

## Übersicht:

- 1 Zahlendarstellung
- 2 Numerik
- 3 Codierung
- 4 Informationstheorie
- ~~5 Boolesche Algebra (Mathematik)~~
- 6 KV und BDD
- 7 Digitalerschaltungen (Kombinatorische Logik)
- 8 Automaten
- 9 Sequentielle Logik (Speicher-Elemente)
- 10 Speicher
- 11 Schaltwerke
- 12 Schaltwerke: Realisierung
- 13 Schaltwerke: Moore vs. Mealy
- 14 Micro 16
- 15 Befehlssatz
- 16 Pipelining
- 17 Speichermanagement
- 18 Multicore
- 19 Speichermodelle

# 01 Zahlendarstellung

## Stellenwertsystem

Basis  $b \geq 2$  ( $b=1$  wäre unüblich)

Ziffern  $[0; b-1]$

## Zahlenumwandlung

Quellsystem  $\rightleftharpoons$  Zielsystem

Binär zu Dezimal

$$(110,1101)_2 = (\dots)_{10}$$

$$1 \cdot 2^2 + 1 \cdot 2^1 + 0 \cdot 2^0 + 1 \cdot 2^{-1} + 1 \cdot 2^{-2} + 0 \cdot 2^{-3} + 1 \cdot 2^{-4} = 6,8125$$

oder im Horner-Schema:

Vorkommateille:

$$((1) \cdot 2 + 1) \cdot 2 + 0 = 6$$

Nachkommateille:

$$(((1)/2 + 0)/2 + 1)/2 + 1)/2 = 0,8125$$

## Dezimal zu Binär

Vorkommateille:

$$a_0 = 6 \bmod 2 = 0$$

$$a_1 = \underbrace{(6-0)/2}_{=3} \bmod 2 = 1$$

$$a_2 = \underbrace{(3-1)/2}_{=1} \bmod 2 = 1$$

$$a_3 = (1-1)/2 \bmod 2 = 0 \rightarrow \text{Abbruch } (110)_2$$

effizienter:

$$a_0 = 6 \bmod 2 = 0$$

$$a_1 = \underbrace{\lfloor \frac{6}{2} \rfloor}_{=3} \bmod 2 = 1$$

$$a_2 = \lfloor \frac{3}{2} \rfloor \bmod 2 = 1$$

$$a_3 = \lfloor \frac{1}{2} \rfloor \bmod 2 = 0 \rightarrow \text{Abbruch } (110)_2$$

Nachkommastelle:

effiziente Verteilung:

$$a_{-1} = \lfloor 0,8125 \cdot 2 \rfloor = \lfloor 1,625 \rfloor = 1$$

$$a_{-2} = \lfloor 0,625 \cdot 2 \rfloor = \lfloor 1,25 \rfloor = 1$$

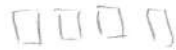
$$a_{-3} = \lfloor 0,25 \cdot 2 \rfloor = \lfloor 0,5 \rfloor = 0$$

$$a_{-4} = \lfloor 0,5 \cdot 2 \rfloor = \lfloor 1 \rfloor = 1$$

$$a_{-5} = \lfloor 0 \cdot 2 \rfloor = \lfloor 0 \rfloor = 0 \rightarrow \text{Abbruch } (0,1101)_2$$

Binär - Hexadezimal - Beziehung

$$2^4 = 16$$



4er Blöcke zur Umwandlung

# Darstellung negativer Zahlen durch Kodierung

## V2-B : Vorzeichen und Betrag

Erstes Bit

- 0 positiv
- 1 negativ

Probleme:

$\pm 0$   
 $+0, 1, 2, 3, -0, -1, -2, -3$

negative Zahlen stehen nach pos

## Einer-Komplement

Bitwise Inversion

Probleme:

$\pm 0$   
 $+0, 1, 2, 3, -3, -2, -1, -0$

Durch doppelte 0 entstehen leicht Überläufe,  
mit +1 korrigieren.

## Zweier-Komplement

Bitwise Inversion + 1 (wenn negativ)

Dadurch eindeutige 0:

$0, 1, 2, 3, -4, -3, -2, -1$

$$\begin{array}{r} 011 \quad (3) \\ 111 \quad (-0) \\ \hline \times 010 \quad (2) \quad \xrightarrow{+1} \quad 011 \quad (3) \end{array}$$

## Exzess

„symmetrischer Exzess“

Kodierung: + kleinste darstellbare Zahl

$$e = 2^{\text{bits}-1} \quad \text{oder} \quad e = 2^{\text{bits}-1} - 1$$

„Asymmetrischer Exzess“

Kodierung: + beliebige Zahl

Rechnen mit Exzess:

$$e = 12$$

$$A = -3$$

$$B = 4$$

$$A_e = A + e = 9$$

$$B_e = B + e = 16$$

$$A_e + B_e = A + B + 2e$$

$$(A+B)_e = A_e + B_e - e$$

## O2: Numerik

Kodierung von Zahlen mit Nachkommastelle

### Festpunkt

$$\text{Gesamtlänge} = \underbrace{1}_{VZ} + \underbrace{g}_{\text{vorkomma}} + \underbrace{n}_{\text{nachkomma}}$$

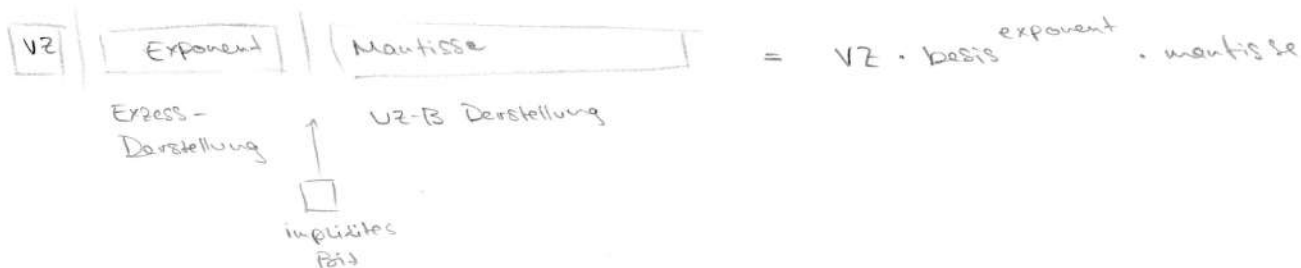
→ Nur  $\in \mathbb{Q}$  darstellbar

Schrittweite / kleinste darstellbare Zahl

$$\Delta x = 2^{-n}$$

### Gleitpunkt

Idee: je größer die Zahl, desto unwichtiger  $\Delta x$ .



### Normalisierte Form

implizites Bit = 1

$$\Delta x = x_{\min} \cdot b^{e_{\min}} = 1 \cdot 2^{e_{\min}}$$

zB:  $b=2$   $e_{\min}=-1$



### Denormalisierte Form

implizites Bit = 0

"Subnormale Zahlen" lassen sich darstellen.

Problem:

- unerwünschte Lücke
- Null nicht darstellbar, muss Kodiert werden

### "Gleitpunkt-Darstellungs-Raum" $\mathbb{F}$

$$\mathbb{F}(b, p, e_{\min}, e_{\max}, denorm)$$

$b$  ... base  $\geq 2$

$p$  ... precision (Mantisse mit implizitem Bit)

$e_{\min}/e_{\max}$  ... kleinster / größter Exponent

denorm ... boolean

### unit of last position (ulp)

$$ulp = (0,000 \dots 0001)_b = b^{-p+1} \leftarrow \text{wegen vorkommastelle}$$

$$\Delta x = ulp \cdot b^{e_{\min}} = b^{e_{\min}-p+1}$$

# IEEE 754

Eigene Parameter für IF.

Single : 32 Bit , 8 Exponent Bits , 24 Mantissenbits inkl impliz.

double : 64 Bit , 11 Exponent Bits , 53 Mantissenbits inkl impliz.

## Kodierungen im Exponenten

VZ	Exponent	Mantisse	Denormalisiert	$\rightarrow e = e_{min}$ obwohl da $e_{min}-1$ steht
-	0...000	-	Denormalisiert	
0	0...000	0	+0	
1	0...000	0	-0	
0	1...111	>0	NaN	
0	1...111	0	+∞	
1	1...111	0	-∞	

Eigentlicher dekodierter Exponent-Wert: Exzessdarstellung.

Beispiel: IF (2, 11, -14, 15, true)  $\rightarrow$  16 Bit - 10 expl. Bit - 1 VZ Bit = 5 Bit

base ... 2  
 inkl. impl.  $\leftarrow$  precision ... 11  
 $e_{min} : -14$   
 $e_{max} : 15$   
 denorm : true

$$Exzess = 2^{(expon. Bits) - 1} - 1$$

$$Exzess = 2^{5-1} - 1 = 15 = (01111)_2$$

Dadurch:

$$e_{min} - 1 = -15 \quad (00000)$$

$$e_{max} + 1 = 16 \quad (11111)$$

## Runden

außerhalb von arithmetischen Operationen

### Abschneiden / truncate

$$\underbrace{(-1,62\overline{6})}_{n=2}_{10} \longrightarrow \square x = (-1,62)_{10}$$

$$\underbrace{(1,10\overline{1})}_{n=2}_2 \longrightarrow \square x = (1,10)_2$$

### Gerichtetes Runden / directed rounding

$$\underbrace{(-1,62\overline{6})}_{n=2}_{10}$$

Aufrunden  $\square x = \max(-1,62, -1,63) = -1,62$

Abzurunden  $\square x = \min(-1,62, -1,63) = -1,63$

$$\underbrace{(1,10\overline{1})}_{n=2}_2$$

Aufrunden  $\square x = \max(1,10, 1,11) = 1,11$

Abzurunden  $\square x = \min(1,10, 1,11) = 1,10$

### Optimale Rundung / round to nearest

$$\hat{x} = \frac{x_1 + x_2}{2} \quad \text{bzw.} \quad \hat{x} = (y_1 + y_2) \cdot 0,1$$

$$\underbrace{(-1,62\overline{6})}_{n=2}_{10} \longrightarrow \frac{(-1,62) + (-1,63)}{2} = -1,625 \neq \hat{x}$$

$\square x = -1,63$

$$\underbrace{(1,10\overline{1})}_{n=2}_2 \longrightarrow \frac{1,10 + 1,11}{10} = 1,101 = x$$

round to even

$$\square x = 1,10 \rightarrow \text{es soll } 0 \text{ stehen}$$

round away from zero

$$\square x = 1,11$$

# Runden mit GRS

guard digit

round digit

sticky bit

(bleibt 1 wenn beim Shiften eine 1 durch geht)

Bei Addition & Subtraktion als sticky Bit rechnen.

## Optimale Rundung:

G R S

0 X X

1 X X

1 0 0

—  
+1

→ Round to even  
wenn LSB von Mantisse

= 0 —

= 1 +1

Round away from zero

+1



# Beispiel: Addition in IEEE 754 (16 Bit)

$$A = (5,58)_{10}$$

$$B = (62,27)_{10}$$

## Umrechnung ins Gleitpunkt-Format

1 Bit VZ

5 Bit Exponent

10 Bit Mantisse explizit

$$\rightarrow \text{excess: } 2^{5-1} - 1 = 15 = (01111)_2$$

### Normalisierung

$$A = (5,58)_{10} = (101,1001010)_2 \cdot 2^0 = (1,011001010)_2 \cdot 2^2$$

$$B = (62,27)_{10} = (11111,0100)_2 \cdot 2^0 = (1,11110100)_2 \cdot 2^5$$

### Exponenten nach Normalisierung kodieren

$$A: 2 = (10)_2 \rightarrow 10 + 01111 = 10001$$

$$B: 5 = (101)_2 \rightarrow 101 + 01111 = 10100$$

### Exponenten angleichen

$$A = m \cdot 2^2$$

$$B = m \cdot 2^5$$

Da  $A < B$  muss A angeglichen werden:

$$10001 + 11 = 10100$$

$$\begin{array}{l} \text{Vorher: } 0 | 10001 | 1 | 0110010100 \\ \text{Nachher: } 0 | 10100 | 0 | 0010110010 \quad 100 \quad \rightarrow \text{rsh}(3) \\ \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \text{GRS} \\ \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \uparrow \\ \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \text{implizit} \end{array}$$

### Addieren

$$\begin{array}{r} A: 0 | 10100 | 010010110010 \quad \text{GRS} \\ B: 0 | 10100 | 11111001000 \quad 100 \\ \hline 0 | 10100 | 1010001111010 \quad 100 \end{array}$$

Normalisieren: rsh(1), Exp+1

$$0 | 10101 | 1100001111010 \quad \text{GRS}$$

### Runden

0xx  $\rightarrow$  unverändert

## Fortsetzung: Subtraktion

- ✓ Umrechnung ins Gleitpunkt-Format
- ✓ Normalisierung, Kodierung des Exponenten
- ✓ Angleichung der Exponenten

$$A = (9,58)_{10}$$

$$B = (62,27)_{10}$$

### Subtraktion

										GRS
B:	0	10100	1	111100	1000					000
- A:	0	10100	0	001011	0010					100
	①	10100	1	11000	10101					100

↑  
Ergebnis bereits vorundiziert

↓

$$A - B = -(B - A)$$

vereinfacht Rechnen

### Runden

...101 100      LSB = 0 -

└──────────┘ → LSB = 1 +1

### Lösung:

$$1 \ 10100 \ 1 \ 11000 \ 10110$$

## Vorgang beim Rechnen mit IEEE Darstellung

- 1) Angleichung der Exponenten
- 2) Mantissen addieren/subtrahieren
- 3) Ergebnis normalisieren
- 4) Runden (Guard Digit, Round Digit, Sticky Bit)

Frage 2  
 Falsch  
 Erreichte Punkte 0,00 von 1,00  
 Frage markieren

Addieren Sie die Zahlen 0100101010111000 und 1011101000111001.  
 Es gilt das Gleitpunktformat  $F(2, 11, -14, 15, true)$  IEEE 754-2008 mit half-precision (16 Bit — wie in der Vorlesung).  
 Verwenden Sie für das Ergebnis ebenfalls dieses Format.  
 Als Rundungsvorschrift verwenden Sie bitte "round to nearest — round to even".

Antwort: 0011100010101001 ✘

Die richtige Antwort ist: 0100101001010100

**Angabe:**

Addieren Sie die Zahlen 0100101010111000 und 1011101000111001  
 Als Rundungsvorschrift verwenden Sie bitte "round to nearest — round to even".

A: 0 10010 1010111000  
 e = 10010 - 01111 = 0011

B: 1 01110 1000111001  
 e = 01110 - 01111 = -0001

Unterschied zwischen Exponenten: 3 und -1: 4

**Exponent von B um 4 erhöhen, Bits um 4 nach rechts shiften:**

A: 0 10010 1010111000 | 000  
 B: 1 10010 0001100011 | 100 < implizites Bit zeigt sich durch shift  
 eigentlich müsste B kodiert werden weil es jetzt denormiert ist  
 aber das ist nur ein Zwischenschritt zum addieren

**Mantissen addieren:**

Da B negativ ist wird B von A abgezogen

A: 1.1010111000 | 000  
 B: 0.0001100011 | 100 < normalerweise ab „round digit“ rechnen aber bei binär auch ab „stick bit“  
 -----  
 L: 1.1001010100 | 100

**Normalisieren:**

Es ist bereits normalisiert

**Runden:**

1.1001010100 | 100 < weitere Rundungsregel  
 "round to nearest — round to even"  
 lsb = 0 → deshalb Mantisse unändert

L: 1.1001010110

**Meine Lösung:**

0 10010 1001010100

**Korrekte Lösung:**

0 10010 1001010100

Frage 3

Falsch

Erreichte Punkte 0,00 von 1,00

Frage markieren

Subtrahieren Sie die Zahl 1101000100001101 von der Zahl 1111010000010010.

Es gilt das Gleitpunktformat  $\mathbb{R}(2, 11, -14, 15, true)$  IEEE 754-2008 mit half-precision (16 Bit — wie in der Vorlesung).

Verwenden Sie für das Ergebnis ebenfalls dieses Format.

Als Rundungsvorschrift verwenden Sie bitte "round to nearest — round to even".

Antwort: 111100000100011

x

Die richtige Antwort ist: 1111010000001111

### Angabe:

Subtrahieren Sie die Zahl 1101000100001101 von der Zahl 1111010000010010

A und B sind negativ

A endet mit 0

B endet mit 1

A - B --> wenn beide negativ sind:  $B - A = - (A - B)$

A: 1 11101 0000010010  
e = 11101 - 01111 = 1110

B: 1 10100 0100001101  
e = 10100 - 01111 = 0101

Unterschied: 14 - 5 = 9

Exponent von B um 9 erhöhen, Bits um 9 nach rechts shiften:

A: 1 11101 0000010010 | 000

B: 1 11101 0000000010 | 101

Mantissen addieren:

A: 1.0000010010 | 000

B: 0.0000000010 | 101

-----  
L: 1.0000001111 | 011 (muss negativ sein)

Normalisieren:

ist bereits normalisiert

Runden:

Bei 011 --> unverändert

Meine Lösung:

1 11101 0000001111

Korrekte Lösung:

1 11101 0000001111

Multiplizieren Sie die Zahlen 1001100011111011 und 1100100000010010.

Es gilt das Gleitpunktformat  $\mathbb{F}(2, 11, -14, 15, true)$  IEEE 754-2008 mit half-precision (16 Bit — wie in der Vorlesung).

Verwenden Sie für das Ergebnis ebenfalls dieses Format.

Als Rundungsvorschrift verwenden Sie bitte "round to nearest — round to even".

Antwort:



Die richtige Antwort ist: 0010010100010001

**Angabe:**

Multiplizieren Sie die Zahlen 1001100011111011 und 1100100000010010

A: 1 00110 0011111011  
e = 00110 - 01111 = -1001  
e = 6 - 15 = -9

B: 1 10010 0000010010  
e = 10010 - 01111 = 0011  
e = 18 - 15 = 3

Neuer Exponent (kodiert): 00110 + 10010 - 01111 = 01001  
6 + 18 - 15 = 9  
Eigentlicher Wert: 9 - 15 = -9 + 3 = -6 (korrekt kodiert)

**Multiplizieren:**  
1.0011111011 \*  
1.0000010010 =  
1.0100010001 | 011

**Runden:**  
011 --> unverändert

✓ **Meine Lösung:**  
0 01001 0100010001  
**Korrekte Lösung:**  
0 01001 0100010001

**Angabe:**

Gegeben ist die Zahl  $(10.011000)_2$ .

Runden Sie die Zahl mittels "round to nearest - round away from zero" auf 3 Nachkommastellen.

Antwort:



Die richtige Antwort ist: 10.011

$$x = 10.011000$$

$$x_1 = 10.011$$

$$x_2 = 10.100$$

$$x' = (x_1 + x_2) * 0.1 = 100.111 * 0.1 = 10.0111 \rightarrow \text{ungleich mit } x$$

Lösung: 10.011

**Anlauf 2:**

**Wichtig: Punkt statt Komma benutzen damit es richtig compiled wird**

Gegeben ist das folgende Bitmuster: 001011111

Interpretieren Sie dieses als eine Festpunktzahl mit 4 Nachkommastellen und geben sie den entsprechenden dezimalen Wert an.

0 0101.1111  
5.9375

Gegeben ist die Zahl  $(10.101101)_2$ .

Runden Sie die Zahl mittels "Abrunden (round min)" auf 3 Nachkommastellen.

$x = 10.101$  (unverändert)

Gegeben ist die Zahl  $(11.100011)_2$ .

Runden Sie die Zahl mittels "round to nearest - round away from zero" auf 3 Nachkommastellen.

$x_1 = 11.100$

$x_2 = 11.101$

$x' = (11.100 + 11.101) * 0.1 = 11.1001$

Lösung : 11.100



**Angabe:**

Multiplizieren Sie die Zahlen 1101010011111000 und 1000110010000010.

Als Rundungsvorschrift verwenden Sie bitte "round to nearest - round to even".

A: 1 10101 0011111000

B: 1 00011 0010000010

**Neuer Exponent (kodiert):**

$$10101 + 00011 - 01111 = 01001$$

$$A_e + B_e = A + B - e$$

**Mantissen multiplizieren:**

$$1.0011111000 *$$

$$1.0010000010 =$$

$$1.0110011001 \mid 011$$

**Runden:**

unverändert

**Meine Lösung:**

0 010010110011001

**Angabe:**

Subtrahieren Sie die Zahl 1011010110001111 von der Zahl 0011001010110100.

Als Rundungsvorschrift verwenden Sie bitte "round to nearest - round to even".

$$A \text{ (pos)} - B \text{ (neg)} = A + B \text{ (pos)}$$

A: 0 01100 1010110100

B: 1 01101 0110001111

**Exponent von A angleichen (+1), Bit shift nach rechts**

Implizite 1 wird sichtbar

A: 0 01100 1101011010 | 000 -->

**Addieren:**

0.1101011010

1.0110001111

-----

10.0011101001

**Normalisieren:**

Exponent +1 und shift nach rechts

01100 + 1 = 01101

10.0011101001 → 1.0001110100 | 100

**Runden:**

Weitere Rundungsregel -> lsb = 0  
unverändert lassen

✓ **Meine Lösung:**

0 01101 0001110100

**Richtige Lösung:**

0 01110 0001110100

**Angabe:**

Addieren Sie die Zahlen  $1\ 01000\ 1000110111$  und  $0\ 11100\ 0101101001$

A (neg)

B (pos)

A:  $1\ 01000\ 1000110111$

B:  $0\ 11100\ 0101101001$

**Exponenten ausgleichen:**

A:  $01000$

B:  $11100$

Exp von A ist um 20 kleiner als B

A:  $1\ 01000\ 1000110111$

A:  $1\ 11100\ 0000000000\ | \ 000$

**Addieren:**

A = 0

Deshalb bleibt b gleich und ist die Lösung und muss nicht gerundet werden

✓ **Meine Lösung:**

B:  $0111000101101001$

# 03: Codierungen

Beispiele für Zeichenkodierungen:

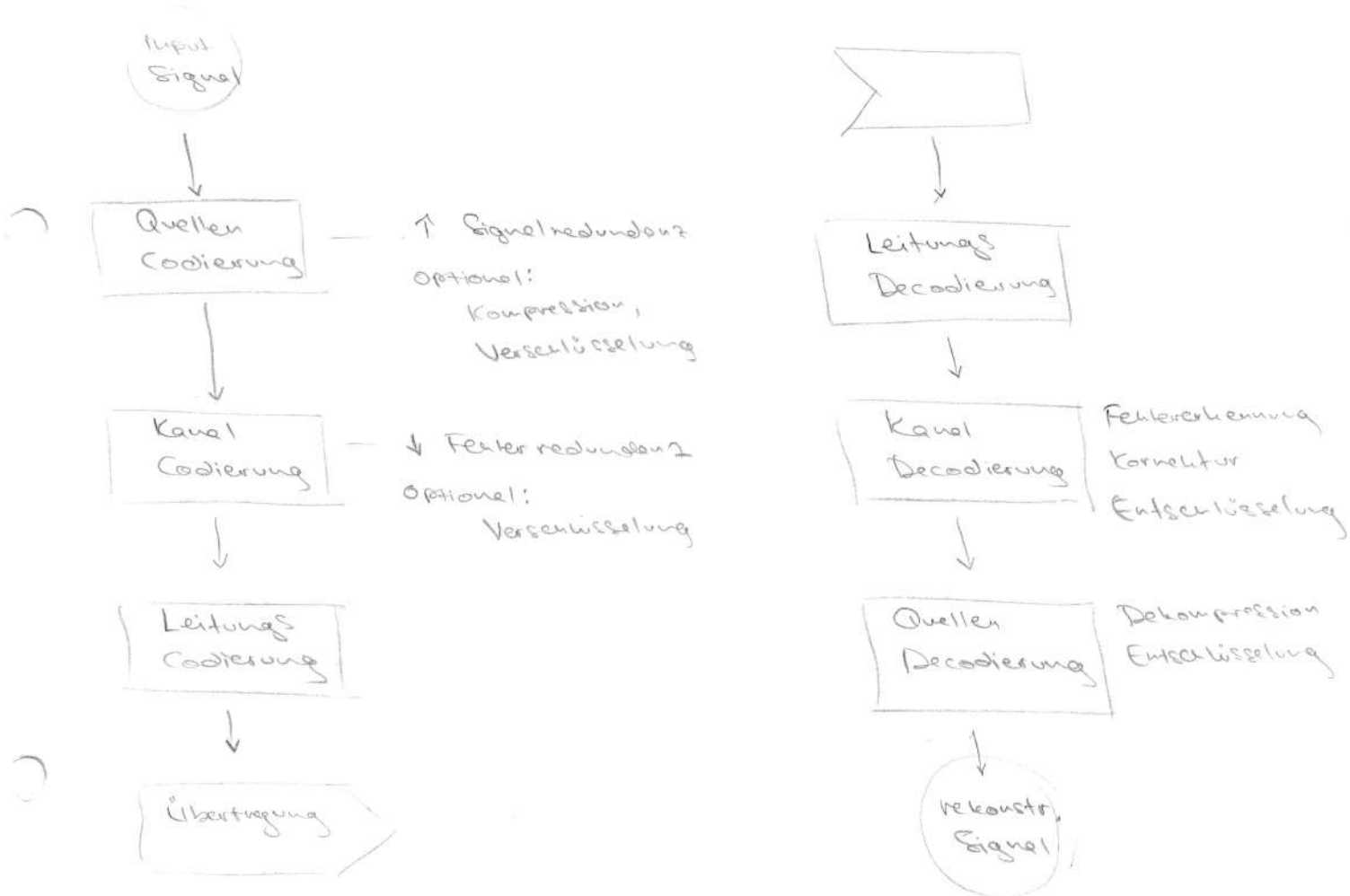
- ASCII
- Unicode
- UTF
- ISO-Latin-1

In Alltag:

- QR-Codes
- Bar-Codes - (ISO 15420)

## Codierungstheorie

Die mehrstufige Codierung: Trade-off  
Sicherheit / Redundanz = Aufwand



In der Quellencodierung

Signal  $\rightarrow$  binärer Code, komprimiert, verschlüsselt

In der Kanalcodierung

binärer Code  $\rightarrow$  Uncodierung, senkung redundanter Bits, verschl.

