, darunter auch die Amdahl Corporation in 1970
 - IBM-kompatible Großrechner
 - I.d.R. schneller und/oder billiger



Geheimes Lager

- **Höchste** Ebene der Speicherhierarchie
- **Schnell** (oft Zugriffszeit ~ 1 Takt)
- Stellt dem Prozessor **idealerweise** die meisten benötigten Daten zur Verfügung
- Speichert (i.d.R.) die **zuletzt** benutzten Daten

Aufbau von Caches: Entwurfsentscheidungen



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Welche Daten werden im Cache **gehalten**?
- Wie werden die Daten **gefunden**?
- Wie werden Daten **ersetzt**?

Schwerpunkt hier auf Loads, Stores werden aber ähnlich gehandhabt

Welche Daten werden im Cache gehalten?

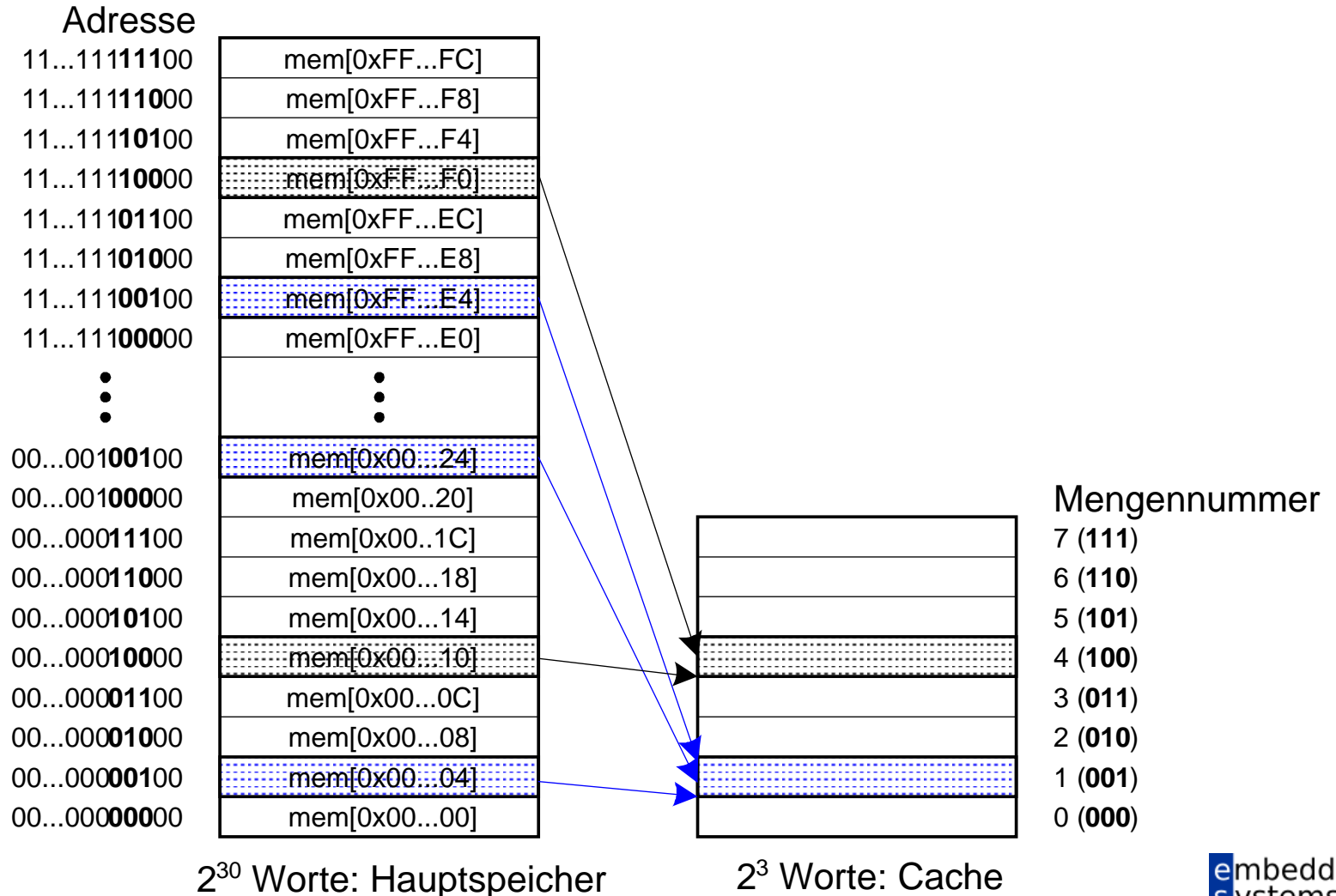
- Idealfall: Cache “**ahnt**” im Voraus, welche Daten der Prozessor benötigen wird und hält diese bereit
- Praxis: Prophezeiungen in der Regel **ungenau**
- Basiere Vorhersagen auf **bisherigem** Verhalten
- **Zeitliche Lokalität**: kopiere gerade **benutzte** Daten in den Cache. Bei nächster Verwendung werden die Daten im Cache gefunden (*Cache Hit*).
- **Räumliche Lokalität**: kopiere benachbarte Daten auch in den Cache
 - Blockgröße: Anzahl von Bytes, die immer **zusammen** in den Cache kopiert werden