In der Leitungscodierung

Optimierung von binärem Code für Übertragung.

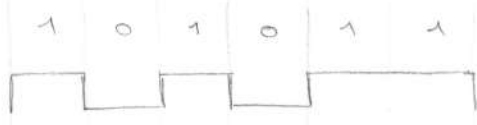
ZB Licht für Glasfaser, Strom, etc

Pegelfolgen = Bitfolgen

# Leitungs - Code NRZ: Non return to zero

benötigt separaten Taktsignal

Man unterscheidet zwischen high und low



NRZ-L low = 0 high = 1



NRZ-M Pegelwechsel = 1



NRZ-S Pegelwechsel = 0



RZ: unipolarer Return-to-Zero-Code

Zweite Takthälfte immer "low"



D-Manchester:

Es muss in jedem Bit-Takt eine Flanke geben

- wenn man nach der Flanke auf derselben Seite ist wie zuvor: 0
- wenn man Seite wechselt: 1

## Arten von Binär codes:

### Codes mit variabler Länge

zB UTF-Code, Huffman Code (Code-Baum) / Morse-Code

### Codes mit fixer Länge (Blockcodes)

#### Nicht linear

zB EAN-13, US-ASCII früher

#### Linear

Durch Summe von Codewörtern bleibt es gültig (Modulo)

zB Hamming-Code

#### Zyklische Codes

Durch schieben / rotieren bleibt Code gültig

zB Polynom-Code (CRC-Codes)

Blockcodes können system. Codes sein

### Systematische Codes (Prüfstellen Codes)

Prüfstellen / redundante Bits

zB. Paritätsbits (auch mehrdimensional)


Paritätsbits für Anteile (Hamming-Code)

Prüfsummen (EAN, Polynom-Codes)

# EAN-13

fixe Länge, nicht linear

13 = 12 Stellen + Prüfziffer

$$\textcircled{z_1} + 3z_2 + z_3 + 3z_4 + z_5 + 3z_6 + z_7 + 3z_8 + z_9 + 3z_{10} + z_{11} + 3z_{12} + \text{Prüfziffer} \equiv$$


Implizit:  
Aus A, B Muster  
(AABA...)

$$\equiv 0 \pmod{10}$$

## Hamming-Distanz

Maß für die Störanfälligkeit von Code

(Bestimmung mit XOR möglich)

zwischen 2 Wörtern:

$$u(00, 01) = 1 = D$$

zwischen mehreren Wörtern:

minimale Hamming Distanz

Bei Abstand  $D$  gilt:

$D-1$  Fehler erkennbar

$k < \frac{D}{2}$  Fehler korrigierbar

## Paritätsbits

$\boxed{m} \quad \boxed{P}$

even-parity = # Einsen gerade bleiben

odd-parity = # Einsen ungerade bleiben

## Hamming-Code

linear, systematisch (fehlerkorrigierend), fixe Länge

Abstand  $D$  immer = 3

1 Bit korrigierbar

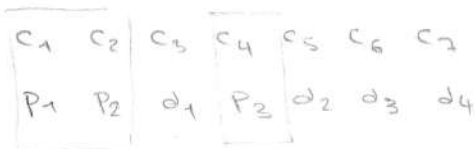
2 Bit erkennbar

Code indizes ab 1 zählen, an jeder 2er Potenz steht ein Prüfbit



Da zB.  $2^7 = 16 + 8 + 4 + 1$  wird  $c_7$  die Prüfbits  $c_{16}, c_8, c_4$  und  $c_1$  beeinflussen.  
 $\begin{matrix} c_{16} & c_8 & c_4 & \text{und} & c_1 \\ \text{"} & \text{"} & \text{"} & & \text{"} \\ P_5 & P_4 & P_3 & & P_1 \end{matrix}$

## Beispiel



$$P_1 = c_1 = (c_3 + c_5 + c_7) \bmod 2 = c_3 \otimes c_5 \otimes c_7$$

$$P_2 = c_2 = (c_3 + c_6 + c_7) \bmod 2 = c_3 \otimes c_6 \otimes c_7$$

$$P_3 = c_4 = (c_5 + c_6 + c_7) \bmod 2 = c_5 \otimes c_6 \otimes c_7$$

↑  
XOR



Fortsetzung des Beispiels:

Sender: Datenwort 0110

Prüfbits

$$P_1 = 0 \oplus 1 \oplus 0 = 1$$

$$P_2 = 0 \oplus 1 \oplus 0 = 1$$

$$P_3 = 1 \oplus 1 \oplus 0 = 0$$

Zu übertragen:

$$\boxed{1 \ 1 \ 0} \boxed{0} \ 1 \ 1 \ 0$$

Empfänger: Störung (1 1) 0 0 0 1 0

$$P_1 = 0 \neq 1$$

$$P_2 = 1 = 1$$

$$P_3 = 1 \neq 0$$

}  $1+1=2 \rightarrow C_5$  wurde gestört  
„Indikatoren“



Bei Empfang Prüfbits selbst neu berechnen und vergleichen

Problem:

Fehlerbündel

# Polynom - Codes

Sender und Empfänger kennen Generatorpolynom  $G(x)$

$$G(x) = x^3 + x^2 + 1 \quad \begin{matrix} (1 & 1 & 0 & 1) \\ & 1 & 1 & 0 \end{matrix}$$

## Sender

Nachrichte Wort: 0 1 1 0 (4 Bit)

Message - Polynom:  $M(x) = 0x^3 + 1x^2 + 1x^1 + 0x^0 = x^2 + x$  ↵

1.  $x^r \cdot M(x)$  wobei  $r = \text{Grad von } G(x)$

$$x^3 \cdot M(x) = x^3(x^2 + x) = x^5 + x^4 \quad 0110000$$

2.  $\frac{x^r \cdot M(x)}{G(x)}$  und  $R(x)$  bestimmen

$$\begin{array}{r} (x^5 + x^4) : (x^3 + x^2 + 1) = x^2 \\ \underline{x^5 + x^4 + x^2} \\ x^2 \text{ Rest} \end{array} \rightarrow R(x) = x^2 + 100$$

3.  $T(x) = x^r \cdot M(x) - R(x)$

Transmission = Übertragung

$$T(x) = (x^5 + x^4) - x^2 = x^5 + x^4 + x^2$$

↑  
Rest abgezogen,  
dadurch Restfrei!

$$\begin{array}{r} 0110 \quad 100 \\ \hline M(x) \quad 1 \end{array} \rightarrow \text{Absenden}$$

(Man kann auch addieren)

Sicherheitsmaß

## Empfänger

0110100 wird fehlerhaft übertragen

1.  $T(x) + E(x)$  Error - Polynom empfangen

2.  $\frac{T(x) + E(x)}{G(x)}$  decodieren

Da  $T(x)$  Restfrei war, zeigt sich ein Rest

# CRC - cyclic redundancy check codes

zyklische Polynom-Codes mit Standard  $G(x)$

ZB: CRC-12 Generator - Polynom 12. Grades

Dadurch 12 Absicherungsbits für 6 Bit  
Nachrichten

Insgesamt:  $6 + 12 = 18$  Bit zu übertragen

12 Absicherungsbits

11 Fehler erkennen

< 6 Fehler korrigieren

$$\text{Fehler} < \frac{\text{Distanz}}{2}$$

$N = \text{message } m + \text{redundancy } r$

Es sind  $2^r$  mögliche redundante Codierungen  
möglich

- 1 Muster für Fehlerfreiheit

-  $2^r - 1$  weitere für weitere Informationen  
über Fehlerquelle

# 04: Informationstheorie

Nach dem Übertragungsmodell:



Das Medium überträgt den Signal.

Analog: Wert - kontinuierlich

Diskret: Wert - diskret

## Sprache

Alphabet: Menge aller Zeichen  $\rightarrow$  Wort

Syntax: Grammatik

Semantik: Bedeutung

Meta-Sprache: zur Definition einer neuen Sprache

## Code und Codierung

Codierung: bijektive Abbildung von 2 Wörtern zwischen 2 Sprachen

Wort  $\mapsto$  Codewort

$\uparrow$  Verarbeitung, Störzicherheit

z.B. Binärcode mit "binary Digits" = Bits

## Informationstheorie nach Claude Shannon

8 Bit = 1 Byte

$2^{10}$  Bit = 1 Kibit (Kilo)

⋮

Informationsgehalt / Überraschungswert

$P$  = Auftrittswahrscheinlichkeit / Erwartungswert

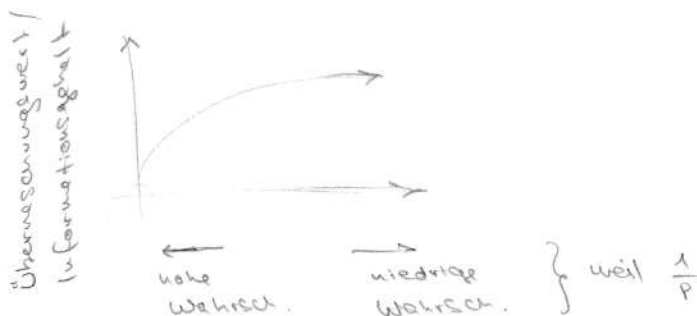
Es muss gelten:

$\uparrow$  Wahrscheinlichkeit =  $\downarrow$  Überraschung / Informationsgehalt  $f\left(\frac{1}{P}\right)$

weil Zeichen voneinander unabhängig  $\sum$  Einzelzeichen = gesamt. Wort

Also:

$h(x) + h(y) = h(xy) \rightarrow$  Erfüllt vom Logarithmus



## Informationsgehalt $h$

$$h = \log\left(\frac{1}{p}\right) = \log(p^{-1}) = -\log(p) \quad [\text{Bit}]$$

↑  
Wahrscheinlichkeit.

## Mittlerer inf. Gehalt $H$

$$H = \sum_i p_i h_i = \sum_i p_i \log\left(\frac{1}{p_i}\right) = - \sum_i p_i \log(p_i) \quad [\text{Bit}]$$

## Code-Wortlänge $L$

ZB	$x \mapsto 1$	$l_x = 1$
	$y \mapsto 01$	$l_y = 2$
	$z \mapsto 00$	$l_z = 2$

## Mittlere Wortlänge $L$

$$L = \sum_i p_i l_i \quad [\text{Bit}]$$

## Redundanz $R$

$$R = L - H \quad [\text{Bit}]$$

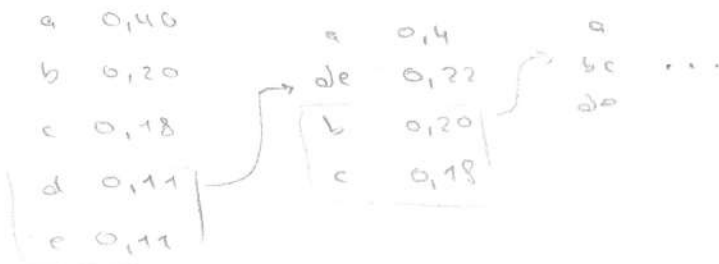
## relative Redundanz $r$

$$r = \frac{R}{L} \quad [\text{Einheitslos}]$$

## Huffman-Code

↓ Redundanz durch Datenverdichtung (Auch Morse Code)

Wenn man ganze Zeichenfolgen kodiert → beliebig kleine Redundanz



## Adaptiver Huffman-Code:

relative Häufigkeit in bisheriger Interaktion.

Nach  $n$  Zeichen Baum neu aufbauen

OG: KV-Diagramme, Binary decision diagrams (BDD)

KNF: Konjunktive Normalform

DNF: Disjunktive Normalform

für bool'sche Funktionen

A	B	C	Ergebnis	Klausel
0	0	0	0	$A \vee B \vee C$
0	0	1	0	$A \vee B \vee \neg C$
0	1	0	1	$\neg A \wedge B \wedge \neg C$
0	1	1	1	$\neg A \wedge B \wedge C$
1	0	0	0	$\neg A \vee B \vee C$
1	0	1	1	$A \wedge \neg B \wedge C$
1	1	0	0	$\neg A \vee \neg B \vee C$
1	1	1	1	$A \wedge B \wedge C$

Wenn Ergebnis 0 dann  
invertiert und als Disjunktion

Wenn Ergebnis 1 dann  
als Konjunktion

Alle Disjunkte mit "und" verbinden  $\rightarrow$  Konjunktive Normalform: KNF

$$(A \vee B \vee C) \wedge (A \vee B \vee \neg C) \wedge (\neg A \vee B \vee C) \wedge (\neg A \vee \neg B \vee C)$$

Alle Konjunkte mit "oder" verbinden  $\rightarrow$  Disjunktive Normalform: DNF

$$(\neg A \wedge B \wedge \neg C) \vee (\neg A \wedge B \wedge C) \vee (A \wedge \neg B \wedge C) \vee (A \wedge B \wedge C)$$

"Vollkonjunktion"

Diese Repräsentation ist nicht optimal.

Wiederholung:

$(a \wedge \neg b \wedge c)$  Konjunktive Form

$(a \vee \neg b \wedge \neg c)$  Disjunktive Form

Negationen nur in atomarer Form!

$$\neg(a \wedge b) = \neg a \vee \neg b \rightarrow \text{nicht atomar}$$

$$(\neg a \vee \neg b) = \neg a \vee \neg b \rightarrow \text{atomar}$$

# Vereinfachung mit KV-Diagrammen

„Karnaugh map“

- Nicht immer minimal
- Nur bis 4 Variablen sinnvoll

	$\neg x_3$	$x_3$	$x_2$	$\neg x_2$	
$\neg x_1$	1	1	0	0	$\neg x_2$
$x_1$	0	1	1	0	$\neg x_2$
$x_1$	0	0	1	X	$x_2$
$\neg x_1$	1	0	1	1	$x_2$
	$x_4$	$x_4$	$\neg x_4$	$\neg x_4$	

Reinholtet alle möglichen Vollkonjunktionen

1 = kommt vor

0 = kommt nicht vor

X = don't care

↑  
eine mögliche minimale  
Blockeinteilung:

Wir suchen minimale Anzahl an Blöcken  
um alle 1en oder 0en abzudecken

- Per Block nur 2, 4, 8, 16, 32 ...

- Nachbarn unterscheiden sich nur  
um eine Variable:  $(\neg x \vee x) = 1$

Davor:

$$\neg x_1 \wedge \neg x_2 \wedge \neg x_3 \wedge x_4$$

$$\neg x_1 \wedge x_2 \wedge \neg x_3 \wedge x_4$$

„minimale disjunktive Normalform“

$$\rightarrow (\neg x_1 \wedge \neg x_2 \wedge \neg x_3 \wedge x_4) \vee (\dots) \vee (\dots) \vee (\dots)$$

↑  
weil 4 Blöcke  
gefunden

# ITE - If then else

$$a \wedge b \equiv \text{ITE}(a, b, 0)$$

$$\neg a \equiv \text{ITE}(a, 0, 1)$$

$$a \vee b \equiv \text{ITE}(a, 1, b)$$

Konversion;  
Wahrheitstabelle  $\leftrightarrow$  ITE

# Shannon-Zerlegung

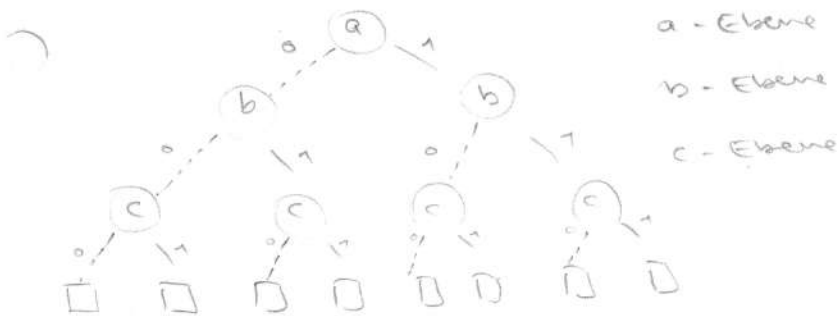
$f: B^n \mapsto B$  bool'sche Funkt.  
"Kofaktor von  $f$ " in der  $x_i = 1$   
 $f_{i,1}: B^n \mapsto B$

oder  
 $f_{i,0}: B^n \mapsto B$

$$f = (x_i \wedge f_{i,1}) \vee (\neg x_i \wedge f_{i,0})$$

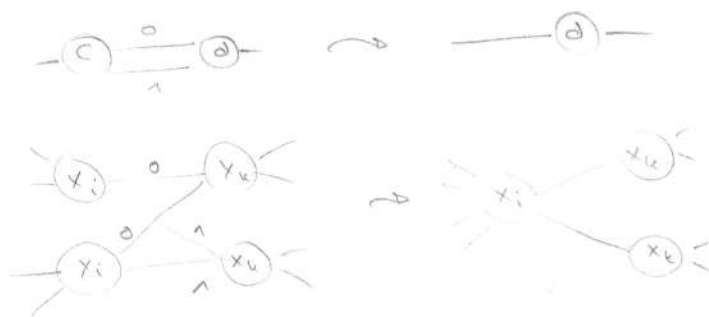
# Binary Decision Diagrams BDDs

Beispiel: Von if, else Baum

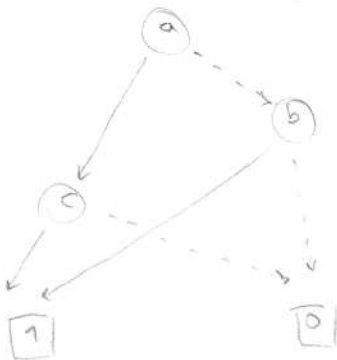


→ Alle Variablen zeigen auf die selbe 1/0

Vereinfachen mit  
delete - rule  
merge - rule  
(nicht minimal)



DNF ablesen, KNF ablesen



DNF  $(a \wedge c) \vee (\neg a \wedge b)$

KNF  $(\neg a \wedge c) \wedge (a \vee b)$



# Reduktion mit Beads

a	b	c	y
0	0	0	0
0	0	1	0
0	1	0	1
0	1	1	1
1	0	0	0
1	0	1	1
1	1	0	0
1	1	1	1

↑  
Abhängig davon  
ob man a,b,c oder  
eine andere Permutation  
hat

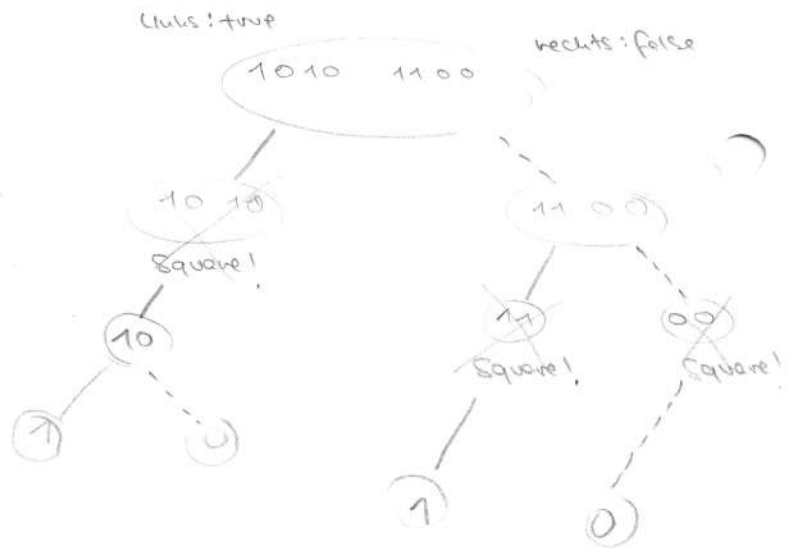
Bead: erste und zweite  
Hälfte stimmen nicht

Square: 7 Bead

Bei "don't cares":

"Greedy: Gleich am Anfang

Lazy: Aufschieben

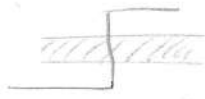


07: Kombinatorische Logik, Digitalschaltungen, kombinat. Digitaltechnik  
 Bool'sche Algebra mit digitalen Systemen

### Diskretisierung

kont. analoge Signale  $\rightarrow$  bool'sche Werte

pos. Logik



log 0 log 1

neg. Logik



log 0 log 1

### Bool'sche Funktionen durch Gatter-Bausteine



Bool'sche Operatoren sind „funktional vollständig“ wenn mit ihnen alle anderen Operat. hergeleitet werden können;

$\{ \wedge, \vee \}, \{ \text{NAND} \}, \{ \text{NOR} \}$

weniger Transistoren,  
billiger herzustellen

$\rightarrow$  Substitution mit NAND und NOR

### Standard-Baugruppen

Encoder: verdichtet Code

n-zu-m-Encoder

n Eingänge  $\geq$  m Ausgänge

Decoder

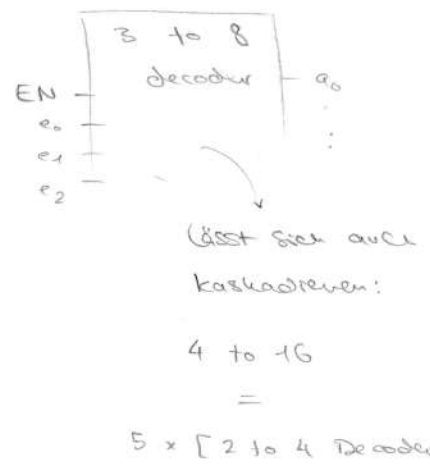
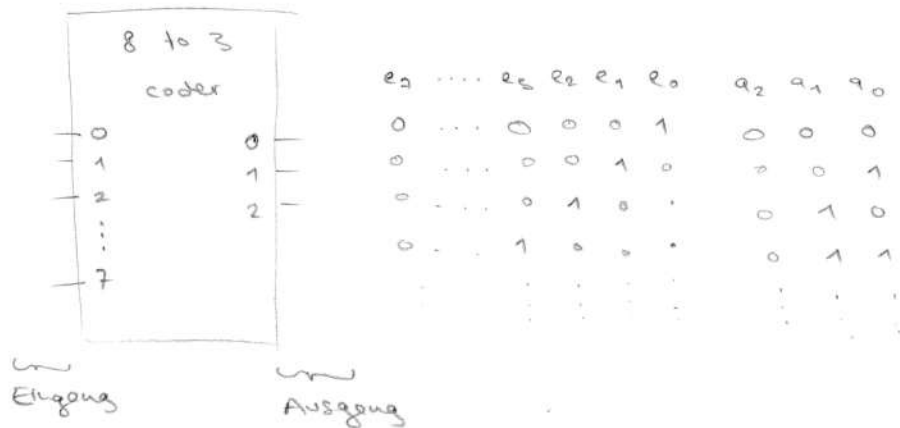
n zu m Decoder

n Eingänge  $\leq$  m Ausgänge

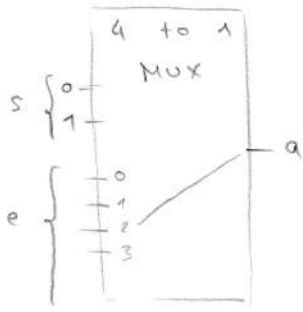
Beispiel:

„1 aus 2“-Code  $\rightarrow$  n Bit Binärzahl

Für Decoder umgekehrt.