- Kapazität (*capacity*, C):
 - Anzahl der im Cache speicherbaren Bytes
- Blockgröße (*block size*, b):
 - Anzahl der auf einen Satz in den Cache geladenen Bytes
- Blockanzahl ($B = C/b$):
 - Anzahl von Blöcken im Cache: $B = C/b$
- Assoziativitätsgrad (*degree of associativity*, N):
 - Anzahl von Blöcken in einer Assoziativitätsmenge (kurz: *Menge*)
- Anzahl von Assoziativitätsmengen ($S = B/N$):
 - Jede Speicheradresse wird auf genau eine Assoziativitätsmenge abgebildet

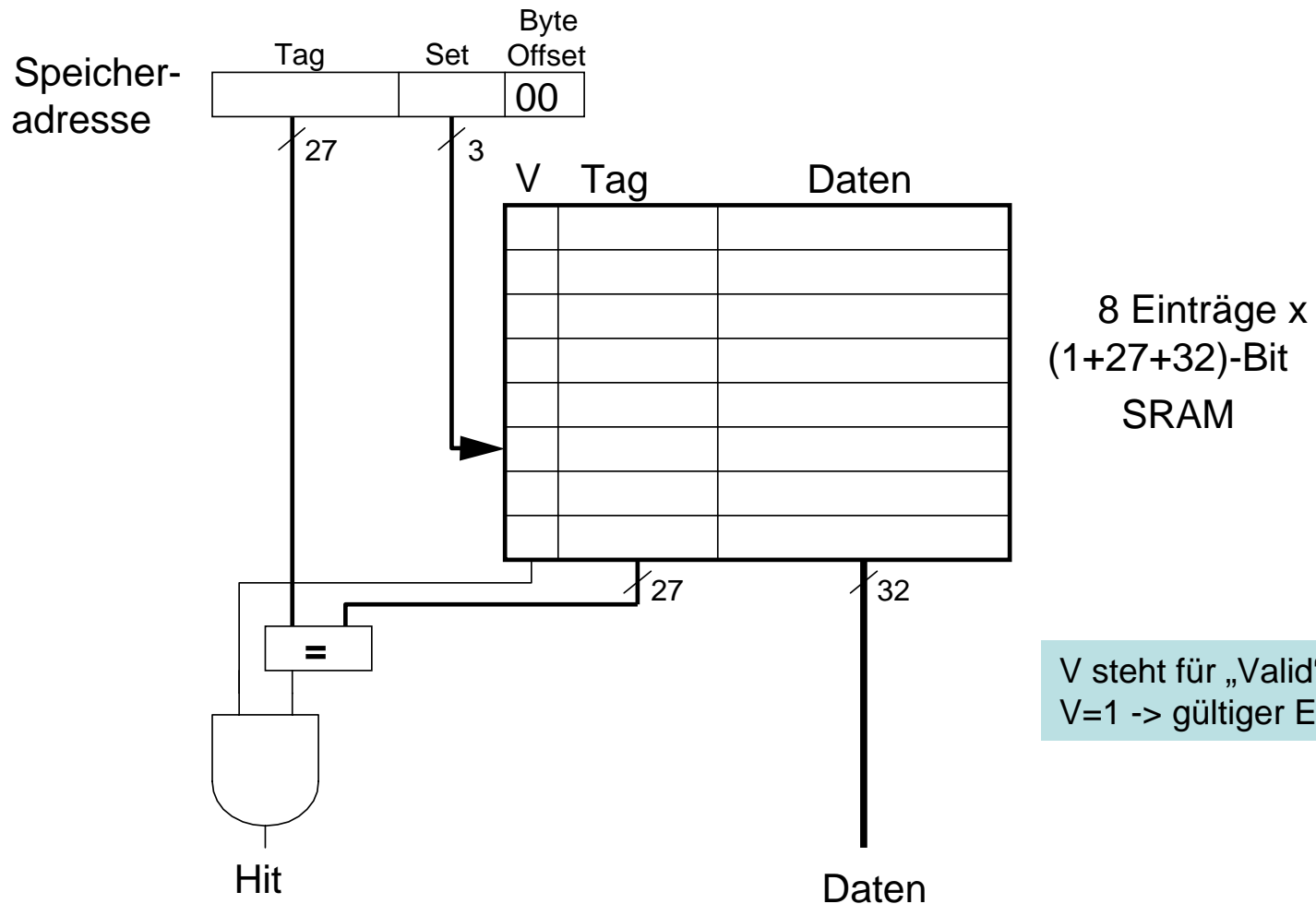
Wie werden Daten gefunden?