# Multiplexer



$S_1$	$S_0$	$a$
0	0	$e_0$
0	1	$e_1$
1	0	$e_2$
1	1	$e_3$

# Demultiplexer

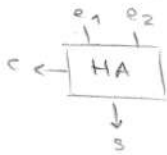
Äquivalent zu Decoder mit Enable-Eingang



Eingang = Enable Input

# Addierer

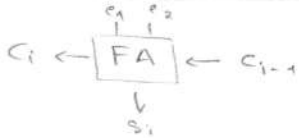
## HA Half Adder



$e_1$	$e_2$	$e_1 + e_2$
0	0	0 0
0	1	0 1
1	0	0 1
1	1	1 0

carry  $(e_1 \wedge e_2)$       sum  $(e_1 \oplus e_2)$

## FA Full Adder

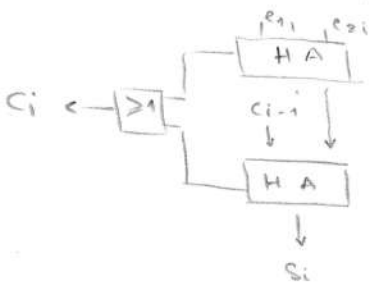


$e_{1i}$	$e_{2i}$	$C_{i-1}$	$C_i$	$S_i$
0	0	0	0	0
0	0	1	0	1
0	1	0	0	1
0	1	1	1	0
1	0	0	0	1
1	0	1	1	0
1	1	0	1	0
1	1	1	1	1

↑  
kanonische DNF

$$S_i = (\neg e_{1i} \wedge \neg e_{2i} \wedge C_{i-1}) \vee (e_{1i} \wedge \neg e_{2i} \wedge C_{i-1}) \vee (\neg e_{1i} \wedge e_{2i} \wedge C_{i-1}) \vee (e_{1i} \wedge e_{2i} \wedge \neg C_{i-1})$$

Löst sich mit Kaskadierung von HA implementieren:

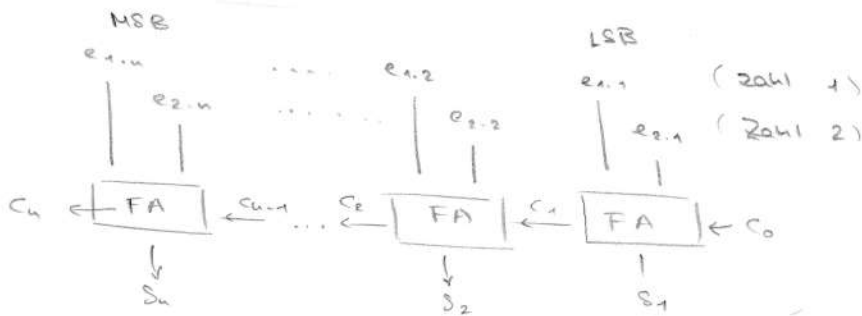


Es gilt:

$$S_i = e_{1i} \oplus e_{2i} \oplus C_{i-1}$$

$$C_i = (e_{1i} \wedge e_{2i}) \vee (C_{i-1} \wedge (e_{1i} \oplus e_{2i}))$$

## Paralleladdierer



Auch Subtraktion möglich:  
eine Zahl invertieren,  
erstem FO:  $C_0 = 1$

# Beispiele: Substitution mit NAND / NOR:

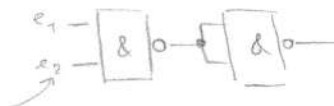
## AND Substitution



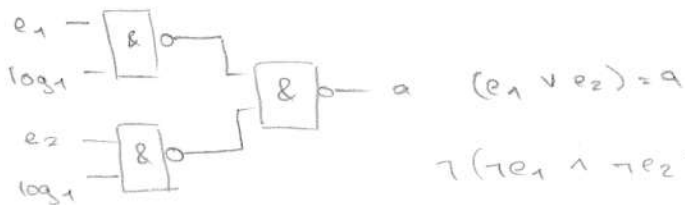
$$(e_1 \wedge e_2) = a$$

$$\neg(\neg(e_1 \wedge e_2) \wedge 1) = \neg(\neg e_1 \vee \neg e_2) \wedge 1$$

$$= e_1 \wedge e_2 \vee 0$$



## OR Substitution



$$(e_1 \vee e_2) = a$$

$$\neg(\neg e_1 \wedge \neg e_2) = (e_1 \vee e_2)$$

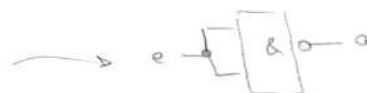


## NOT Substitution



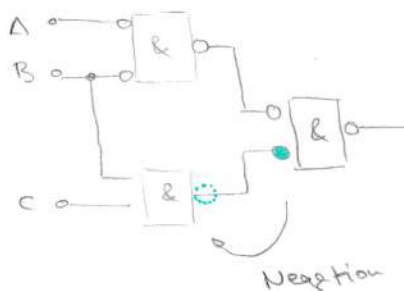
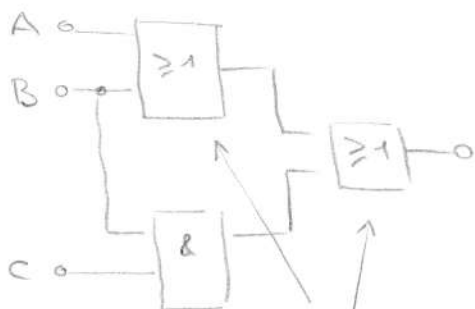
$$\neg e = a$$

$$\neg(e \wedge 1) = \neg e \vee 0$$

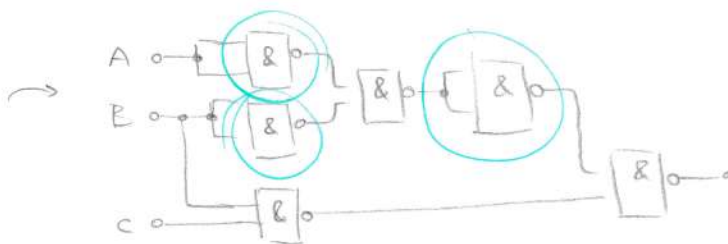
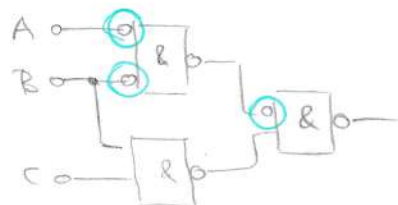


## Beispielfunktion: Ersetzen der Gatter mit NAND

$$f(A, B, C) = (A \vee B) \vee (B \wedge C)$$



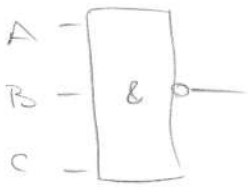
OR ersetzen:  
 $(X_1 \vee X_2) = \neg(\neg X_1 \wedge \neg X_2)$



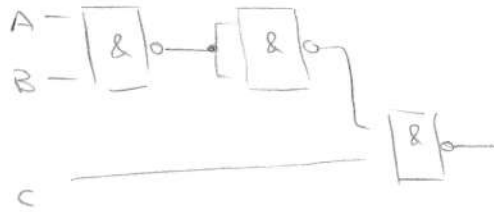
Jetzt wurden alle OR sowie AND mit NAND ersetzt

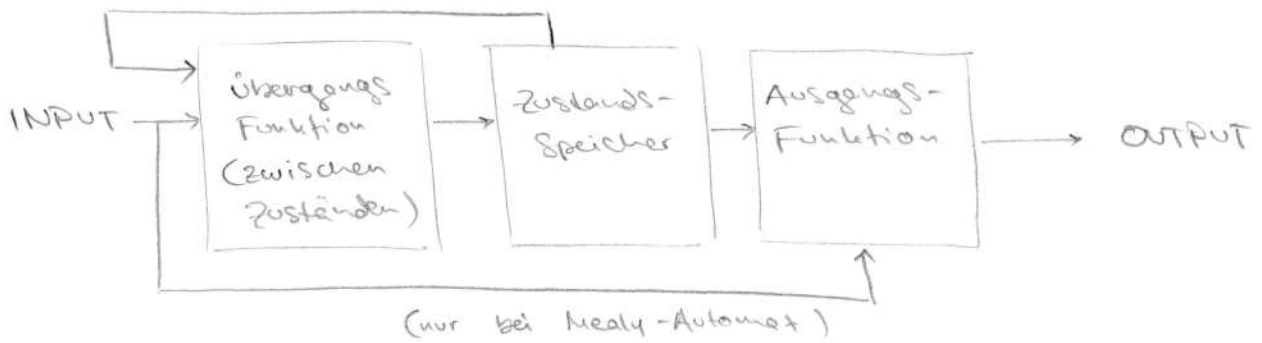
Aber wir haben noch 3 Negationsgatter

### 3 NAND Gatter Eingänge



≡





### Formale Definition

Systeme mit Zustand (Schaltwerk) werden mit Automaten beschrieben.

Endliche Automaten:

- endliche Zustandsmenge  $Q$
- Startzustand  $s \in Q$
- Endzustandsmenge  $F \subseteq Q$
- Übergangsfunktion  $\delta$
- Endliches Eingabealphabet  $\Sigma$

### Moore-Automat

- Ausgabealphabet  $\Omega$ , und Ausgabefunktion  $\lambda: Q \mapsto \Omega$

### Mealy-Automat

- Ausgabealphabet  $\Omega$ , und Ausgabefunktion  $\lambda: Q \times \Sigma \mapsto \Omega$
- $\uparrow$                      $\uparrow$   
 Zustand            Eingabesymbol

### ○ Deterministische Automaten:

- Nachfolgezustand genau definiert
- $$\delta: Q \times \Sigma \mapsto Q$$

### Nicht-Deterministische Automaten

- $\exists$  ein Zustand-Eingabepaar ohne Nachfolge-Zustand
- $$\delta \subseteq Q \times \{\Sigma \cup \{\epsilon\}\} \times Q \quad \rightarrow \exists \text{ Zustände ohne Eingaben}$$
- $\uparrow$   
Leerwort

### ○ Vollständige Automaten

- $\forall (z, x) \in Q \Rightarrow \exists z' : \delta(z, x) = z'$

Wir gehen davon aus:

Automat vollständig & deterministisch

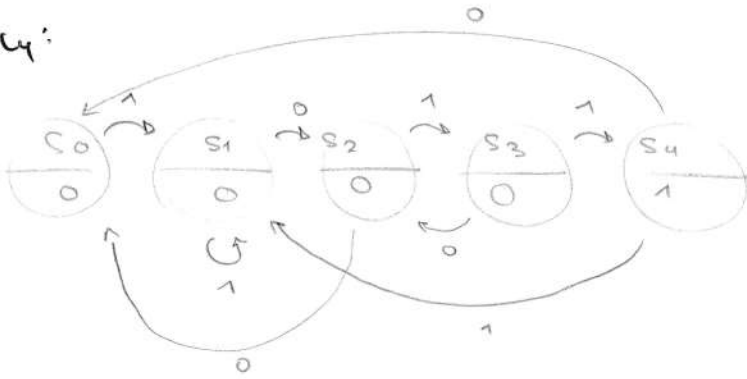




Aufgabenstellung:

Output soll 1 sein, wenn Input Bitfolge 1011 ist,  
 (Überlappungen (1011011) auch bewerten!)

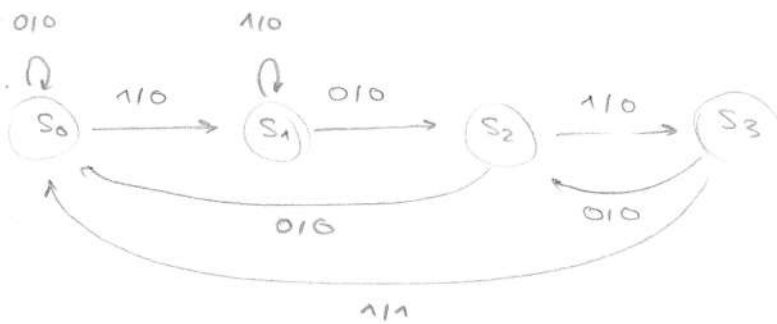
Mealy:



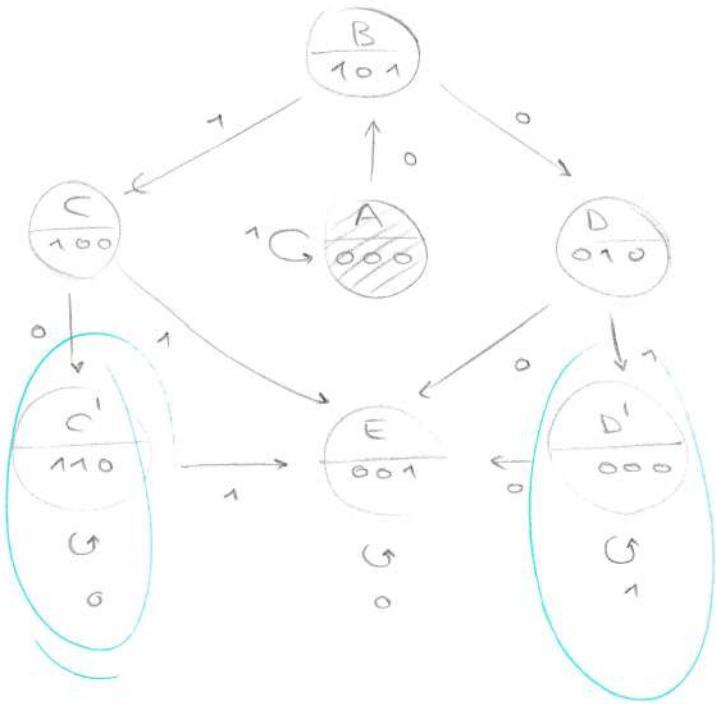
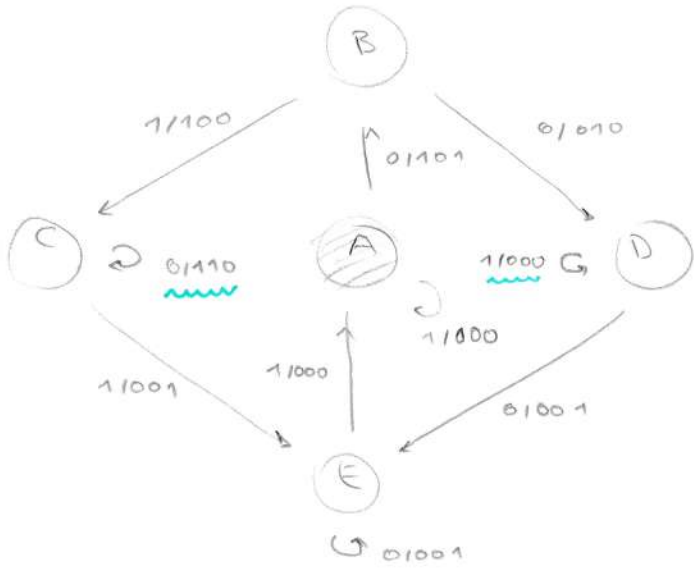
Zustand	gelesene Folger
S <sub>0</sub>	-
S <sub>1</sub>	1
S <sub>2</sub>	10
S <sub>3</sub>	101
S <sub>4</sub>	1011

Bei jedem Übergang wird ein Bit vom Input „konsumiert“.

Moore:



# Konversion / Transformation

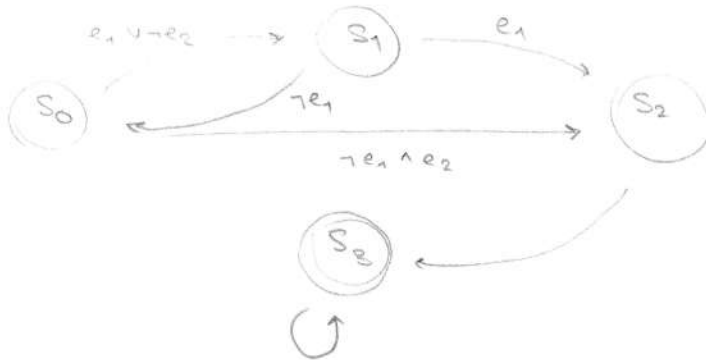


# Zustandsgraph - Beispiel

$$Q = \{s_0, s_1, s_2, s_3\}$$

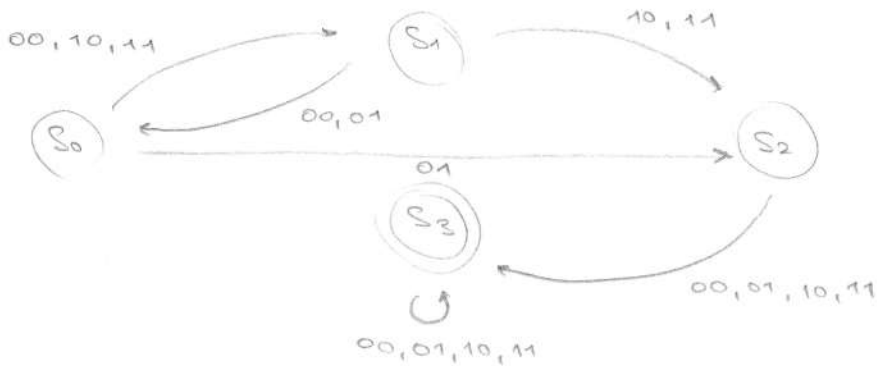
$$s = s_0$$

$$F = \{s_3\} \quad (\text{Wenn kein } F \text{ mit 2 Ringen eingezeichnet, dann } \forall \text{ Knoten} = \text{Endzust.})$$



Wir können  $e_1, e_2$  mit Bitmuster angeben:

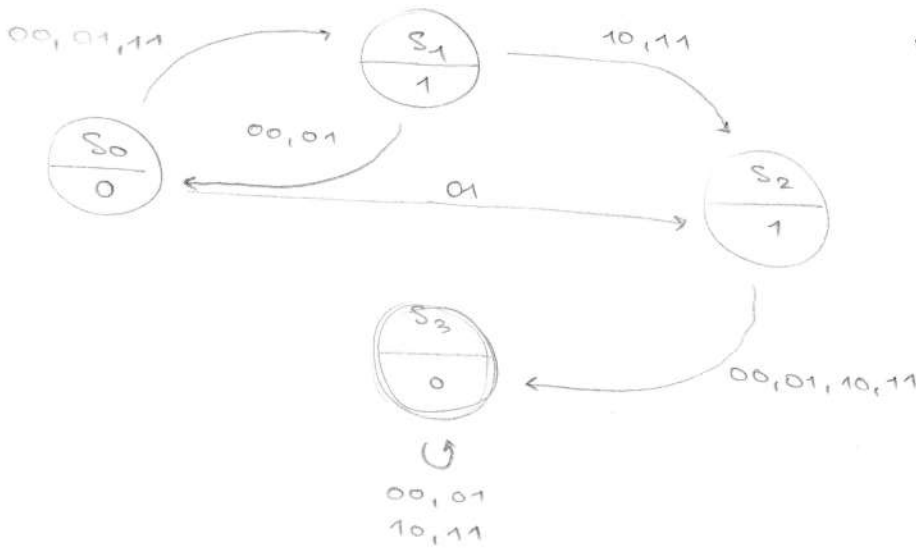
Reihenfolge:  $e_1 e_2$



# Moore Schaltwerk

Zustand bestimmt Ausgabe:

$$\lambda = \{ (S_0, 0), (S_1, 1), (S_2, 1), (S_3, 0) \}$$



Reihenfolge der Eingabe:  $e_1, e_2$

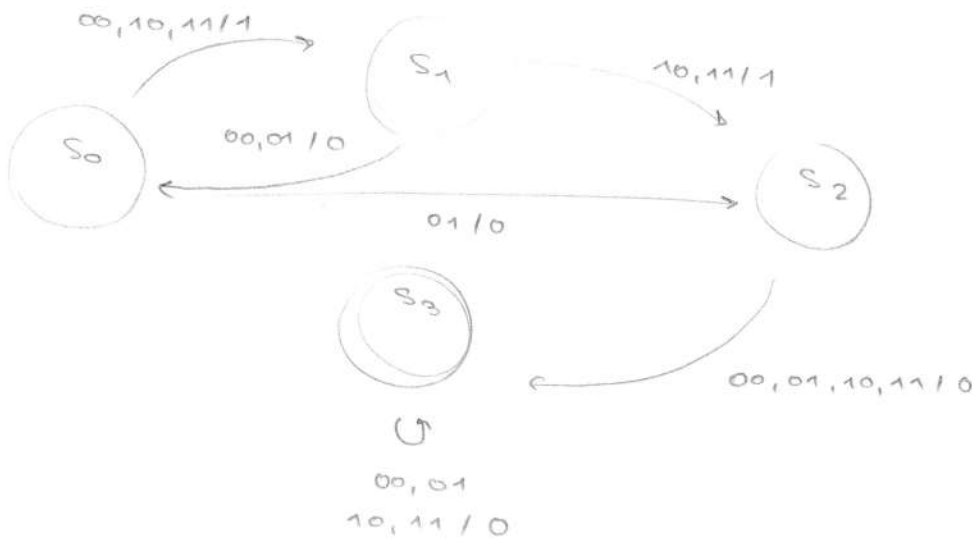
Ausgabe:  $a$

↑  
immer  
Ausgeben!

# Mealy Schaltwerk

Zustand UND Eingabe bestimmen Ausgabe:

$$\lambda = \{ ((S_0, 00), 1), \dots, ((S_0, 01), 0), \dots \}$$



# 09 : Sequenzielle Logik

## Kombinatorische Logik

Output vom aktuellen Input-Zustand abhängig

Keine Speicherung von Informationen / Zuständen.

## Sequenzielle Logik

Output vom aktuellen Input und bisherigem Ereignisverlauf (Historie) abhängig

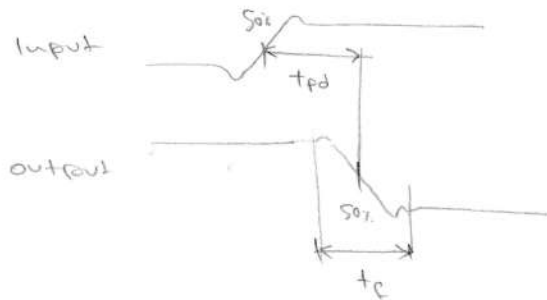
Mit Speicherelementen

## Zeitverhalten

Idealisiert:



Real:



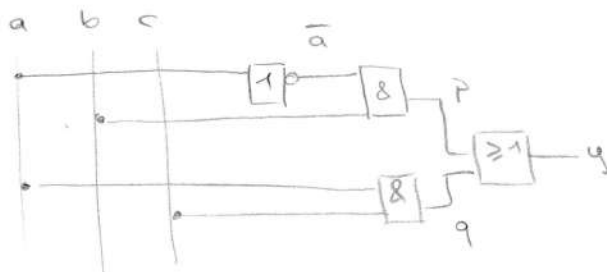
$t_{pd}$ : Durchlaufzeit

$t_f$ : Flankensteilheit

## Electronics Glitch / logical hazard

Wenn man keinen Takt hat:

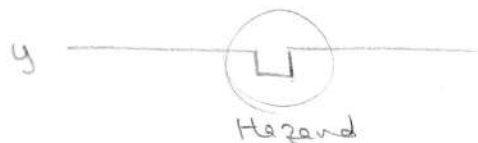
output kann kurzfristig falsch sein



$$f(a, b, c) = (\neg a \wedge b) \vee (a \wedge c)$$

Wenn b und c wahr sind, muss y unabhängig von a wahr bleiben.

Wenn  $a=1 \rightarrow a=0$  wird q vor p true sein



# Speicherelemente: Latches und Flipflops

Speichern 1 Bit (so lange Spannung anliegt)

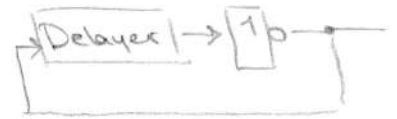
Taktzustand → Latches

Taktflanken → Flipflops

## Takt-Signal (Clock-Signal) [Hz]

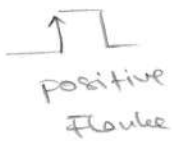
Synchronisiert (vor allem große) Schaltungen

- synchrone Takung (periodisch)
- asynchrone Takung



Stabiliter:  
Quarzoszillator

Flanke = Zustandswechsel



positive  
Flanke



negative  
Flanke

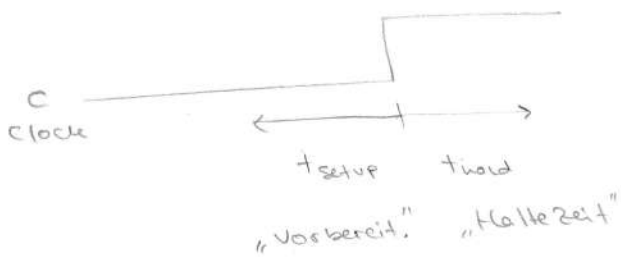


durch pos  
flanke  
getriggert



durch neg  
flanke  
getriggert

## Flipflop / Zeitspezifikation beim Takten:



ZB bei D-Flipflop:  
Eingang muss in Vorher, und Haltezeit  
stabil bleiben.  
sonst metastabile Zustände

## Glitches und Hazards

Wenn am Ausgang kurz falscher Wert steht, erzeugt durch das nicht-Einhalten von  $t_{\text{setup}}$ ,  $t_{\text{hold}}$  und generell durch  $t_{\text{pd}}$

## Metastabiler Zustand

unkontrolliertes Togglen.

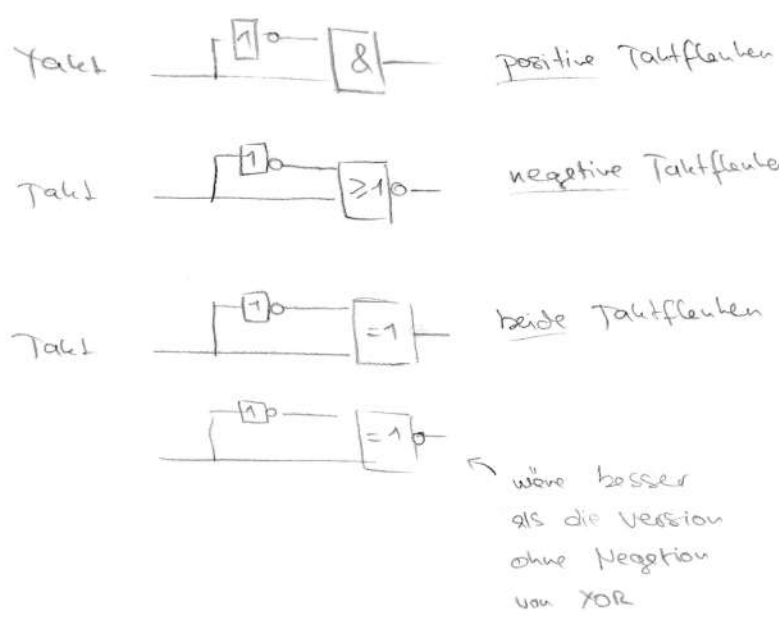
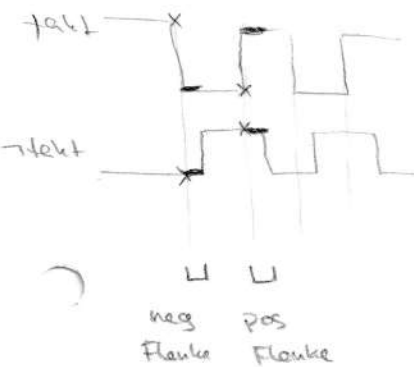
ZB wenn man bei RS-Latch gleichzeitig von  $R=S=1$  auf  $R=S=0$  wechselt  
 unerlaubte Eingabe

## Bistabile Kippstufe

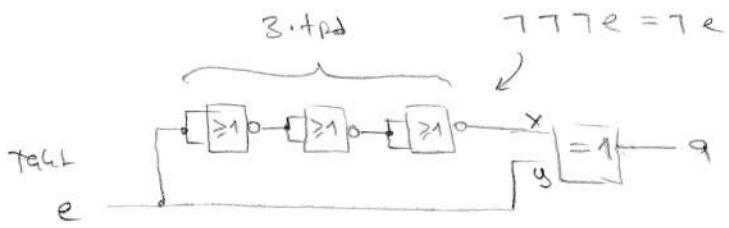
Schaltungen die 1 Bit speichern können besitzen 2 stabile Zustände  
 ZB Set, Reset

## Monostabile Kippstufe

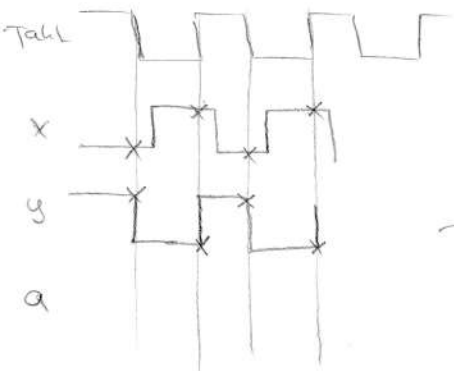
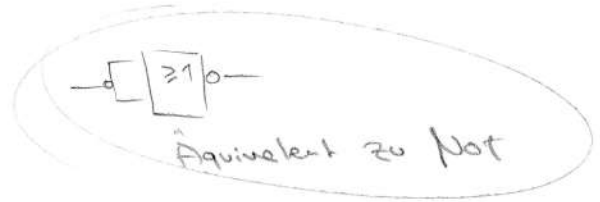
Nur ein Zustand ist stabil  $\rightarrow$  so entsteht ein Taket: es fällt immer noch der Flanke zurück zum stabilen Zustand



# Lösung des Beispiels:



	NOR	XOR
00	1	0
01	0	1
10	0	1
11	0	0

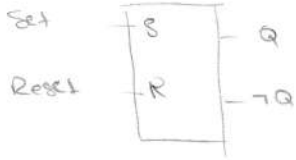


für negative Flanke  $x=0$   $y=1$   $\rightarrow 1$   
 für positive Flanke  $x=1$   $y=0$   $\rightarrow 1$

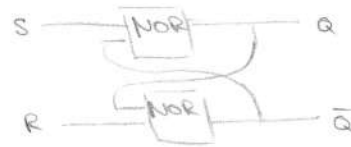


# Speicherelemente

## RS-Latch

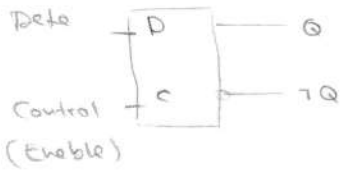


S	R	State
0	0	Hold
0	1	Reset = 0
1	0	Set = 1
1	1	Nicht erlaubt → Metastabilität



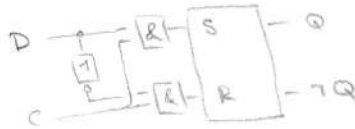
Wenn  $R=S=1$  und dann  $R=S=0$   
dann ständiger Wechsel

## D-Latch



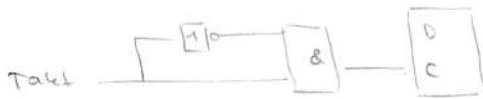
Baut auf RS-Latch auf:

$$Q = \overline{Q} \iff Q = \overline{\overline{Q}} = \overline{0} = 1$$



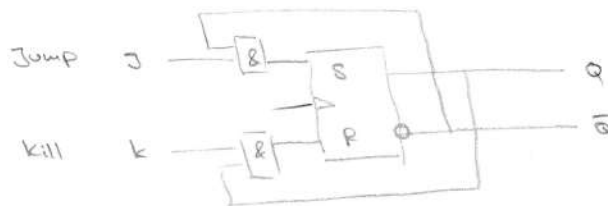
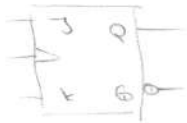
$D=1 \rightarrow$  Set  
 $D=0 \rightarrow$  Reset  
 $C=0 \rightarrow$  Hold

## D-Flip-Flop



Bewusst erzeugter Hazard: Erkennung von Taktflanke

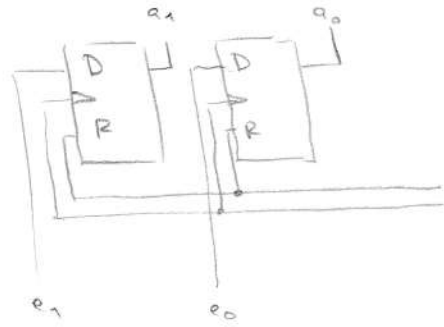
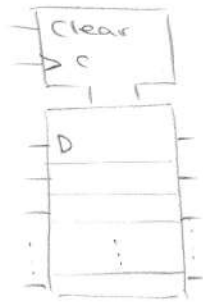
## J-K Flip-Flop



J	K	State
0	0	Hold
0	1	0
1	0	1
1	1	Toggle

# Anwendungen

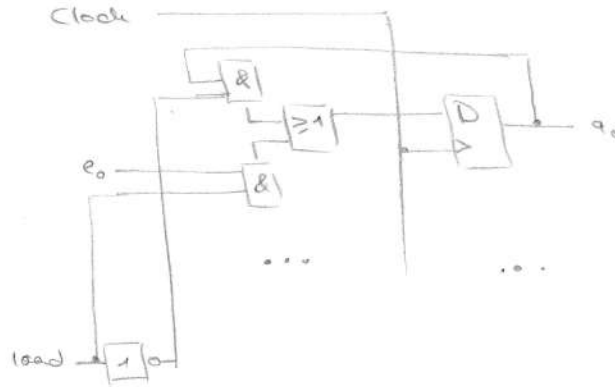
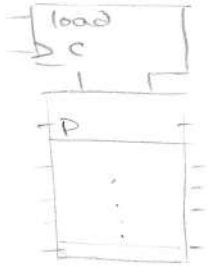
## Register



bei RS Latch in D-Flipflop ist R noch erreichbar

Clear  
Clock

## Parallel-Register



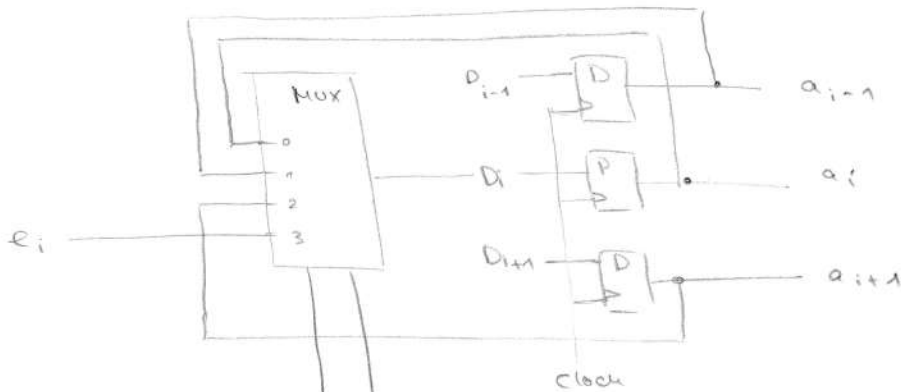
wenn load: Zustand übernehmen, ansonsten load

## Schiebe-Register



Anwendung: paralleles Addieren

## Universal-Register

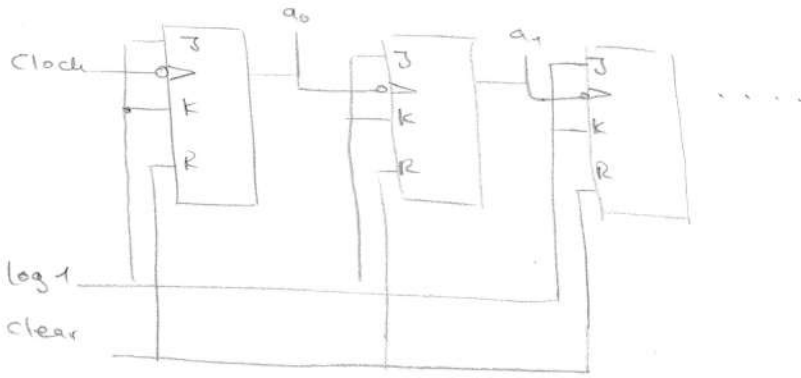


	$s_0$	$s_1$	
0	0	0	—
1	0	1	shift left
2	1	0	shift right
3	1	1	parallel load

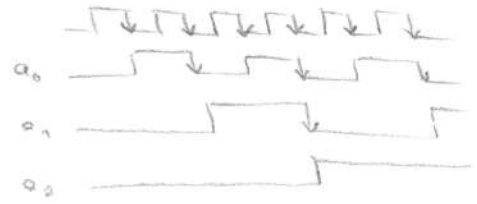
Alle Mux werden gleichzeitig mit selber Angabe geschaltet

# Zähler

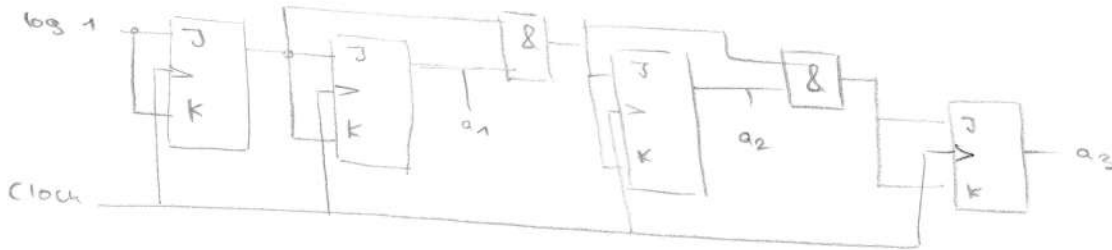
(asynchron getaktet)



Dadurch, dass  $J=K=1$  wird getoggled bei jeder negativen Flanke:



(synchron getaktet)



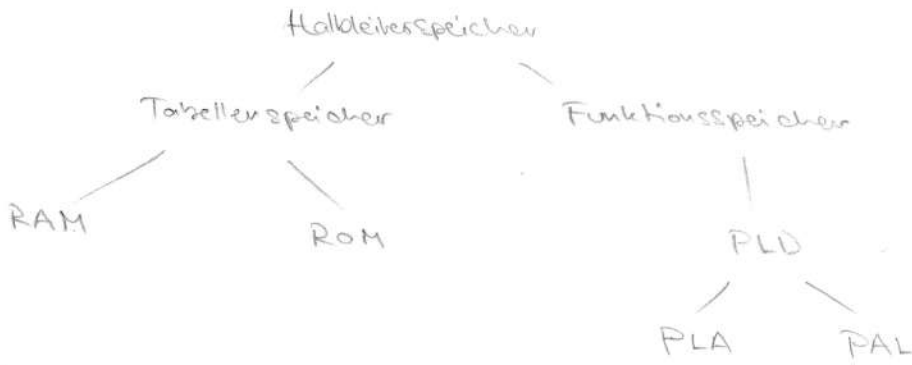
↑  
toggled  
mit jedem  
Takt

↑  
toggled  
wenn  $a_0 = 1$

↑  
toggled  
wenn  $a_0 \& a_1 = 1$

↑  
toggled  
wenn  $a_0 \& a_1 \& a_2 = 1$

# 10: Speicher



Tabellenspeicher speichern Werte in Zeilen die ansprechbar sind (Gedächtnis).

Funktionspeicher speichern bool'sche Funktionen (kein Gedächtnis)

Funkt. Speicher lassen sich als Tabell. Speicher implementieren, vice versa

## Funktionspeicher: PLD

PLD programmable logic device

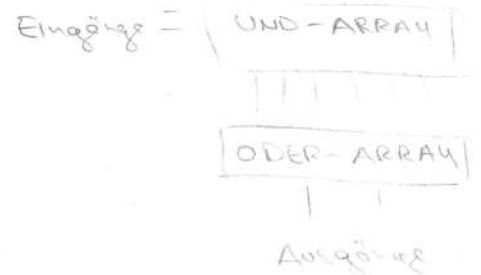
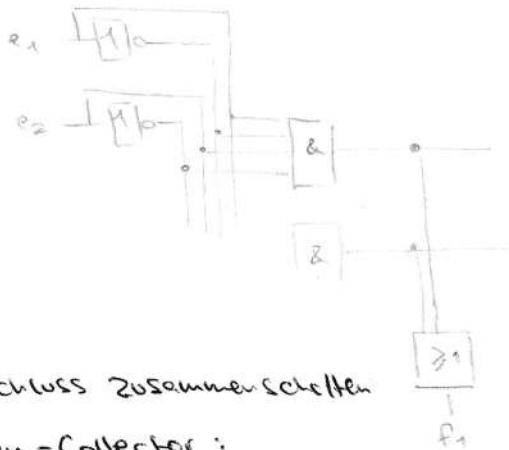
PLA programmable logic array — UND-ODER Matrizen

PAL programmable array logic — ODER ist undefiniert

Vom Anwender 1x programmierbar:

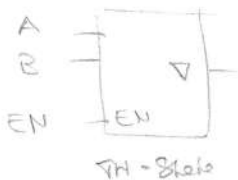
Diode / Fuse (~~—~~) verbindet zwei Leitungen und kann durchgebrannt werden mit versch. Strömstärken

Funktion als disjunktive Normalform DNF:



Ausgänge ohne Kurzschluss zusammenschalten

Tri-State und Open-Collector:



Leitet nicht wenn enable auf 0 ist



# ASIC

Application specific integrated circuit

- Schnell, effizient, nur 1 Funktion
- teuer

→ ZB für Handy: GSM-Protokolle fix

ES gibt auch **FPGA** (field programmable gate array) wie PLA aber effizienter.

## Tabellenspeicher: ROM

Read-only-memory

- Permanente Speicherung, nicht flüchtig
- Mehrfachprogrammierbar

1k → # Datenwörter, # Zeilen (1024)

8 → Bitbreite



Beispiel: 16 x 4 ROM

Mit einem 4 zu 16, wird eine Zeile gelesen

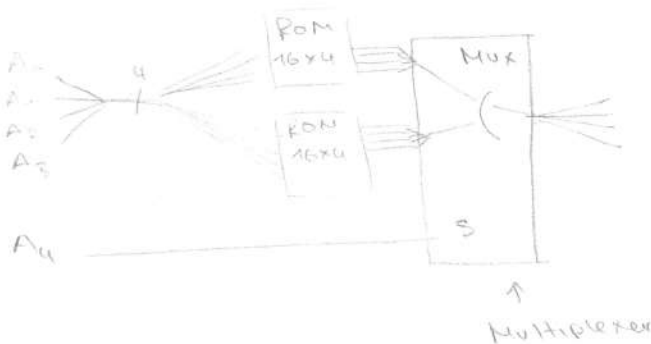
↑            ↑  
# Eingang   # Ausgang



Zum Programmieren wird ein Anti-Fuse benutzt (Beim verschmelzen wird Kontakt erzeugt.)

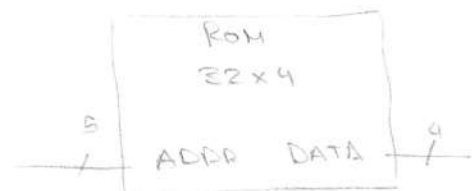
## Erweiterung von ROM: Multiplexer

Statt einem 32 x 4 ROM kann man zwei 16 x 4 ROMs verwenden

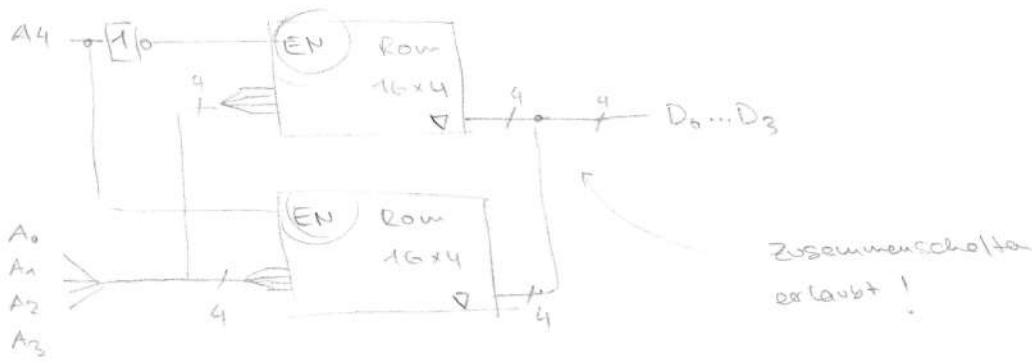


ROM 1 für Zeile 0 bis 15  
ROM 2 für Zeile 16 bis 31

Dadurch:

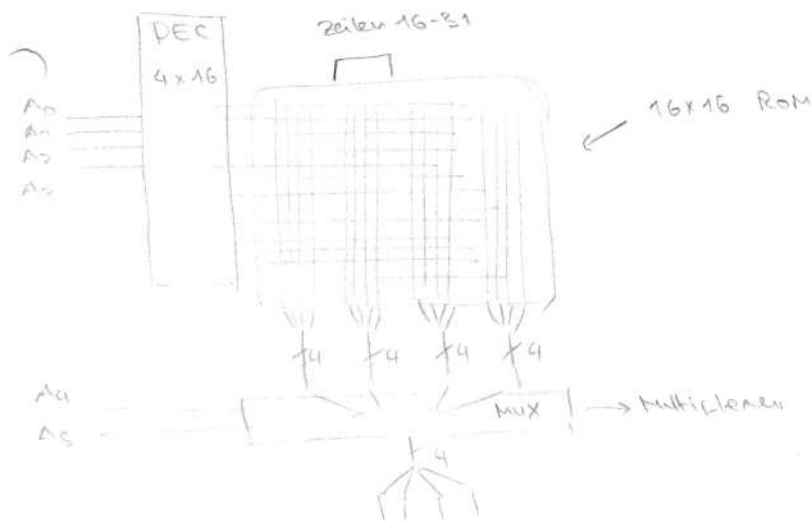


## Erweiterung von ROM: Tri-State



## Erweiterung von ROM: Multiplexer

Aus 16x16 ROM → 64x4 ROM



## Anwendungen von ROM heute:

### Flash EEPROM

USB, SSD, sehr schnell

### EEPROM

Speicher von Betriebssystemen, Embedded Systemen → Daten die selten geändert werden

## ROM vs FAT:

- Nicht flüchtig
- langsam
- hohe Kosten pro GB
- begrenzte # an Löschzyklen

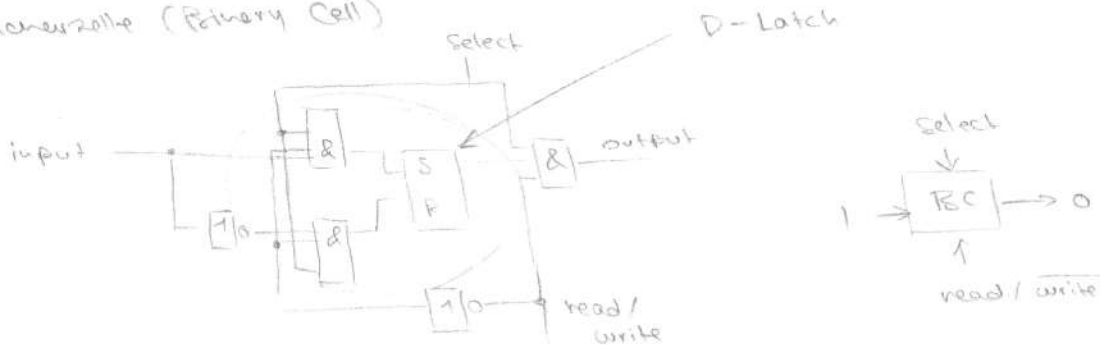
# Tabellenspeicher : RAM

random-access-memory (historisch bedingt)

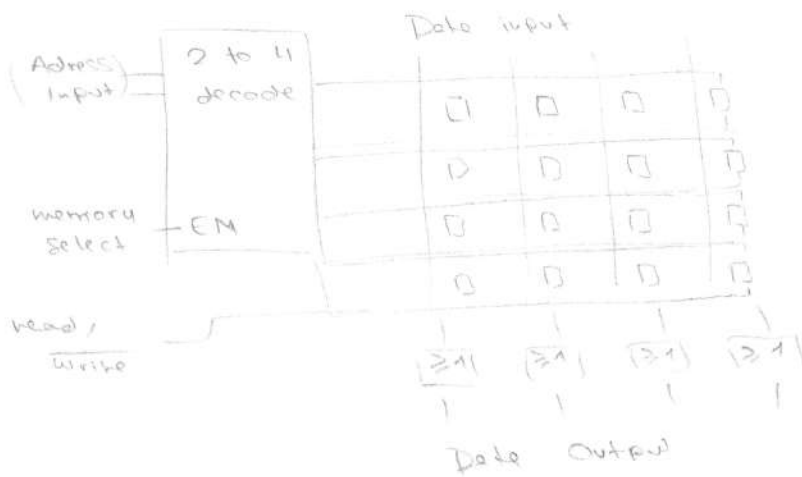
DRAM : flüchtig, Speicherung im Kondensator

SRAM : nicht flüchtig, aus CMOS Transistoren

## Speicherzelle (1T1C1)



Matrix aus binary cells:



Kaskadierung wie bei ROM mit decoder möglich.

# 1.1: Schaltwerke

## Schaltwerke

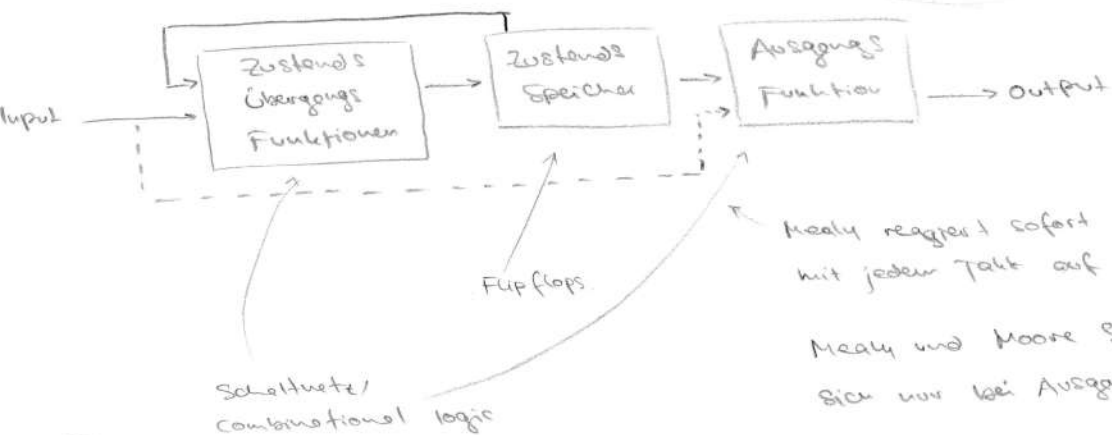
### Combinational Logic

logische Funktionen ohne Speicher / internen Zustände

### Sequentielle Schaltwerke

#### Sequential Logic

Ausgang abhängig von internem Zustand und Eingang.



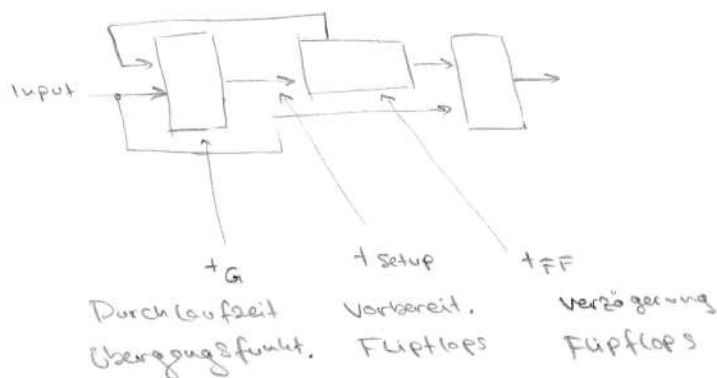
Mealy und Moore Schaltwerke unterscheiden sich nur bei Ausgangsfunktion.