- Cache ist organisiert als S Mengen (*sets*)
- Jede Speicheradresse wird genau auf **eine** Menge abgebildet
- Caches werden klassifiziert nach **Anzahl von Blöcken** in einer Menge:
 - **Direkt abgebildet (*direct mapped*)**: Ein Block pro Menge
 - **N -fach Mengenassoziativ (*N-way set associative*)**: N Blöcke pro Menge
 - **Vollassoziativ (*fully associative*)**: Alle Cache Blocks in einer Menge
- Beispielorganisationen für einen Cache mit:
 - Kapazität $C = 8$ Worte
 - Blockgröße $b = 1$ Wort
 - Damit ist Blockanzahl $B = 8$

Direktabgebildeter Cache

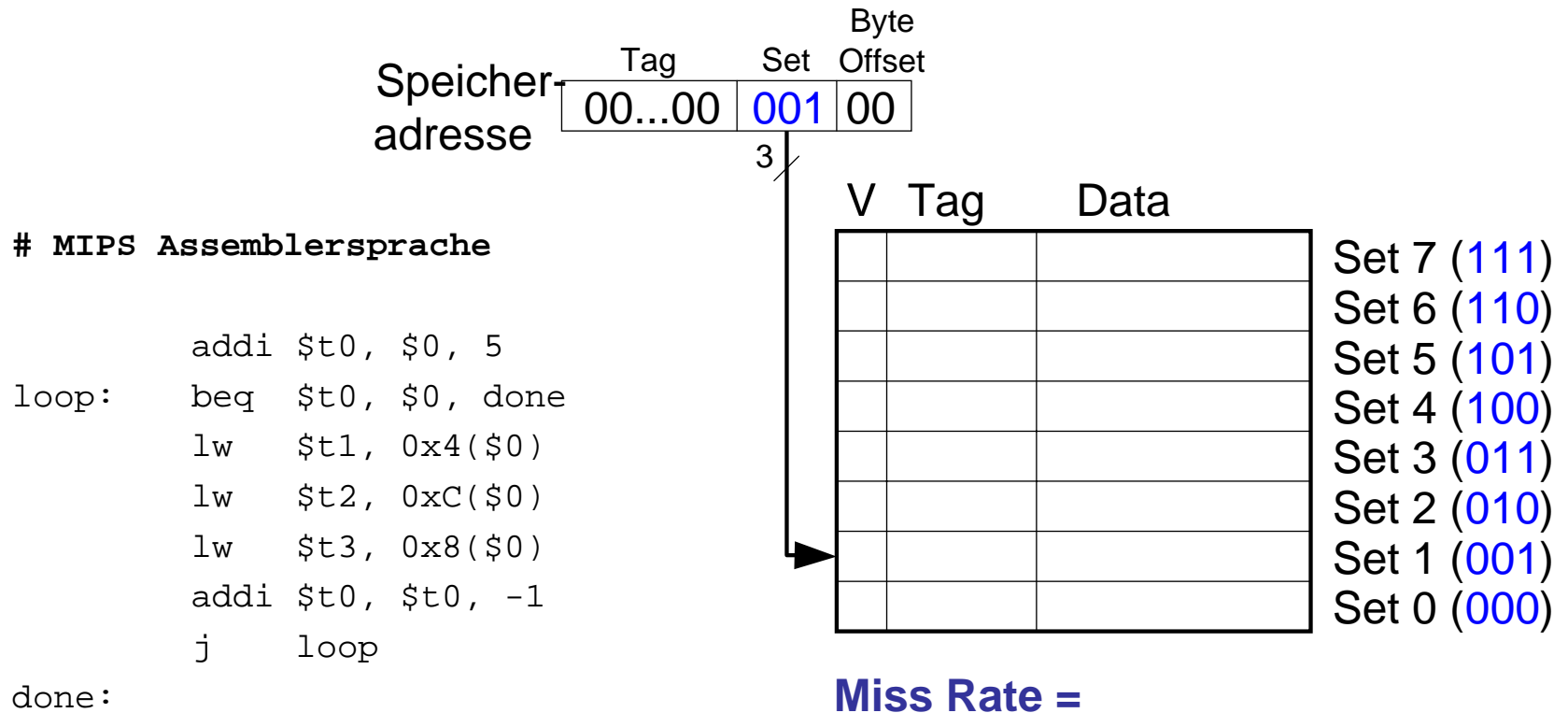


Direktabgebildeter Cache: Hardware

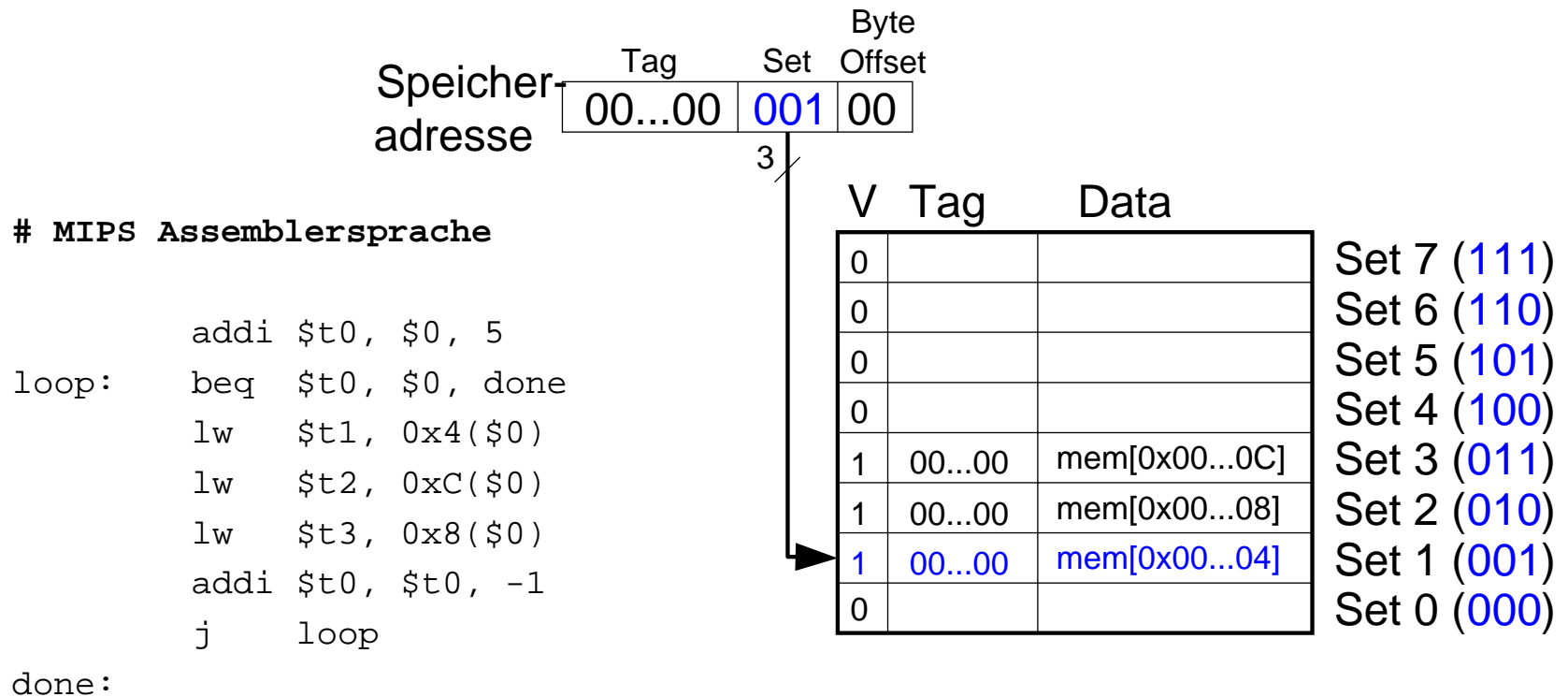


V steht für „Valid“
V=1 -> gültiger Eintrag

Leistung eines direkt abgebildeten Caches



Leistung eines direktabgebildeten Caches

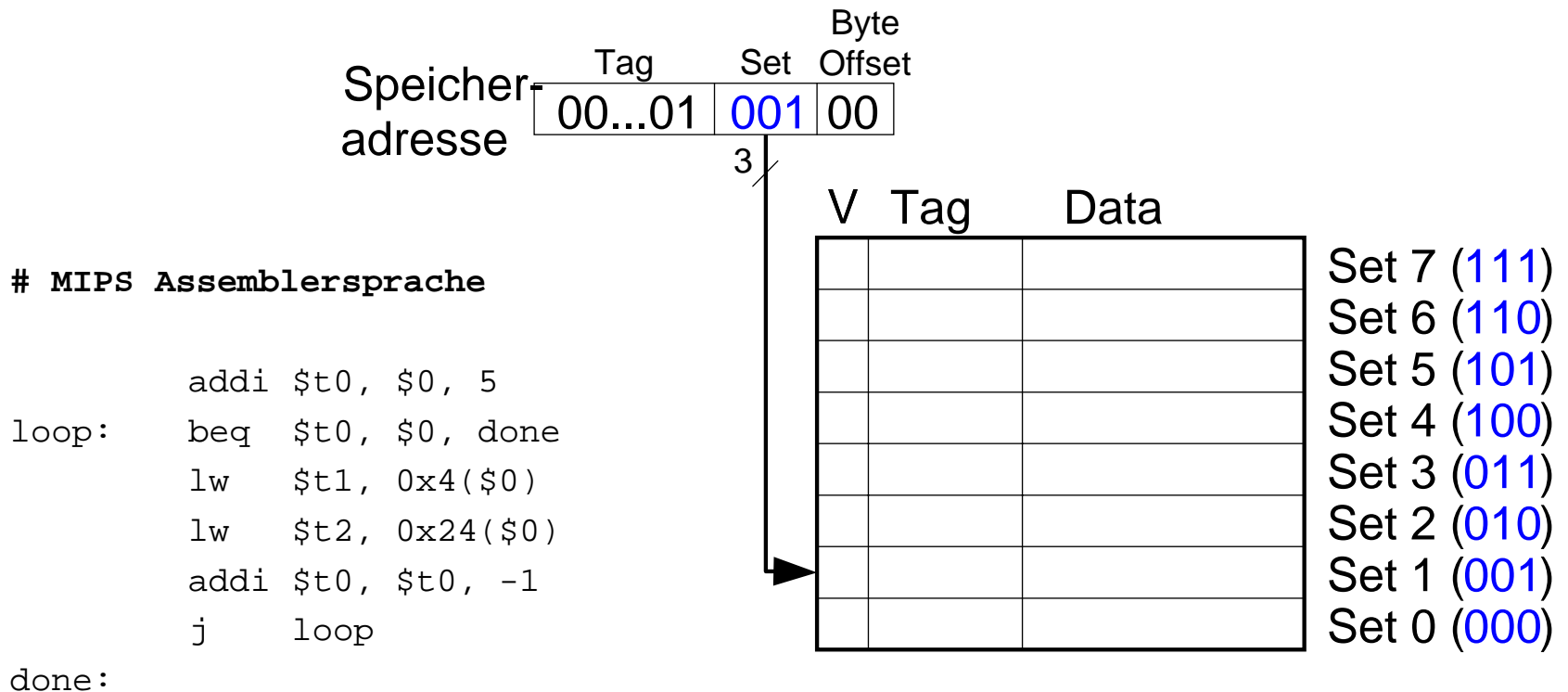


$$\text{Miss Rate} = 3 / 15 \\ = 20\%$$

Temporale Lokalität

Unvermeidbare (compulsory) Misses

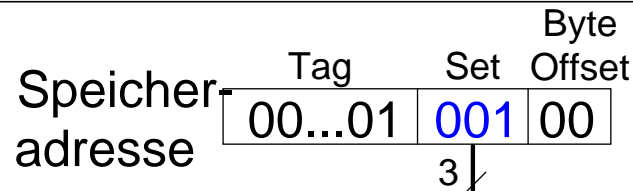
Konflikte bei direktabgebildeten Caches



Konflikte bei direktabgebildeten Caches

MIPS Assemblersprache

```
    addi $t0, $0, 5
loop: beq  $t0, $0, done
      lw   $t1, 0x4($0)
      lw   $t2, 0x24($0)
      addi $t0, $t0, -1
      j    loop
done:
```



V	Tag	Data	
0			Set 7 (111)
0			Set 6 (110)
0			Set 5 (101)
0			Set 4 (100)
0			Set 3 (011)
0			Set 2 (010)
1	00...00 00...01	mem[0x00...04] mem[0x00...24]	Set 1 (001)
0			Set 0 (000)

Miss Rate = 10/10
= 100%

Conflict Misses