## Taktung, Vermeidung von Meta-Zuständen

### Annahmen:

- immer synchron: Änderung nur bei Taktsignal
- Eingangssignal auch synchron mit Takt, weil mealy sofort auf input reagiert

Zustand muss aktualisiert werden bevor nächste Eingabe folgt:



### Minimaler Taktabstand:

$$T_{\min} = t_{FF} + t_G + t_{\text{setup}}$$

### Maximaler Takt $f_{\max}$ :

$$f_{\max} = \frac{1}{T_{\min}}$$



# Zeitverhalte

Abstraktion:

NBT Gatter

Eingang



Ausgang



Realität:

ES braucht Zeit bis Ladungsträger fließen

Eingang



Ausgang



$t_{pd}$

propagation delay  $t_{pd}$

Durchlaufzeit

Von der Taktflanke bis zur Ausgangsänderung, quasi Reaktionszeit

rise / fall time  $t_f$

Aufstiegs / Abfallzeit

Steilheit der Flanke

Weiters:

Bei Flipflops (zB D-Flipflop) muss Eingang während Taktflanke stabil bleiben, weil sonst metastabile Zustände entstehen können (unkontrolliertes Folgen)

Setup time  $t_{setup}$

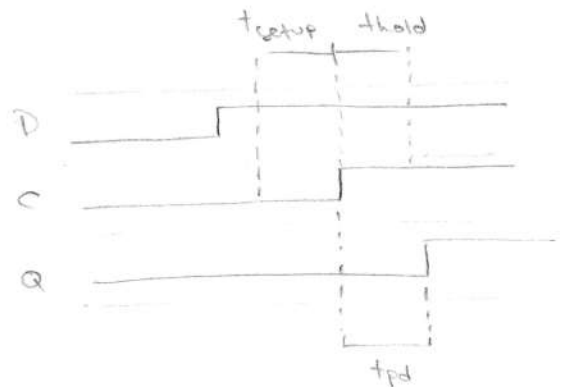
Vorbereitungszeit

Zeit vor Taktflanke

hold time  $t_h$

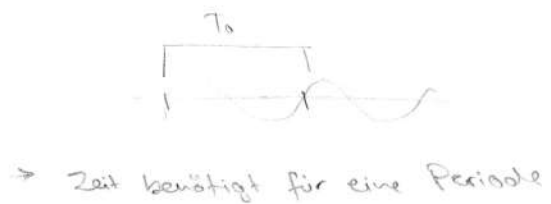
Haltezeit

Zeit nach der Taktflanke



# Maximale Taktfrequenz

$$\text{Frequenz [Hz]} = \frac{1}{T_0 \text{ [s]}}$$



SI-Präfixe:

T	Tera	$10^{12}$
G	Giga	$10^9$
M	Mega	$10^6$
k	kilo	$10^3$
h	hekto	$10^2$
da	deka	$10^1$
		$10^0$
d	dezi	$10^{-1}$
c	Zenti	$10^{-2}$
m	milli	$10^{-3}$
μ	mikro	$10^{-6}$
n	nano	$10^{-9}$
p	piko	$10^{-12}$

Sollten  
benutzt

Beispiele:

$$1 \mu\text{s} = 10^{-9} \cdot \text{s}$$

$$1 \text{ MHz} = 10^6 \cdot \text{Hz} \quad : 1 \text{ Mio Zyklen pro Sekunde}$$

## Beispiele zur Berechnung der maximalen Taktfrequenz von Flipflops

Benötigte Angaben:

Gatter Durchlaufzeit	15 ns
Flipflop Durchlaufzeit von Takt bis Ausgang	45 ns
t <sub>setup</sub> für Flipflops	5 ns
	65 ns

$$\begin{aligned} &\rightarrow \frac{1}{65 \cdot 10^{-9}} = \\ &= 15\,384\,615,384 \dots = \\ &= 15,4 \text{ MHz} \end{aligned}$$

Weiteres Beispiel:

Gatter t <sub>pd</sub>	30 ns
Flipflop t <sub>pd</sub>	40 ns
Flipflop t <sub>setup</sub>	10 ns
	90 ns

$$\rightarrow \frac{1}{90 \cdot 10^{-9}} =$$

$$= 11,1 \text{ MHz}$$

↓

$$11,1 \text{ MHz} < 25 \text{ MHz}$$

$$\begin{aligned} &15,4 < 25 \text{ MHz} \\ &\text{Max flipf.} \quad \text{Maximale Taktfrequ.} \end{aligned}$$

↑  
keine obere  
Schranke!

# Systematische Schaltwerkentwicklung

1. Aufgabenstellung verstehen
  - prüfe Vollständigkeit

2. Entwurf des Zustandsgraphen (ES gäbe auch andere Wege Zustandsübergänge zu modellieren)
  - prüfe Vollständigkeit für alle mögl. Inputs
  - prüfe Eindeutigkeit

Deterministic finite state machine

state machine: Zustandsabhängiger konzeptueller Automat

finite: Endliche # an Zuständen

deterministic: Eindeutige Nachfolge Zustände

3. Minimierung der Zustands #
  - automatisiert

4. Zustandscodierung
  - 1 aus  $n$  oder „dichte Codex.“

5. Übergangs- und Ausgangsfunktion bestimmen
  - aus Zustandsgraph ableiten
  - mit KV-Diagramm vereinfachen

6. Dokumentation der Gesamtschaltung

7. Berechnung von  $f_{max}$  Taktfrequenz

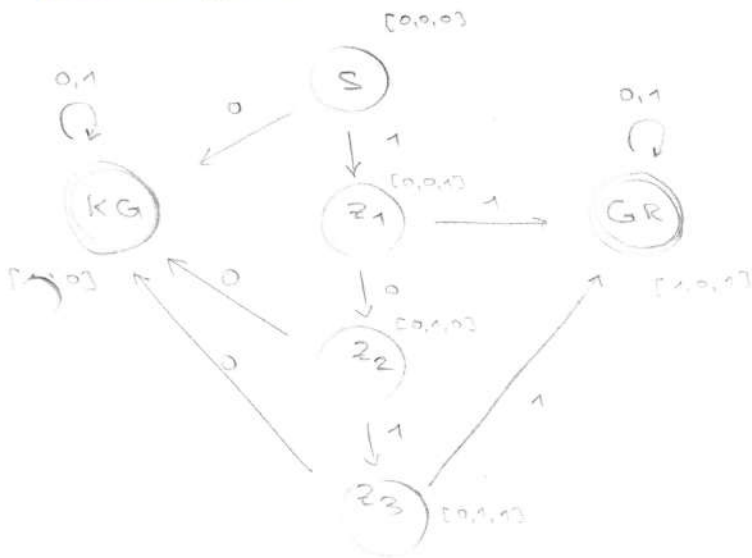
# Systematische Schaltwerkentwicklung

## Beispiel

### Angabe:

Schaltung mit Eingang E, durch die 4 Fkt. geschaltet werden (HSE zuerst)  
 Soll erkennen ob die Zahl (basis 2)  $\leq (10)_{10} = (1010)_2$  ist.

### Zustandsgraph:



### Zustandscodierung:

Zustand	Beschreib.	D <sub>2</sub>	D <sub>1</sub>	D <sub>0</sub>
S	B3 (HSE)	0	0	0
Z <sub>1</sub>	B2	0	0	1
Z <sub>2</sub>	B1	0	1	0
Z <sub>3</sub>	B0 (LSE)	0	1	1
KG	$\leq 10$	1	0	0
GR	$> 10$	1	0	1

Bezieht sich auf Stelle in der Zahl

### Zustandsübergangstabelle

alter Zustand			E	neuer Zustand		
D <sub>2</sub>	D <sub>1</sub>	D <sub>0</sub>		D <sub>2</sub> '	D <sub>1</sub> '	D <sub>0</sub> '
0	0	0	0	1	0	0
0	0	0	1	0	0	1
0	0	1	0	0	1	0
0	0	1	1	1	0	1
0	1	0	0	1	0	0
0	1	0	1	0	1	1
⋮						

Eindeutige Codierung für jeden Zustand  
 ideale Codierung

→ Mit KV Diagrammen vereinfachen

$$D_2' = D_2 \vee (D_0 \wedge E) \vee (D_1 \wedge \neg E) \vee (\neg D_0 \wedge \neg E)$$

$$D_1' = (\neg D_0 \wedge D_1 \wedge E) \vee (D_0 \wedge \neg D_1 \wedge D_2 \wedge E)$$

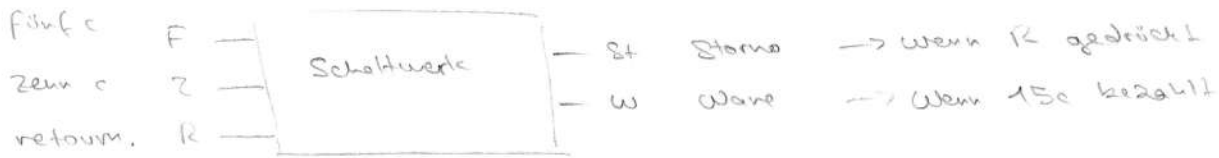
$$D_0' = (\neg D_2 \wedge E) \vee (D_0 \wedge D_2)$$

Abschließend:

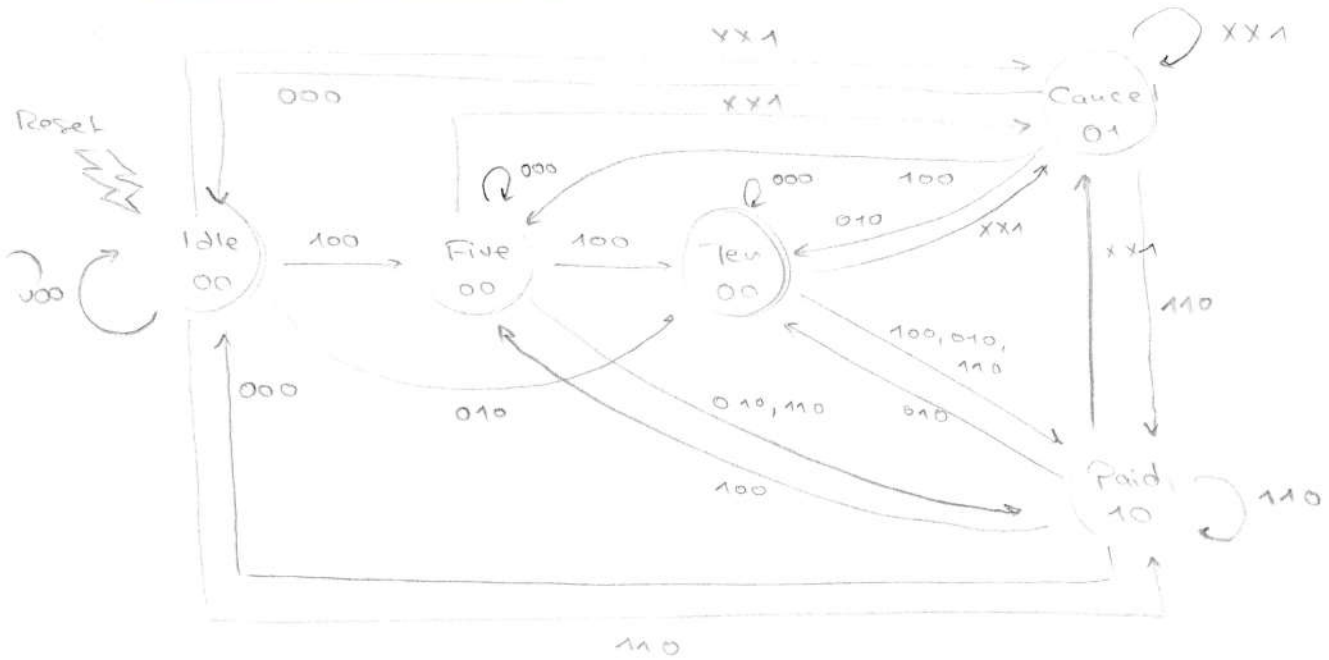
Entwurf des Schaltwerks

# 12 : Realisierung von Schaltwerken

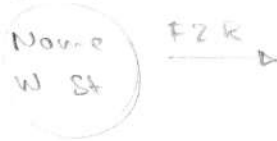
Beispiel: Münzautomat



## Moore Zustandsgraph



Notation:



Wir wandeln den Zustandsgraphen, von der Zustandscodierung in eine Tabelle für die Übergangsfunktion um!

Input	F	0	...	0	...
	Z	0	...	1	...
	R	0	...	0	...
Aktueller Zustand	Idle	Idle	...	Ten	...
	Five	Five	...	Paid	...
	Ten	Ten	...	Paid	...
	Paid	Paid	...	Ten	...
	Cancel	Cancel	...	Ten	...

→  
neuer Zustand

# Zustandscodierung

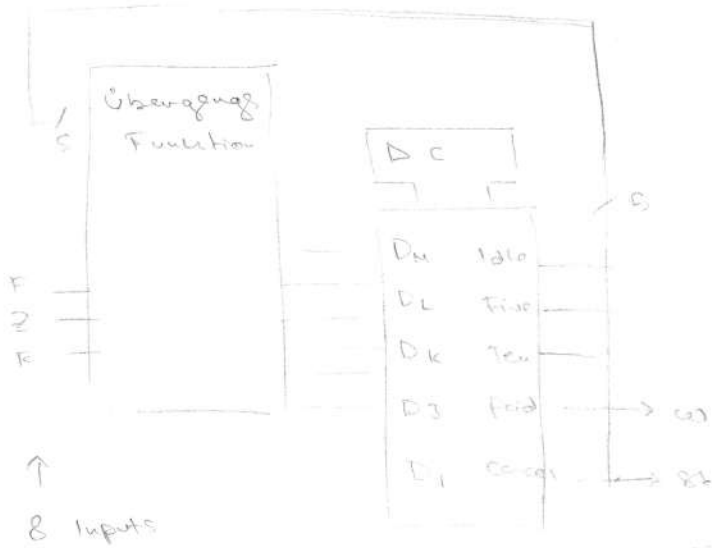
1 aus n - Codierung

I	J	K	L	M	Zustand
0	0	0	0	1	Idle
0	0	0	1	0	Five
0	0	1	0	0	Ten
0	1	0	0	0	Paid
1	0	0	0	0	Cancel

"Dichte Codierung"

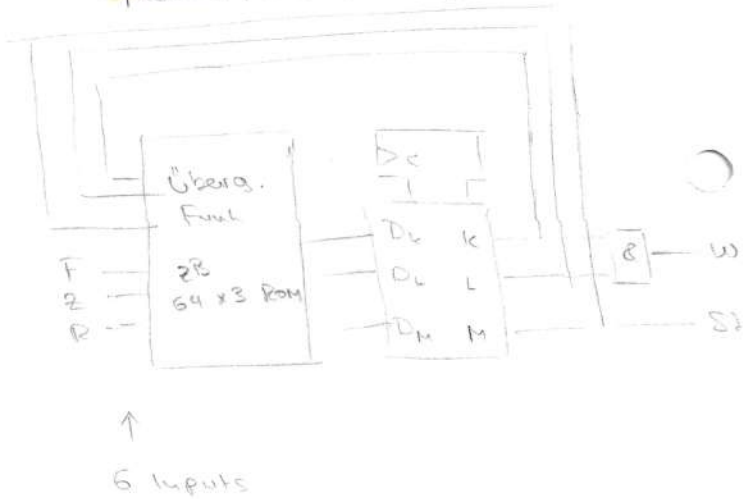
K	L	M	Zustand
0	0	0	Idle
1	0	0	Five
0	1	0	Ten
1	1	0	Paid
0	0	1	Cancel

Moore: 1 aus n cod.

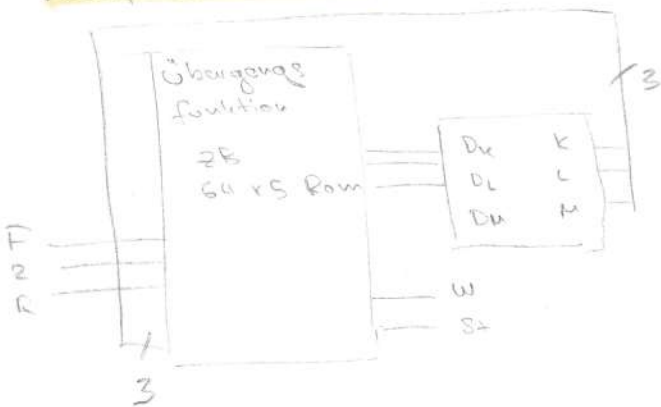


aus Effizienzgründen diese Anordnung

Moore: Dichte Cod.

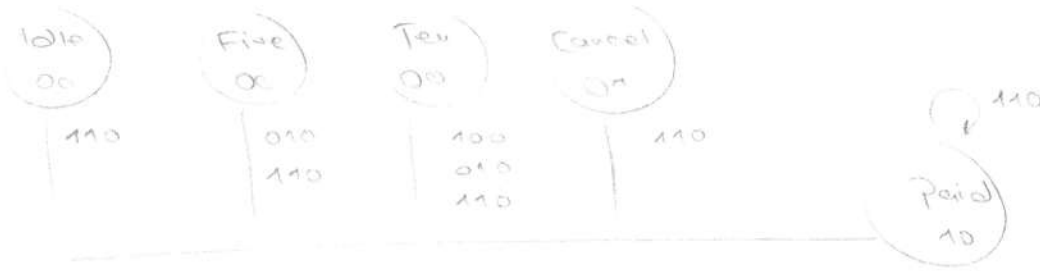


Mealy: Dichte Cod.



# Realisierung der Übertragungsfunktion

Funktionspeicher (Einzelne Gatter oder PLA, PAL)



Tabellenspeicher (ROM)

Angenommen: Moore und Deusse Encodierung

$$D_{paid} = (Idle \wedge F \wedge Z \wedge R) \vee (Five \wedge Z \wedge R) \vee \dots$$

$2^6 = 64$ , deshalb  $64 \times 3$  ROM



## Zeitlicher Ablauf

Stare Folien für synchrone Signale (Eingänge & Zustand)

Auch möglich: asynchr. Eingang

## Weiteres Beispiel:

Bitzähler (Zähler mit jedem Takt) → keine Eingabe; Autonomer Automat

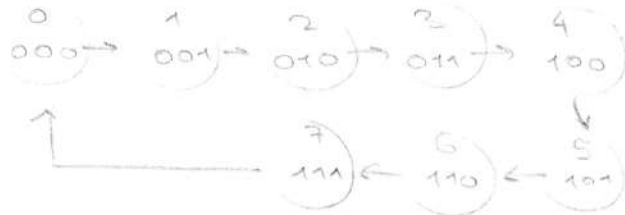
Zustandsspeicher: JK (Jump-Kill) - Flipflop

Übergangsfunktion: mit Gatter

Zustandsdiagramm:

(wird anhand bei 7)

C	B	A
MSB		LSB



## Übergang und Ausgangsfunktion

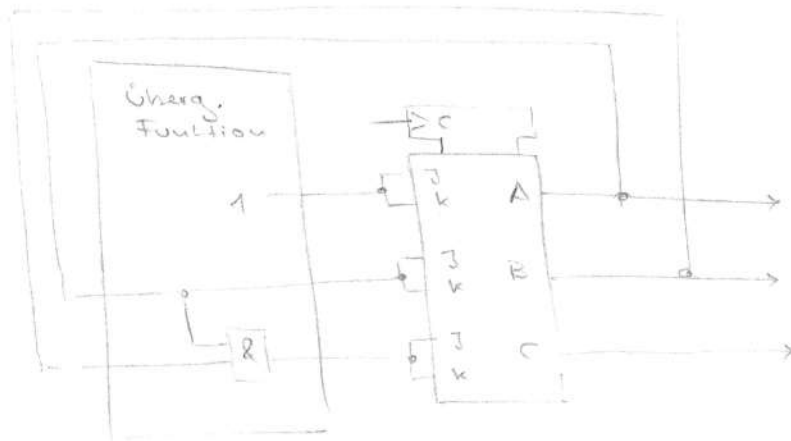
- Ausgangsfunktion Zustand = JK-Flipflop-Ausg.
- JK-Flipflop zum kontrollierten Lenken

JK	
00	hold
01	0
10	1
11	toggle

bei  $J=K=0$  hold  
 $J=K=1$  toggle

Dichte Codierung

Zustand	
0	...000
1	...001
2	...010
3	...011
...	...
7	...111



## Berechnung d. Maximaler Taktfrequenz

- pro Gatter 15ns
- Durchlauf Takt bis Ausgang 45ns
- Vorbereitungszeit 5ns
- Haltezeit 2ns
- max Taktfrequenz JK-Flipflop 25 MHz

Berechnung:

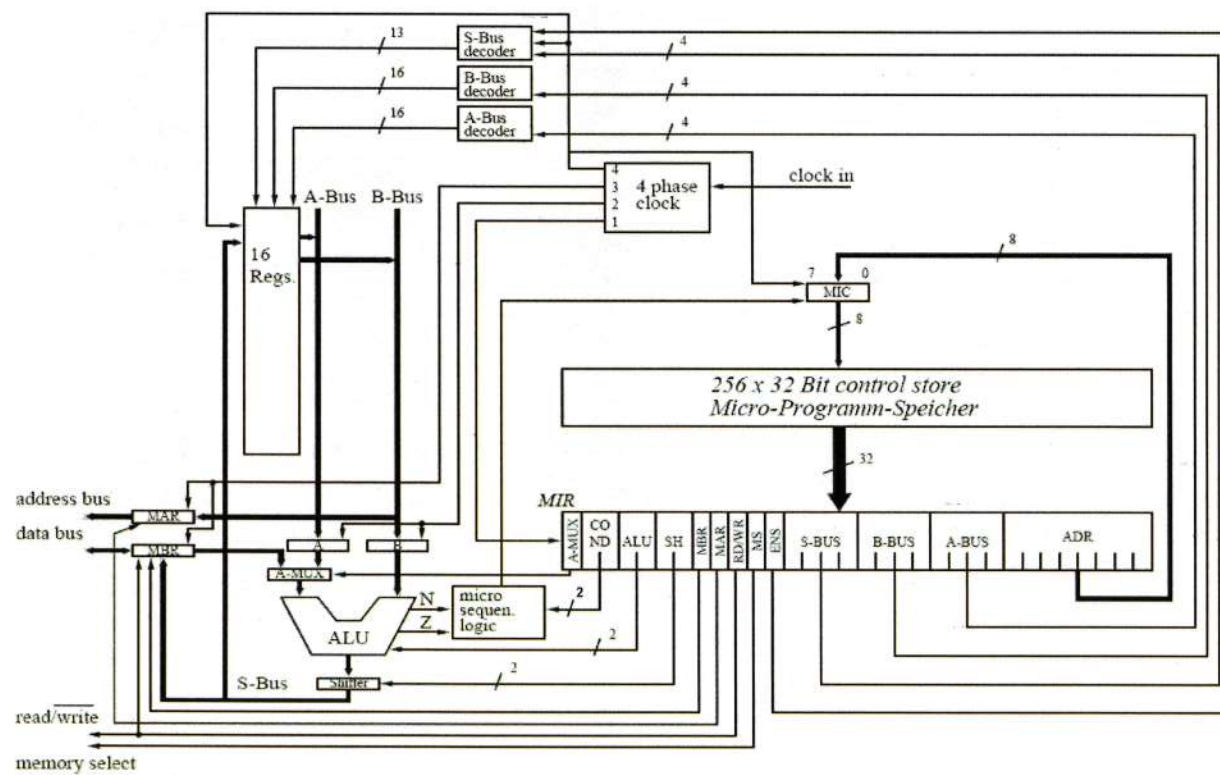
Durchl. Flipflop:	45ns
Durchl. Überg.:	15ns
	(nur 1 Gatter)
Vorbereitung flip-flop:	5ns
<hr/>	
	65ns

$$\frac{1}{65\text{ns}} = 15,4 \text{ MHz}$$

15,4 MHz < 25 MHz  
 nicht beschränkt



# Übersicht



# 14: Micro 16 (Mikroprocessor)

Prozessor Schaltwerk zur Verarbeitung von Daten  
besteht aus Rechenwerk und Leitwerk

Rechenwerk Führt Rechenoperationen aus

Leitwerke Steuerung der Reihenfolge mit der Befehle ausgeführt werden

## ALU Arithmetic Logic Unit

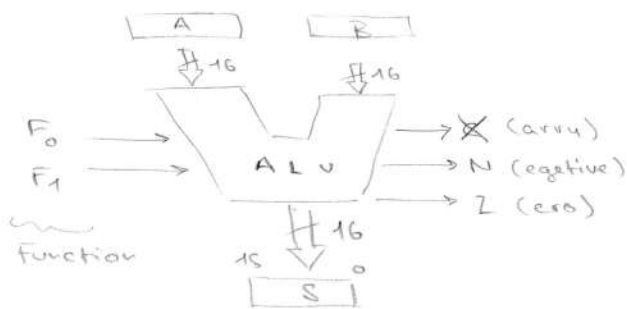
Funktionen:



- Durchschalten ohne Änderung
- Bitweise UND
- Bitweise Komplementbildung
- Parallele Addition zweier Datenwörter → Carry bit am Ende ausgelassen (nicht returned)

Es gibt viele Wege zu addieren!

- Carry ripple adder
- Carry skip adder
- Carry select adder
- Carry look ahead adder

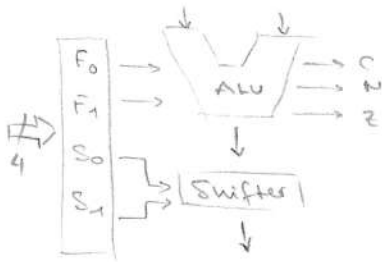


Micro Instructionen:

$F_0$	$F_1$		Symbolisch
0	0	unverändert A	$R \leftarrow A$
0	1	$A + B$	$R \leftarrow A + B$
1	0	$A \wedge B$	$R \leftarrow A \wedge B$
1	1	$\neg A$	$R \leftarrow \neg A$

# Shifter

Folgt nach ALU



S<sub>0</sub> S<sub>1</sub>

0	0	keine Änderung
0	1	shift left
1	0	shift right
1	1	-

Bezieht sich auf Output

SH ← R
SH ← lsh(R)
SH ← rsh(R)
-

# Register File und Steuerleitungen

Multiplexer schaltet Inhalt von beliebigem Register zu A/B

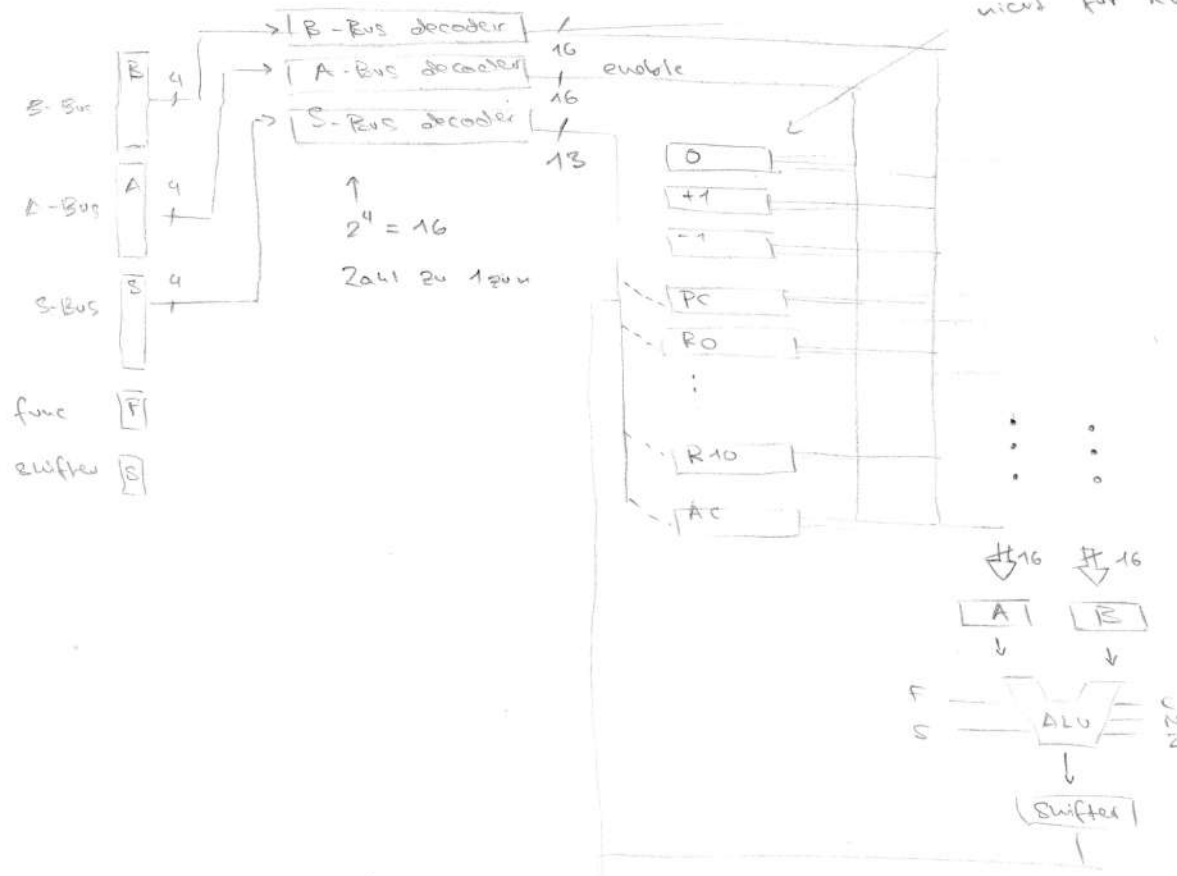
Bus parallele Leitungen von  $N$  Registern zu A/B

A-Bus : nach A

B-Bus : nach B

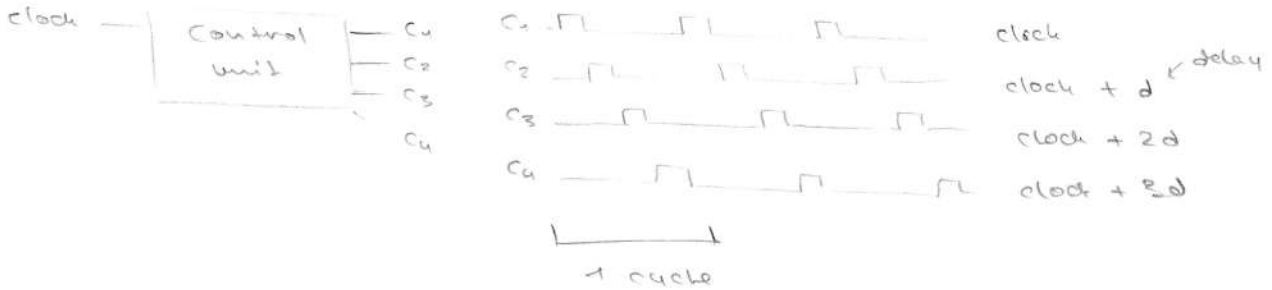
S-Bus : von Shifter zurück

nicht für Rückgabe! (0, +1, -1)



$2^4 = 16$   
Zahl zu 16-Bit

# Control Unit



Ablauf:

C<sub>1</sub>: Microinstruktionen laden

  Daten auf A/B-Bus legen

  Auswahl ALU-Funktion, Shift

C<sub>2</sub>: Register A/B mit Daten laden

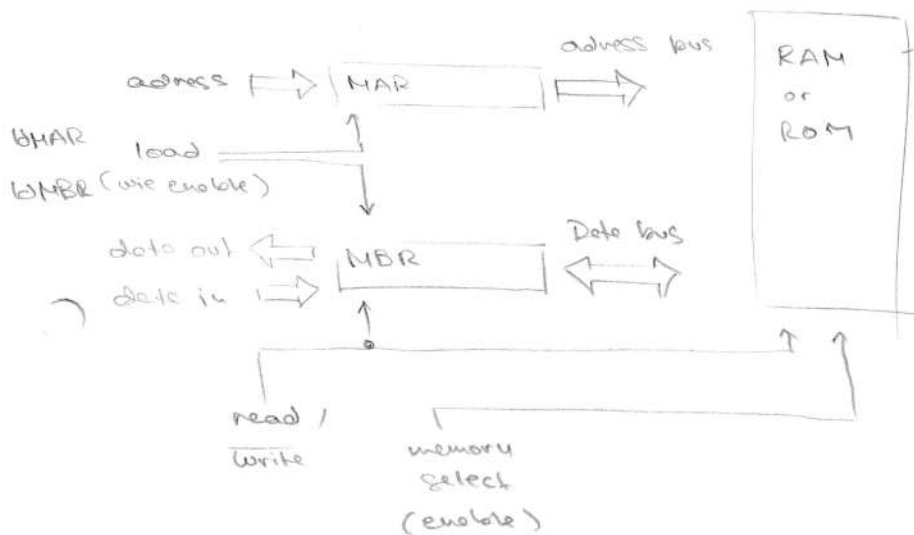
C<sub>3</sub>: ALU und Shifter Zeit geben

C<sub>4</sub>: Ergebnis von S-Bus in das Zielregister laden

## Speicheranbindung

MAR Memory Address register

MBR Memory Buffer register



## Speicherzugriff

- Handshake Verfahren (zu kompliziert)
- Speicherzugriffszeit abwarten (genau 2 Taktzyklen)

1. Adresse in MAR laden

2a) schreiben: MBR befüllen

2b) lesen: MBR inhalt in Register speichern

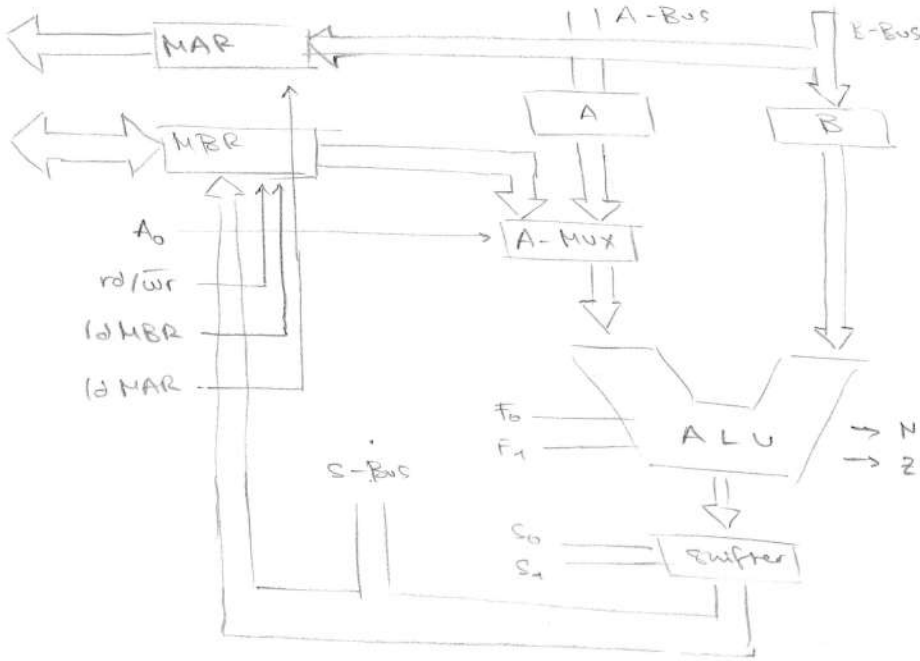
2a) schreiben: load für MAR und MBR auf 1 stellen

read/write auf 0

2b) lesen: load nur für MAR auf 1

read/write auf 1

# ALU Schnittstelle zu RAM/ROM

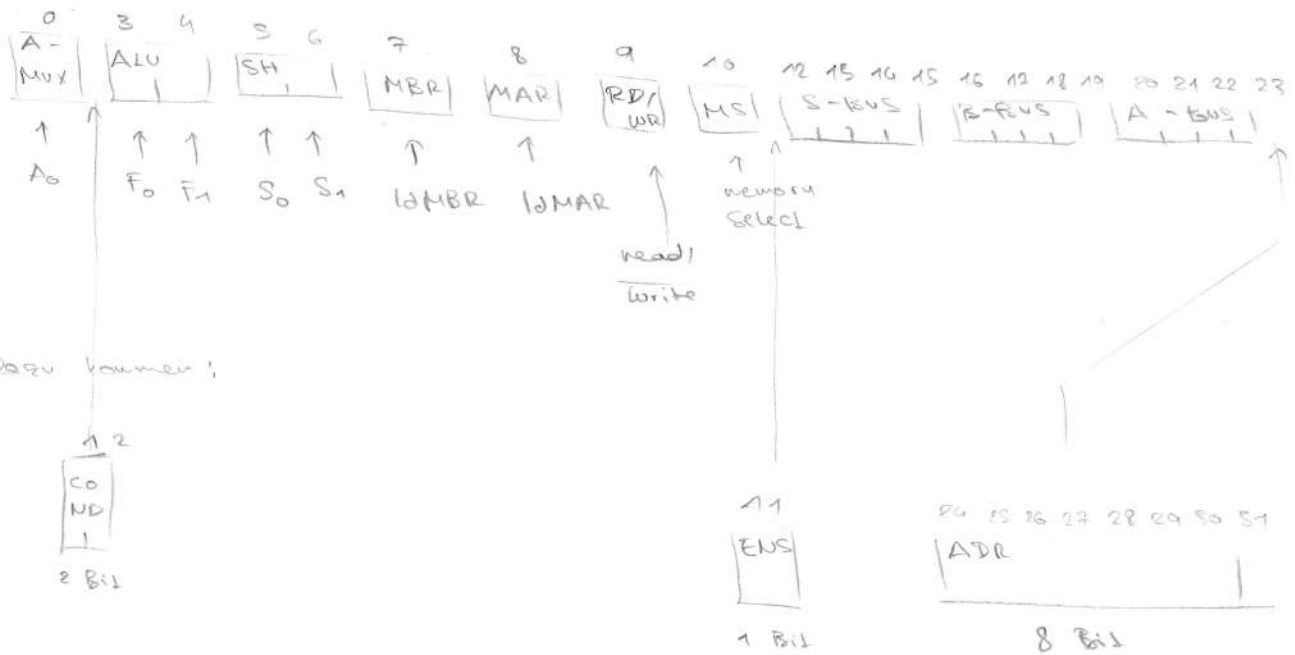


Man sieht:

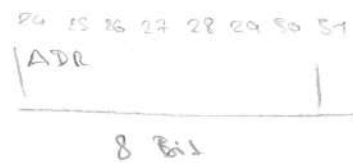
- MAR wird immer mit B-BUS Eingang gefüllt.
- A-Mux entscheidet ob MBR output oder A-BUS nach ALU geschickt wird
- S-BUS schickt immer nach MBR

## Micro Instruction Register

(Wir lassen Carry Bit von ALU aus, damit wir genau 32 Bit haben)

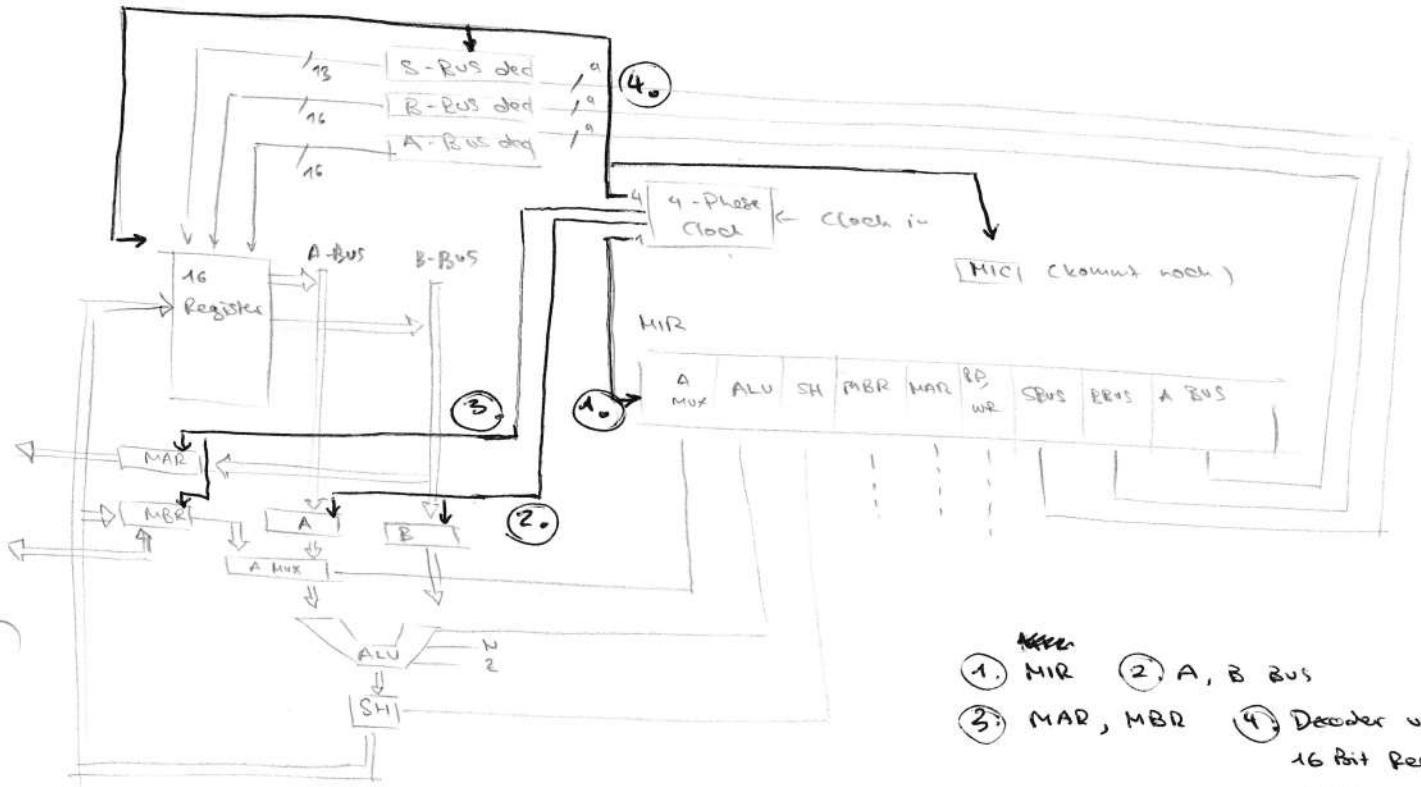


Dazu kommen:



$$|[0; 31]| = 32$$

# Vereinfachte Micro 16 Architektur (bisher - unvollständig)



- ①. ~~MIR~~ MIR
- ②. A, B Bus
- ③. MAR, MBR
- ④. Decoder und 16 Bit Register, MIC

Dazu kommen:

256 x 32 Bit control Store  
 Mikroprogramm Speicher } Micro-Code ROM

ADR

8 Bit  $\rightarrow 2^8 = 256$  Adressen des Speichers

ENS

Enable S-Bus-Decoder

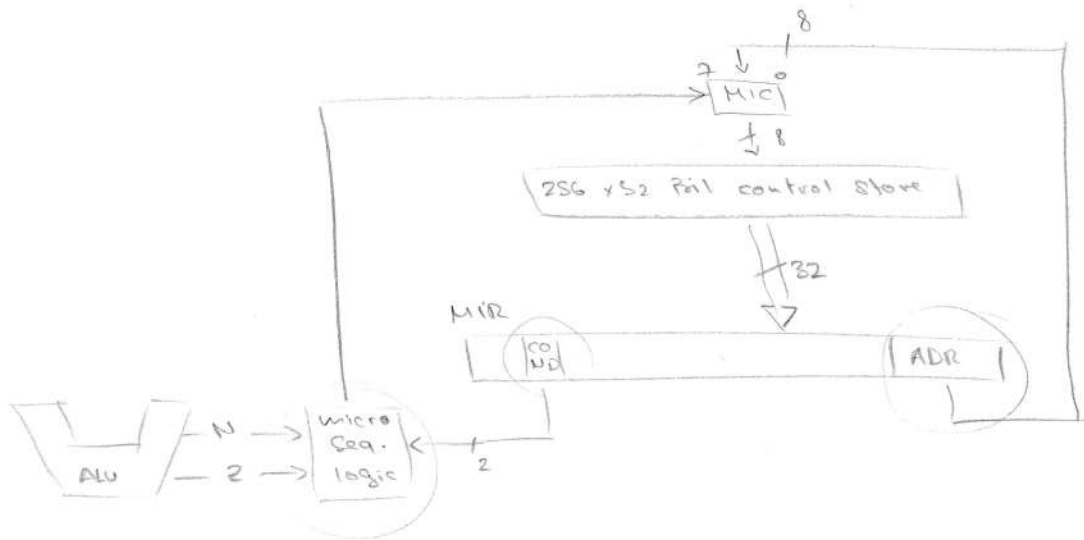
COND

2 Bit  $\rightarrow$  Conditions für micro sequencing logic

00	kein Sprung	-
01	Sprung zu ADR in ROM wenn N=1	if N goto ADDR
10	Sprung zu ADR in ROM wenn Z=1	if Z goto ADDR
11	Sprung zu ADR in ROM	goto ADDR

MIC - Micro Instruction controller

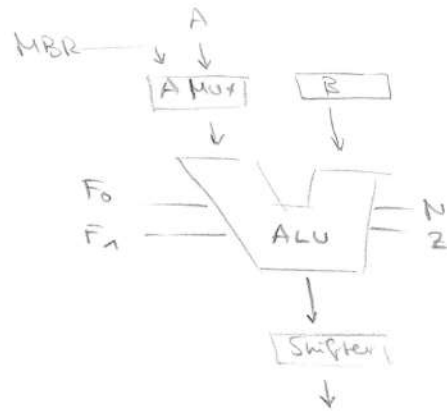
Gibt Input an Micro Code ROM



MIC wird immer mit einer ADR geladen,  
 micro seq. logic bestimmt aber auf Grund von N, Z  
 ob zu dieser Adresse gegengru werden soll

## Micro instructions für die ALU

$F_0 F_1$		
00	A durchschalten	$R \leftarrow A$
01	$A+B$	$R \leftarrow A+B$
10	$A \wedge B$ Bitweise	$R \leftarrow A \wedge B$
11	$\neg A$	$R \leftarrow \neg A$



## Shifter

$S_0 S_1$		
00	keine Änderung	$SH \leftarrow R$
01	shift left	$SH \leftarrow lsh(R)$
10	shift right	$SH \leftarrow rsh(R)$
11	ungültig	—

## Speicheranbindung

Lesen:

$MAR \leftarrow R_i ; rd$   
 $rd$   
 $R_k \leftarrow MBR$

Schreiben:

$MAR \leftarrow R_i ;$   
 $MBR \leftarrow R_j ; wr$   
 $wr$

MAR: Memory address register  
 MBR: Memory buffer register

Zum lesen:

$read / \overline{write} = 1$   
 memory select = 1  
 $ADR \Rightarrow MAR$  und  $load\ MAR = 1$   
 Bits von MBR holen

Zum schreiben:

$read / \overline{write} = 0$   
 memory select = 1  
 $ADR \Rightarrow MAR$  und  $load\ MAR = 1$   
 Daten  $\Rightarrow MBR$  und  $load\ MBR = 1$

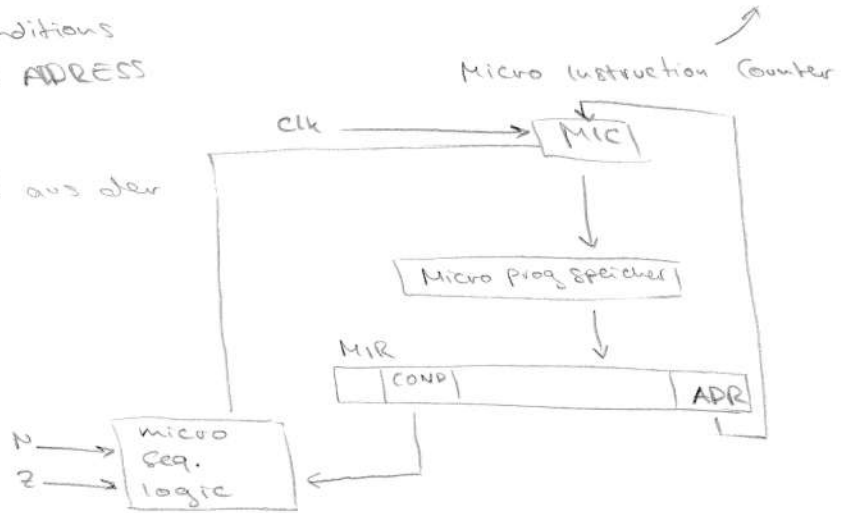


# Micro Sequencing Logic - Conditions

Entscheidet basierend auf conditions und N, Z ob MIC wirklich zu ADDRESS gehen soll oder nicht

(führt ausserdem die command aus der nächsten Zeile aus)

mit clock immer +1



COND

00

01

10

11

- [if N goto ADDR]
- [if Z goto ADDR]
- goto ADDR

Als Bitfolge : 32 Bits

(memory select)

A-Mux	COND	ALU	SH	MBR	MAR	RD/WR	MS	ENS	S-Bus ADR	B-Bus ADR	A-Bus ADR	ADR
1	2	2	2	1	1	1	1	↑ 1	4	4	4	8

Enable S-Bus Decoder

Lesen-Zugriff : ENS = 0

Schreib-Zugriff : ENS = 1

## 15 Prozessorarchitekturen, Befehlssatz

Prozessoren haben eigene Architekturen  
von diesen leiten sich Befehlssätze ab

Mit Befehlssätzen schreibt man Programme (Mikroprogramme → in Mikro 16: ROM)

Ziel:

wichtigere erweiterbare Befehlssätze (z.B. Multiplikation) → Mehr Maschineninstruktionen  
die mehr können als  
Prozessor

Gen 5: Very high level language VHLL

Gen 4: Fourth gen language (anwendungsbezogen)

Gen 3: Höhere Programmiersprachen

Gen 2: Assemblerbefehle

Gen 1: Binäre Instruktionen

} compiler  
} interpreter

Abhängig von Prozessorarchitektur

## Maschinenbefehle

Aritmetische Befehle

Addition, Subtraktion, Multiplikation, Division, Inkrement, Dekrement

Arithmetische Shift

Logische Befehle

AND, OR, XOR, Invert

Shiftbefehle

Shift, rotate

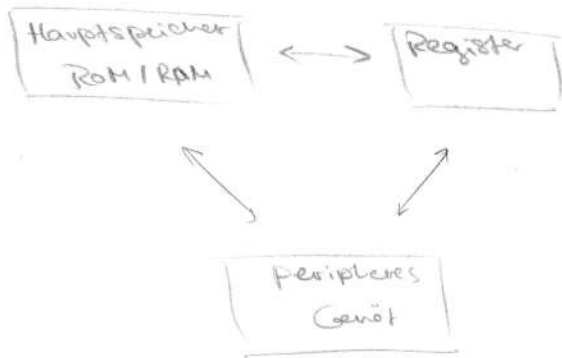
Datentransferbefehle

Mov, I/O Befehle → Daten aus Peripherie durch Ports, Stackmanipulation

(Flow Control) Kontrollbefehle

Funktionsaufrufe, Sprünge

## Befehle für Datentransfer



- MOV
- Datenwert kopieren, zB String und Umadressieren
- PIO (Programmable I/O)  
Peripheres wie I/O behandelt
- DMA (direct memory access)  
Eigener Controller verwaltet I/O Zugriffe  
CPU setzt nur DMA-Kommando ab  
Interrupt nach Beendigung

## Adressierung

- Register Mode  $R1 \leftarrow R2$
- Immediate Mode  $R1 \leftarrow 1029$
- Direct addressing Mode: goto 128  
 $RB \leftarrow \text{memory} [(500)_{16}]$
- Register-indirect Mode:  $RB \leftarrow \text{memory} [RS]$   
← Pointer  
oder auch Pointer und Displacement  
 $[(500)_{16} + RS]$   
↑  
ADD (RS) +
- Programm-Counter-Relative-Addressing-Mode
- Indirect Addressing Mode  
 $RB \leftarrow \text{memory} [\text{memory} [(500)_{16}]]$

## Befehle für Datentransfer: Zwischen Register und Stack

Für Stack wird RAM benutzt und mit Pointer auf Stack head gerichtet  
(Alte Elemente nicht gelöscht)

push-16 (R1) { 1. MAR  $\leftarrow$  SP (Stackpointer)  
2. MBR  $\leftarrow$  R1; wr  
3. wr  
SP  $\leftarrow$  SP + (-1)

push-16 (R3) ← speichert Ereignis in R3

## Control flow Befehle

Man unterscheidet zwischen bedingten & unbedingten Sprüngen



branch if zero flag (boolean)

branch if equal

⋮

- Funktion / Prozedur / Unterprogramm

Fasst code zu einer Einheit mit Input und output zusammen

Zu merken:

- Wo ist Funktion
- wohin weiter nach Funktion
- Argumente
- Rückgabewert (an Stack oder Register)

Vorgang:

Call Subroutine

PC = Return  
Adress

1. Registerinhalte und PC (Program Counter Register) speichern (z.B. an Stack)
2. PC ← Startadresse der Funktion
3. Abarbeiten der Funktion

Return from Subroutine

1. Registerinhalte & PC wiederherstellen
2. Mit PC wert fortfahren
3. Ablauf fortsetzen bis Programmende

Ein Stack kann dazu dienen lokale Variablen und Aufrufhierarchie zu speichern,

Gefahr: Stackoverflow / Stackunderflow (<0)

Ermöglicht rekursive Aufrufe,

# Interrupts

asynchrone Unterbrechung des Programmablaufs  
ausgelöst durch spez. Hardware: Interrupt Controller

Auslöser: INTERN (Software) → Exception

- Division by Zero
- Illegaler Befehl
- Stack overflow

EXTERN (Hardware) → Error

- Keyboard Taste gedrückt  
(Damit Prozessor nicht ständig warten muss auf keyboard)
- Speicherzugriff von DMA beendet

Eigene Routinen für verschiedene Interruptions  
Interrupt Service Routine - ISR  
für Fehler - Behandlung

Fehler (Error / Exception) selbst erkennen oder durch codierten Interrupt  
via Parallelbus error-Type erhalten

Wichtig: Interrupt-Adresse ist immer an der selben Stelle und wird mit  
Interrupt-Vektor bestimmt.

- Prozesse müssen möglichst kurz sein
  - Interrupts haben Prioritäten (falls sie von einem anderen Interrupt unterbrochen werden)
  - Sie können auch aufgeschoben / verzögert werden in spezifischer wichtiger code Teilen
- ↳ Dann: flag setzen → Interrupt pending als Status im Statusregister.

## 2 Wichtige Architekturen:

RISC Reduced instruction set CPU (Processor) → Bessere Taktfrequenz

- Schnelle Zugriffe auf Speicher um Routinen auszuführen, wenige Befehlsblöcke
- Weniger Adressierungarten
- Viele Register
- Speicherhierarchien

→ ARM (Advanced RISC Machines)  
MIPS

CISC Complex instruction set CPU

- Gegenteil von RISC → langsamer und nicht alle Befehle oft benötigt
- Unterschiedl. Datentypen und Adressierungarten
- Befehle haben unterschiedliche Ausführungszeiten
- Steuerung durch Mikroprogramm

↳ IBM  
Intel

# 16 Pipelining

## MIPS Microprocessor without interlocked pipeline stages

Wiederholung:

- Maschinengebiete zur Erweiterung der Prozessorfunktionalität (mehr als Mikroinstruktionen)

Beispiel: Stack, dient um Rücksprungsadressen bei Funktionsaufrufen zu speichern

- Interrupts  
Synchron oder Asynchron

## Speicher-Adressierung

Direkt (immediate)

$R1 \leftarrow \text{memory}[0x\text{CAFE}]$

Indirekt

$R2 \leftarrow 0x\text{CAFE}$




$R1 \leftarrow \text{memory}[R2]$


Prä und Post Inkrement / Dekrement

$R1 \leftarrow \text{memory}[(R2) +]$

↑  
egal ob Prä oder post; Danach immer gleich

## Taktung bisher in Micro16

1.  (MIR)
2.  (REG)
3.  (Adr)
4.  (WB)

 "4 Phasen Redundanz"  
← 2ns →

Je nach Befehl kann eine der 4 Phasen redundant sein:

- 4 Phasen / Takte pro Befehl
- 2ns Taktung
- 8ns pro Befehl
- 125 Mio Instr. pro Sek (MIPS)

Unnötige Phasen lassen sich nicht überspringen

## Segmentierung von MIPS Datenpfad

IF	Instruction Fetch
ID	Instruction Decode and source register read
EX	Execute
MEM	Memory access
WB	Write back

Dadurch: 5 stufige Pipeline, 2ns Taktung

inklusive dem Einlaufen brauchen 6 Befehle (ohne Hazards) 10 Takte.

2ns-Taktung

Absolute Ausführungszeit

6 Befehle in 10 Takten je 2ns  $\rightarrow$  20ns

Realer Durchsatz

$$6 \text{ Befehle, } 20 \text{ ns} \rightarrow \frac{6}{20 \cdot 10^{-9}} = 300 \text{ MIPS}$$

oder 3,33 ns / Befehl

Theoretischer Durchsatz

$$1 \text{ Befehl, } 2 \text{ ns} \rightarrow \frac{1}{2 \cdot 10^{-9}} = 500 \text{ MIPS}$$

(kein Einlaufen)

## Performance Verbesserung:

- Ausführungszeit für einzelne Befehle im Datenpfad bleibt gleich:  
System selbst wird nicht in Verarbeitung schneller
- Durchsatz verbessert sich im best case k-fach

## Bottlenecks / Begrenzungen:

- Balance der Pipeline Stufen  
sie sollten alle möglichst gleich lang brauchen
- Einlaufen der Pipeline
- Hazards
  - Datenabhängigkeiten
  - Kontrollstrukturabhängigkeiten
  - Strukturelle Abhängigkeiten (nicht in MIPS - keine Architekturfehler)

# 16 Pipelining : Parallelverarbeitung zur Beschleunigung

Bisher aus Micro-16 Control Unit:

1. (MIR) Befehl vom MIR fetchen
2. (REG) Register AB für ALU enablen
3. (Adr) MAR und MBR enablen
4. (WB) S-Bus Decoder, Scratchpad / Register Files und <sup>micro instr. counter</sup> MIC für nächste Mikroinstruktion aktivieren.

Angenommen wir haben das Micro16-Programm:

Phasen die beansprucht werden:	①	②	③	④	
$R6 \leftarrow R6 + R5$	✓	✓	X	✓	(braucht MAR und MBR nicht)
$MAR \leftarrow R4; rd$	✓	✓	✓	X	(schreibt nichts im Speicher)
$rd$	✓	X	X	X	(wartet nur)
$R2 \leftarrow 1ch(R1)$	✓	✓	X	✓	(braucht MAR und MBR nicht)

Bei Micro-16:

4 Phasen Rechenwerk mit 4 Takten pro Befehl (4 ns) → 8ns pro Befehl  
 bei 2ns Taktung (4-2ns)  
 > 125 Mio Instruct. p. sec (MIPS)  
 Flops = floating point operation

## Ermöglichung der Parallelisierung durch neue Architektur:

MIPS: Microprocessor without interlocked pipeline stages

- Eigene Instruktionen: MIPS Assembly language
- 5 Stages:

1. IF Instruction fetch
2. ID Instruction decode
3. EX Execute / Address Calculation
4. MEM Memory Access
5. WB Write Back





## Befehle aus Vorlesung:

lw \$10, 20(\$1)

$R10 \leftarrow \text{Memory} [R1 + 20]$

sub \$11, \$2, \$3

$R11 \leftarrow R2 - R3$

lw \$1, 100(\$0)

$R1 \leftarrow \text{Memory} [R0 + 100]$

→ siehe MIPS Assembler Instructions

Aber nicht alle Befehle beanspruchen alle 5 Pipeline stages gleich stark!

Befehl	IF	ID	EXE	MEM	WB	$\Sigma$
lw	2ns	1ns	2ns	2ns	1ns	8ns
sw	2ns	1ns	2ns	2ns		7ns
[add, sub, and, or slt]	2ns	1ns	2ns		1ns	6ns
beq	2ns	1ns	2ns			5ns

↓  
(branch if equal)

Sowohl  $< 2ns$ , beanspruchen Sie jedoch zur Vereinfachung 2 ganze ns.



## Durchsatz:

5 Stufig, 2ns Taktung  
Ein Befehl 10ns  
insgesamt 10 Takte

Real 3,33 ns pro Befehl

Theoretisch 2ns pro Befehl

↓  
Wichtig: mehr Stufen  $\neq$  mehr performance!

(sinkt ab  $k=15$ )

„Super pipelining“

## Bottlenecks:

- Auffüllen der Pipelines „Einlaufen“
- Balance der Stufen
- Hazards (Abhängigkeiten)

## Sind mehr pipeline Stufen besser?

Intel Pentium: doppelte # Stufen → doppelter theoretischer Durchsatz durch Verdopplung der Taktfrequenz/Rate  
(jede Stufe in 2 aufteilen)

## Problem:

deutlich mehr Hazards

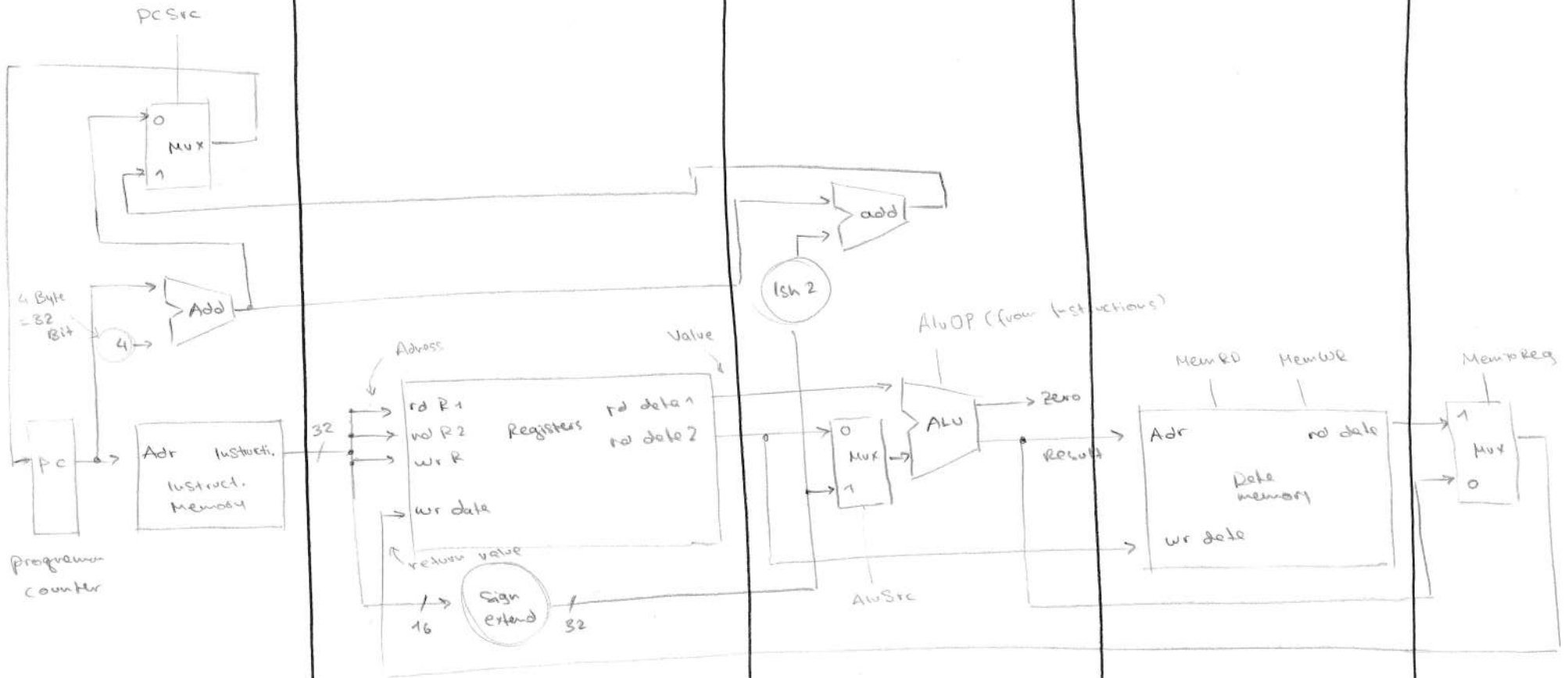
Instruction Fetch

Instruction Decode and read registers

Execute, Address Calculation

Memory Access

Write Back Registers



(Bitte Risc-V)

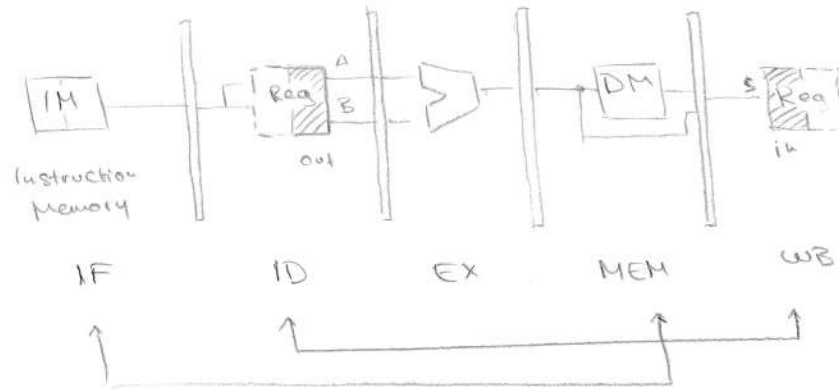
IF / ID

ID / EX

EX / MEM

MEM / WB

Erhöhte Abstraktion:



Dependencies:

- Gleicher Speicher für Write-Back und Instruction Decode (Annahme: in der ersten Takthälfte schreiben, dann lesen)
- Gleicher Speicher für Instructionen und Daten

# Hazards

## Strukturelle Hazards

Pipeline Stufe nicht für alle Befehle verfügbar

→ Lösung: Architektur Neuentwurf, Stall

## Control Hazards

Nachfolgebefehl hängt von Verzweigung ab

→ Lösung: prediction, delayed branching (code umändern), stall

## Data Hazards

Berechnung abhängig von Vorgänger-Befehl, falsche Ausführungsreihenfolge

→ Lösung: forwarding, Code-Optimierung, stall

RAW - Read after Write  
gelesen bevor richtig beschrieben

$$\begin{array}{l} R2 \leftarrow R1 + R3 \quad \leftarrow 2. \\ R4 \leftarrow R2 + R1 \quad \leftarrow 1. \end{array}$$

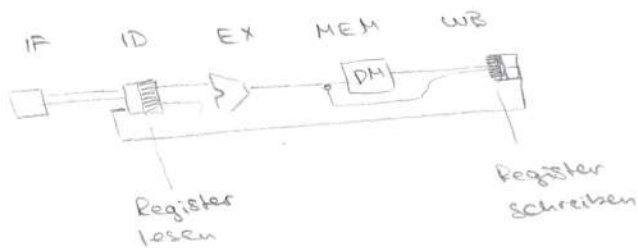
WAR - Write after Read  
(über)schrieben bevor richtig gelesen

$$\begin{array}{l} R2 \leftarrow R1 + R3 \quad \leftarrow 2. \\ R3 \leftarrow 1 \quad \leftarrow 1. \end{array}$$

WAW - Write after Write  
vom vorherigen Befehl falsch überschrieben

$$\begin{array}{l} R2 \leftarrow R1 + R3 \quad \leftarrow 2. \\ R2 \leftarrow 1 \quad \leftarrow 1. \end{array}$$

Die Data Hazards entstehen wegen:



Gemeinsamer Registersatz:  
1. Takt Hälfte Schreiben  
2. Takt Hälfte Lesen

Beispiel: Stall bei RAW - Read after Write

write  $R2 \leftarrow R1 + R3$   
 read  $R4 \leftarrow R2 + R1$  } Problem: Read before write

Takt	IF	ID	EX	MEM	WB
1	$R2 \leftarrow R1 + R3$				
2	$R4 \leftarrow R2 + R1$	$R2 \leftarrow R1 + R3$			
3	○	$R4 \leftarrow R2 + R1$	$R2 \leftarrow R1 + R3$		
4	○	$R4 \leftarrow R2 + R1$		$R2 \leftarrow R1 + R3$	
5	○	$R4 \leftarrow R2 + R1$			$R2 \leftarrow R1 + R3$ ✓ ausgeführt
6	□	○	$R4 \leftarrow R2 + R1$		
7	△	□	○	$R4 \leftarrow R2 + R1$	
8	◇	△	□	○	$R4 \leftarrow R2 + R1$

# Hazards

Strukturelle Hazards (Hardware Problem) → für MIPS irrelevant

Mehrere Stufen brauchen dieselben Ressourcen → Stalling, Architektur ändern

## Control Hazards

Nachfolgebefehl hängt vom Ausgang des Sprungs ab → prediction, delayed branch

## Data Hazards (am wichtigsten)

Berechnung benötigt Ergebnis des Vorgängers → forwarding, cache-optimierung, Stalling

RAW: read after write

1.  $R2 \leftarrow R1 + R3$  schreiben
2.  $R4 \leftarrow R2 + R1$  lesen → falsch

lesen bevor Register befüllt wurde mit richtigem Wert

WAR: write after read

1.  $R2 \leftarrow R1 + R3$  lesen ← irrtümlich mit  $R3 \leftarrow 1$  ausgeführt
2.  $R3 \leftarrow 1$  schreiben

Schreiben bevor gelesen wurde

für MIPS irrelevant

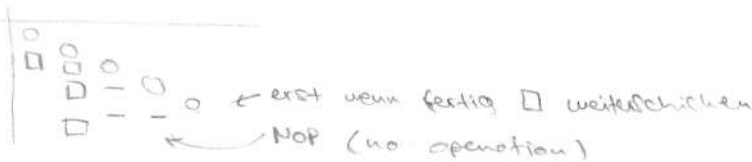
WAW: write after write

1.  $R2 \leftarrow R1 + 3$
2.  $R2 \leftarrow 1$

Vom Vorgänger überschreiben

## Stall / Pipeline stall

Man lässt bubbles entstehen → Entweder durch Architektur oder Compiler



## Delayed branching

Branch Befehle zurückschicken und in Zwischenzeit andere Befehle exekutieren

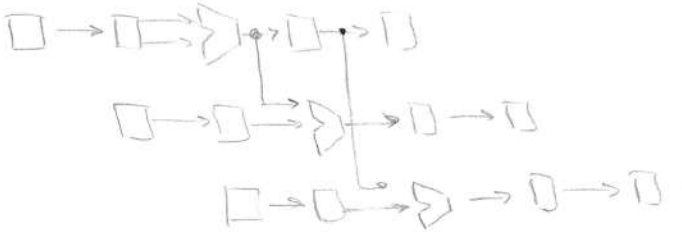
## Code Optimierung

Dort wo keine Abhängigkeiten, Delay erzeugen und andere Befehle währenddessen machen

# Data Forwarding

add \$t0, \$t1, \$t2  
sub \$t3, \$t0, \$t1  
or \$t4, \$t0, \$t2

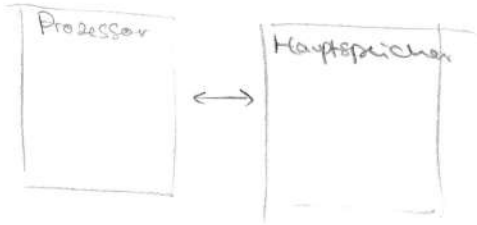
$R_0 \leftarrow R_1 + R_2$   
 $R_3 \leftarrow R_0 + R_1$   
 $R_4 \leftarrow R_0 + R_2$



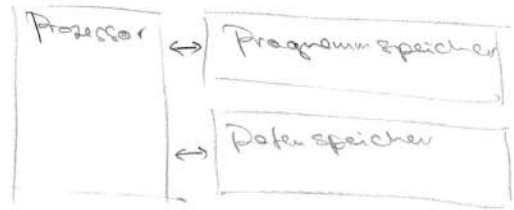
leitet Daten sofort weiter um Hazards zu vermeiden

# 17 Speichermanagement, Caches und Hierarchien

## Von-Neumann-Architektur



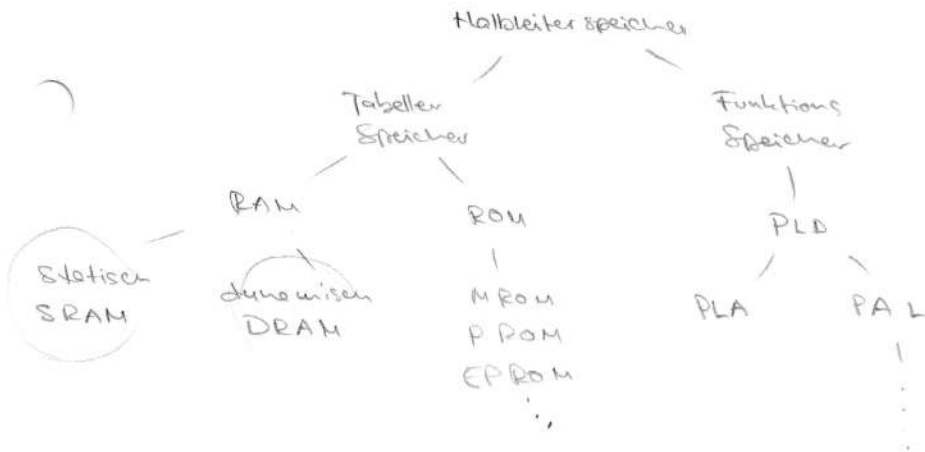
## Harvard Architektur



(Bei Micro 16:

Microinstruktionen - Programmspeicher  
 Datenspeicher als ROM / RAM mit MAR)

## Speicherbausteine (Wiederholung)



### SRAM

- Speichert in Latches
- + Sehr schnell (kurze Zugriffszeiten)
- sehr teuer

### DRAM

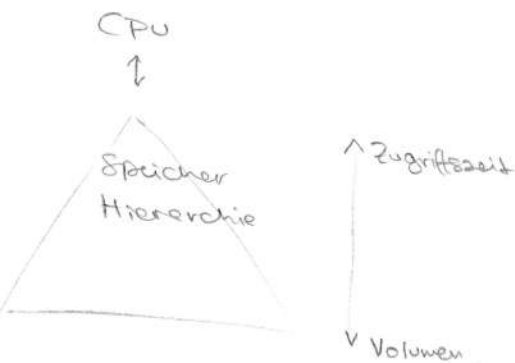
- Speichert in Kondensatoren
- hat refresh-cycles

↓  
 Ziel: Preis-Performance-Balance

## Beobachtung:

- Temporal locality - oft dieselben Speicherzugriffe hintereinander
- Spatial locality - oft benachbarte Speicherzugriffe hintereinander

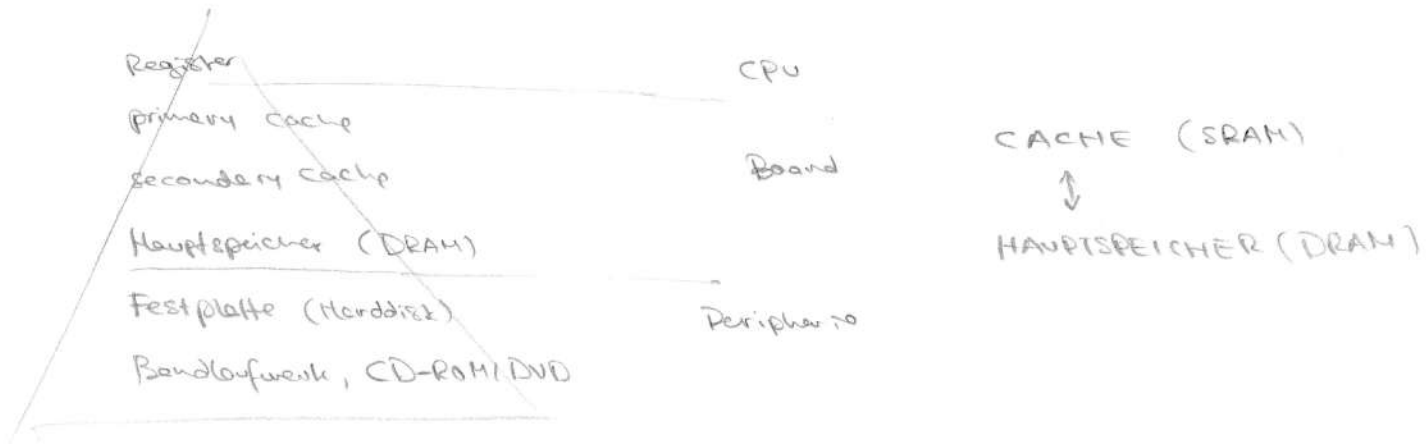
## Lösung



Dadurch erzeugt man gepufferte Speichercontroller

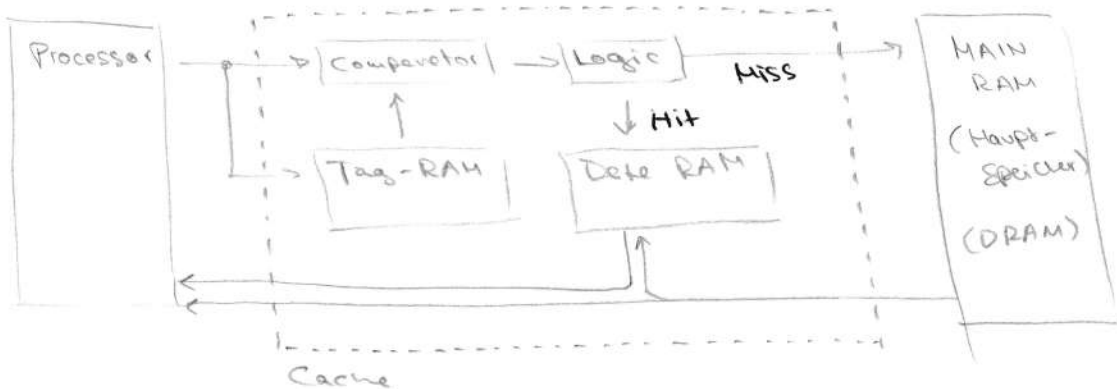


## Typische Hierarchie



## Cache

Schneller Zugriff der CPU auf Daten / Befehle



## Cache-Performance

$$t_{eff} = h \cdot t_{cache} + (1-h) \cdot t_{main}$$

$t_{eff}$  ... effektive Speicherzugriffszeit

$h$  ... Trefferquote (Hit-Rate)

$t_{cache}$  ... Zugriffszeit auf Cache

$t_{main}$  ... Zugriffszeit auf Main

$$t_{AMA} = t_{hit} + m \cdot t_{penalty}$$

$t_{AMA}$  ...  $\emptyset$  Zugriffszeit

$t_{hit}$  ... Hit time

$m$  ... Miss Rate  $(1-h)$

$t_{penalty}$  ... Miss penalty

↓  
 Refüllen bei Miss

1. Stell der Pipeline
2. Adresse an Speicher management
3. Daten von Main auf Cache

Wichtige Parameter:

Hit-Time

Hit-Rate → Miss-Rate =  $1 - \text{Hit Rate}$

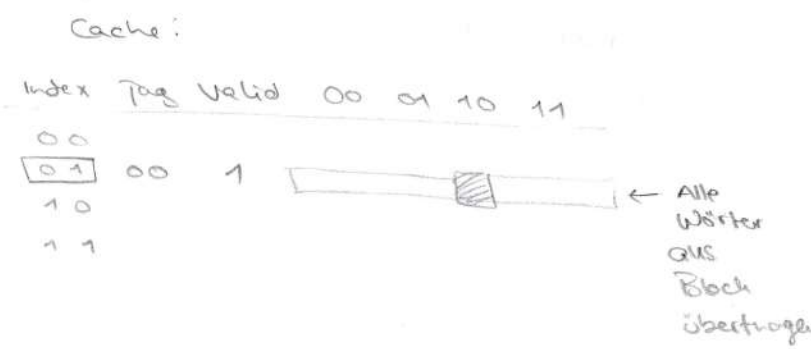
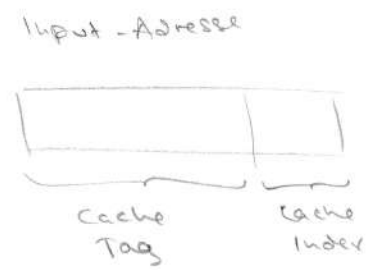
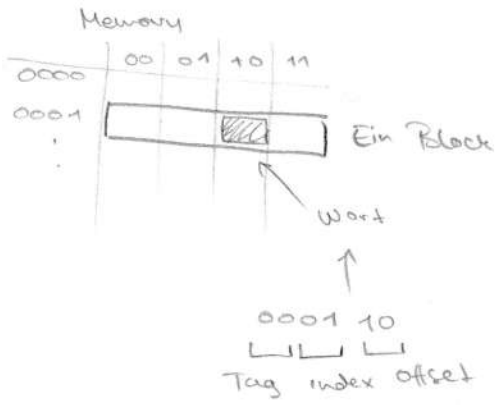
Miss-Penalty

Cache Performance wichtig, weil sonst bottle-neck für CPU

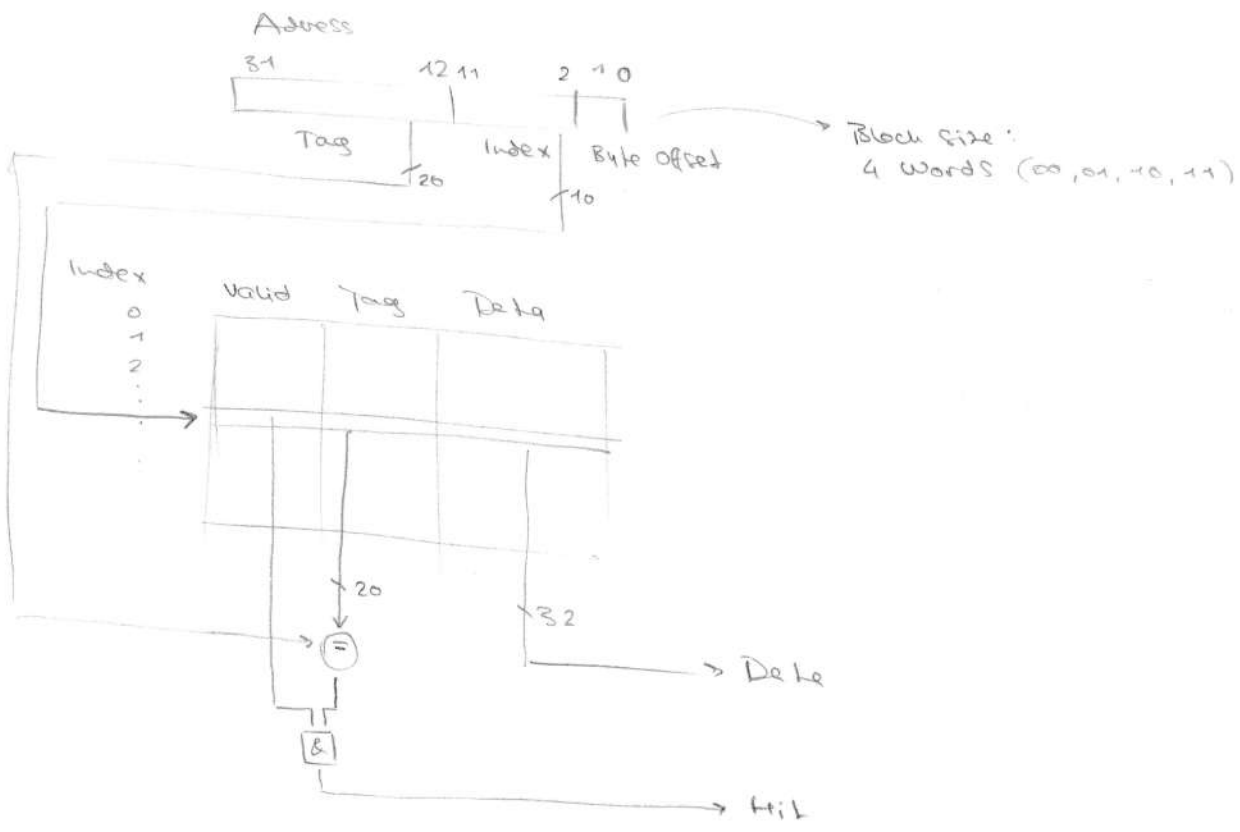
# Cache Verwaltungsstrategien

## Direct Mapped

- Ein Teil der Adresse bestimmt Cache Position (wie Modulo Rechnung)  
↳ niederwertige Bits: „Cache-Index“
- Restlichen Bits zur eindeutigen Identifikation als „Cache-Tag“
- „Valid-Bit“ zur Erkennung bisher unbewiteter Positionen
- oft ineffizient einzelnes Wort / Byte zu cachern, deshalb werden Blöcke gecached



Detaillierter:



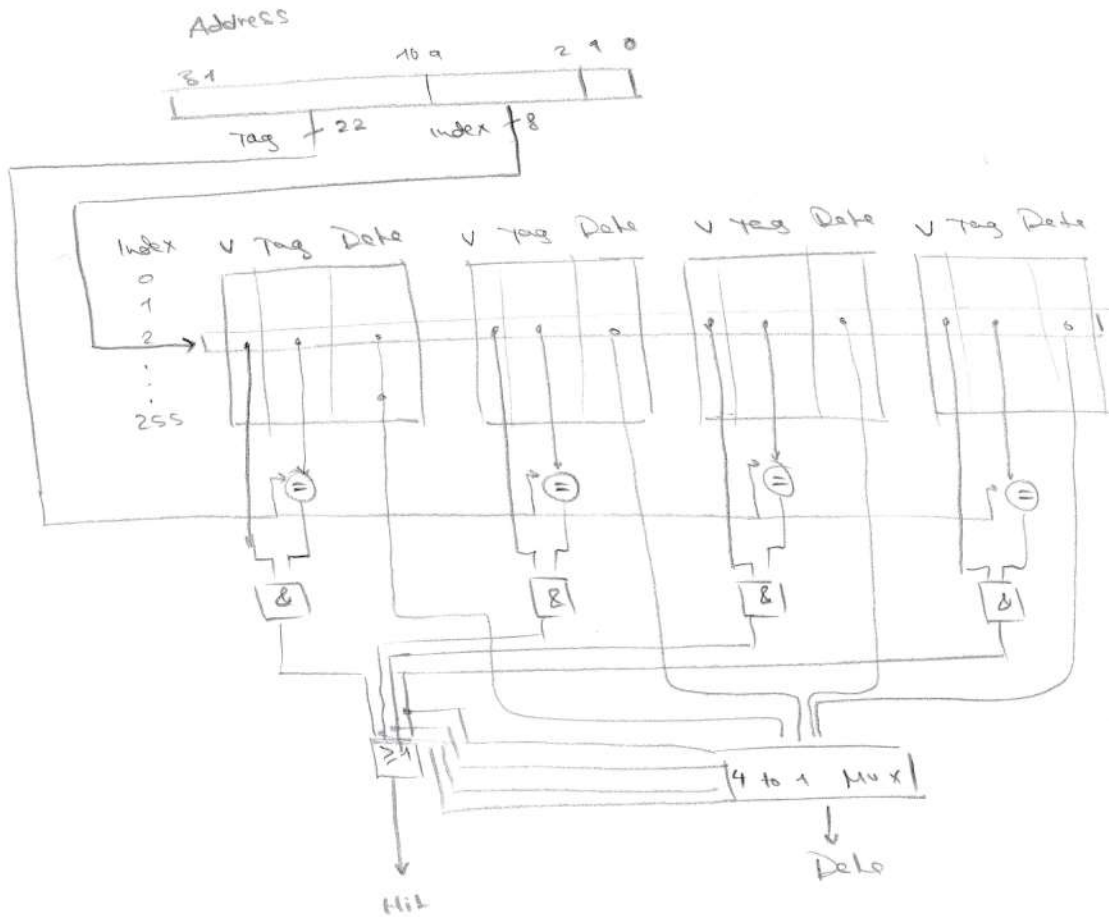
## Fully associative Cache

Jeder Block darf in jeder beliebigen Position im Cache gelegt werden  
(kompliziertester Aufbau)

kein Index  $\rightarrow$  volle Adresse als Tag

## N-way Set associative Cache

Jeder Block bekommt einen Set zugewiesen der fully associative ist



$$\# \text{ Blöcke} = \# \text{ Sets} \cdot \underbrace{\# \text{ ways}}_{\text{index ways}}$$

# Cache Verwaltungsstrategien im Vergleich

## Direct-mapped cache

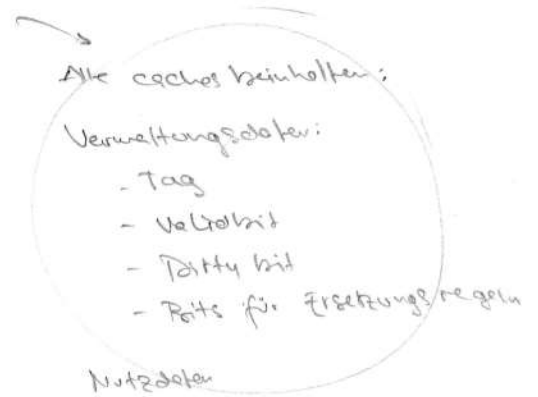
- + Eindeutige Search-Position im Cache für jeden Block
- + einfache Cache Verwaltung ohne Multiplexer
- ~ niedrige Hit-Rate

## Fully associative Cache

- Block kann überall sein im Cache (Suche aber mit XOR Block! nicht einzeln)
- komplizierte Verwaltung, großer Multiplexer
- + optimale Hit-Rate

## N-way Set associative Cache

- + Suche beschränkt → Pro Block gibt es N mögliche Positionen (im Set)
- + Einfache Ersetzungsregeln (LRU)
- + handhabbarer Multiplexer
- ~ vernünftige Hit-Rate



Beispiel: Assoziativität bei 8 Einträgen:

1-way set associative  
(= direct mapped)

Set	Way 0	
	Verw.	Data
0		
1		
2		
3		
4		
...		
7		

8-way set associative  
(= fully associative)

Set	Way 0		...	Way 7	
	Verw.	Data		Verw.	Data
0					

4-way set associative

Set	Way 0		Way 1		Way 2		Way 3	
	Ver	D	V	D	V	D	V	D
0								
1								

2-way set associative

Set	Way 0		Way 1	
	V	D	V	D
0				
1				
2				
3				

## Umgang mit Hit & Miss bei Write-Zugriffen

Bei Lesezugriffen:

Hit → fertig

Miss → nachladen und nach Regeln ein Set, Index überschreiben

Bei Schreibzugriffen

Hit → Cache und DRAM müssen konsistent sein

Miss → Nachladen optional

## Write-Miss

(WT) Write through

aktuellere cache und DRAM gleich setzen

+ Datenkonsistenz immer garantiert

- Häufige Zugriffe auf DRAM → Performance Verlust

(CB) Copy Back / write back

aktuellere cache und markiere cache-Block mit "dirty" tag,  
aktuellere Hauptspeicher wenn Block auf Cache entfernt wird

- Datenkonsistenz nicht garantiert

- Read miss wird langsamer wenn Blöcke entfernt werden

+ Seltener Zugriffe auf DRAM → Performance Gewinn

Buffered write through

(Vorteile von WT und CB)

neue Werte werden in Cache und zweiten schnellen Zwischenspeicher eingetragen,  
falls Puffer voll (nicht wenn cache Block entfernt wird), muss Prozessor arbeiten

## Write-Miss

Write-Around

Schreibe direkt in Speicher bei miss

(meist mit WT verwendet)

Fetch-on-write

Fülle einfach cache und lade restlichen Blockwörter aus DRAM nach

(simuliert dann quasi write hit, danach)

Vermohtungsarten:

- Valid Bit

- Dirty Bit

- Ersetzungsstrategie Bits

# Cache Replacement Strategien

Auswahl zum ersetzen:

Direct Mapped Cache : Eindeutig

Set Associative Cache : Innerhalb von Set

Fully Associative Cache : beliebige Auswahl

} → braucht replacement Strategie

## LRU - Least recently used

Am längsten (Zeit) nicht verwendet

(mit Queue)

- aufwändig zu realisieren

## LFU - Least frequently used

Am wenigsten verwendet

(mit Counter)

häufig verwendet → erst nach vielen Misses ersetzbar

"reference bit"

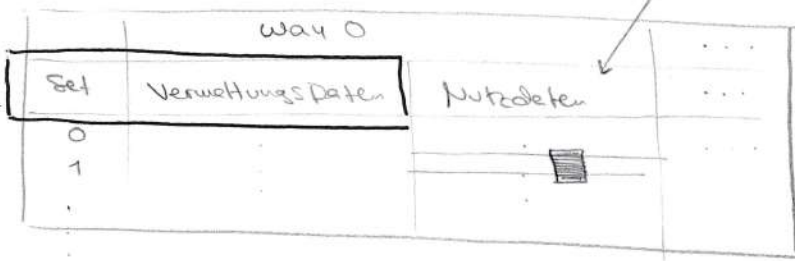
Random

FIFO (first in - first out)

simple

Übersicht: Cache - Aufbau:

Ganze Blöcke mit jedem Wort übertragen



Angenommen bei 8 Einträgen:

- 1 way, 8 Sets → direct mapped
- 8 ways, 1 Set → fully associative
- Sonst → n-way set associative

} # Blöcke = #sets · #ways  
Pro block x Wörter

Index/Set

Spricht Zeile an

Tag

XXXX[Index]  
↑  
Macht eindeutig

Valid-Bit

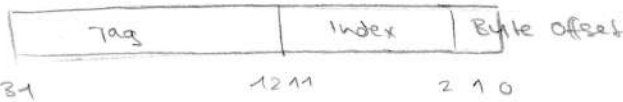
sagt ob benutzt

Dirty-Bit

Siehe copy back  
us  
write through

Ersetzungsregel-Bits

Adresse:



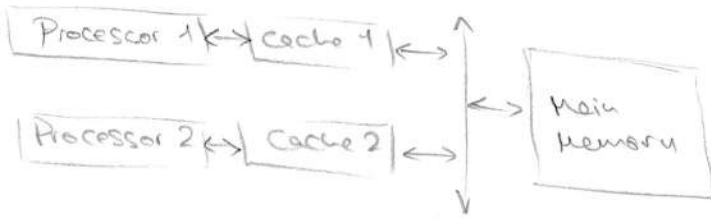
Bestimmt aus Adresse

Cache Position

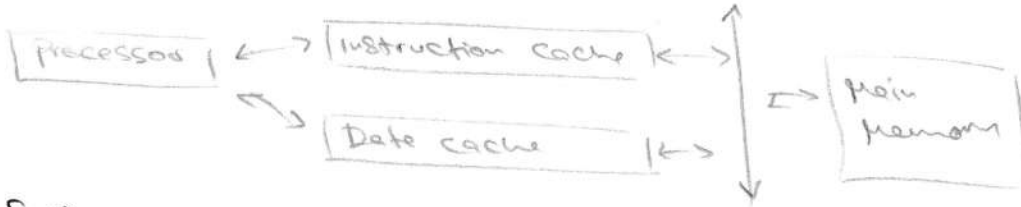
↑  
hier: 4 Wörter pro Block: (00, 01, 10, 11)

# Exkurs: Multiprozessor-Systeme

Multiprozessor-Systeme mit Caches



Harvard-Architektur



## Exkurs: Speicherverwaltung

BUS ... früher: Backpanel mit sockets

- Data-Transfer-Bus: zwischen CPU und Memory

- Address
- Data
- Memory

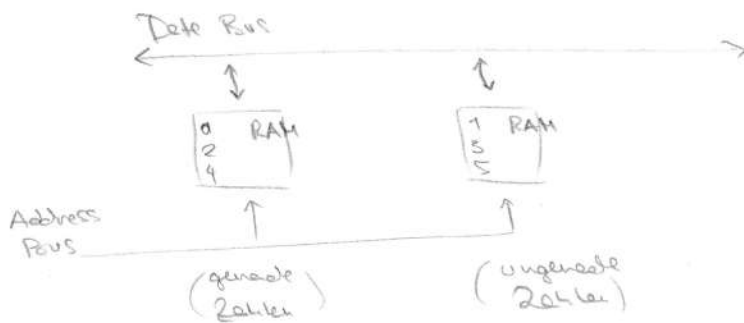
} System Bus

- Arbitration-Bus: für mehrere Nutzer

- Interrupt-Bus

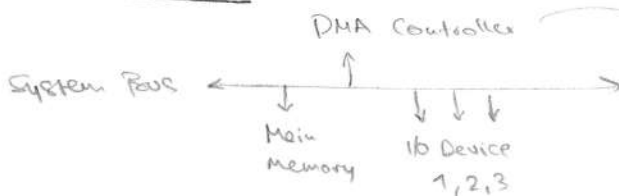
## Interleaved Memory

Speicher aufteilen in Segmente um Zugriffe zu beschleunigen



(2-fach interleaved Memory)

## Direct Memory Access



ist zuständig für Kommunikation zwischen Prozessor & peripheren Geräten  
ähnlich wie Bridge verschiedene System-Busse miteinander verbindet



## 18 Multicore

Clock-Speed und Performance kann nicht mehr steigen.

Deshalb mehrere Cores (CPUs) auf einem Chip → Parallelisierung — Schwache Speichermodelle (im nächst. Kapitel)

parallelis. Software      Multicore Systeme

### Amdahls Gesetz

über Parallelisierung

$t_p$  parallelisierbare Laufzeit

$t_s$  sequentielle (nicht parallelisierbare) Laufzeit

Laufzeit vor Parallelisierung

$$t_{tot}^v = t_s + t_p$$

$$\text{Speedup } S = \frac{t_{tot}^v}{t_{tot}^n}$$

Laufzeit nach Parallelisierung

$$t_{tot}^n = t_s + \frac{t_p}{c} \leftarrow c \text{ threads}$$

$$\text{Ausnutzung } a = \frac{t_{tot}^v / c}{t_{tot}^n} \leftarrow \begin{array}{l} \text{maximale} \\ \text{Auslast.} \\ \text{echte} \\ \text{Auslast.} \end{array}$$

Erkenntnis:

Große Speedup verbesserung bei gleichem Problem SCHWERER als bei schwierigerer Problemstellung!  
↑ höhere Laufzeit

### Vektor-Rechner

meist bei Bildverarbeitung benutzt

ALU rechnet mit Blöcken (Vektoren) (zB 4 Gleitkommazahlen à 32 Bit)

+ vereinfachte Synchronisation

+ reduzierte Kontrollbandbreite

+ effizient bei paralleler Strukturen (zB Bildern)

+ erweitertes CPU Instruktionssatz

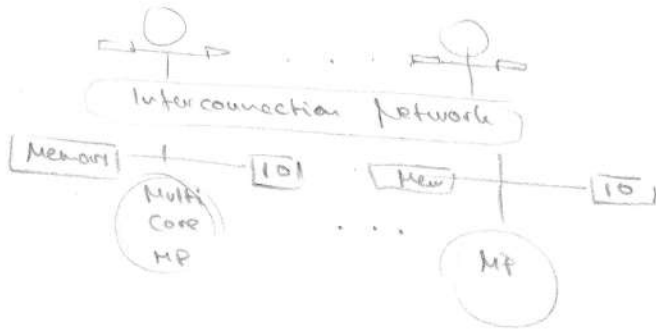
### GPU Graphic processing unit (nicht LVA relevant)

(Grafik-Karten)

ganz andere Architektur als CPU → für parallele Berechnungen geeignet

eigene Programmierertechniken

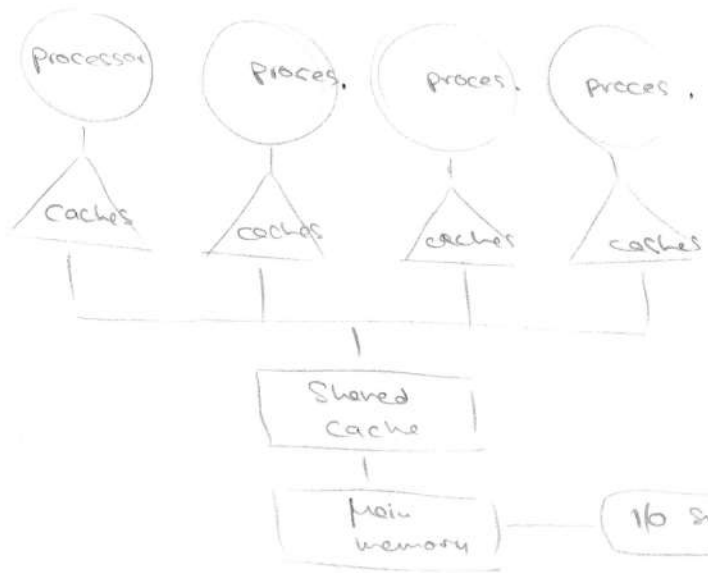
# DSM - Distributed Shared Memory



- Ein Speicher pro CPU
- Problem: NUMA  
Non uniform memory access / Latency  
(Speicherzugriff auf einen Speicher außerhalb der eigenen CPU)

# SMP - Symmetric Multiprocessors

(Für TGI das einzig relevante)



- keine # an Prozessoren
- Cache Hierarchy

- Uniform memory access / Lat
- gleicher physikalischer Adressraum für  $\forall$  Prozessoren
- Synchronisation mit Shared variables

# 19 Speichermodelle

Sequentially consistent world:

## Interleavings

Alle möglichen Statementreihenfolgen in threads

Thread 1 (a, c)  
Thread 2 (b, d)

→ lokale Ordnung muss erhalten bleiben:  
 $a < c, b < d$   
(a vor c), (b vor d)

Aber - threads selbst können an beliebigen Zeitpunkten ausführen:

(a, c, b, d) (b, a, c, d)  
(a, b, c, d) (b, a, d, c)  
(a, b, d, c) (b, d, a, c)

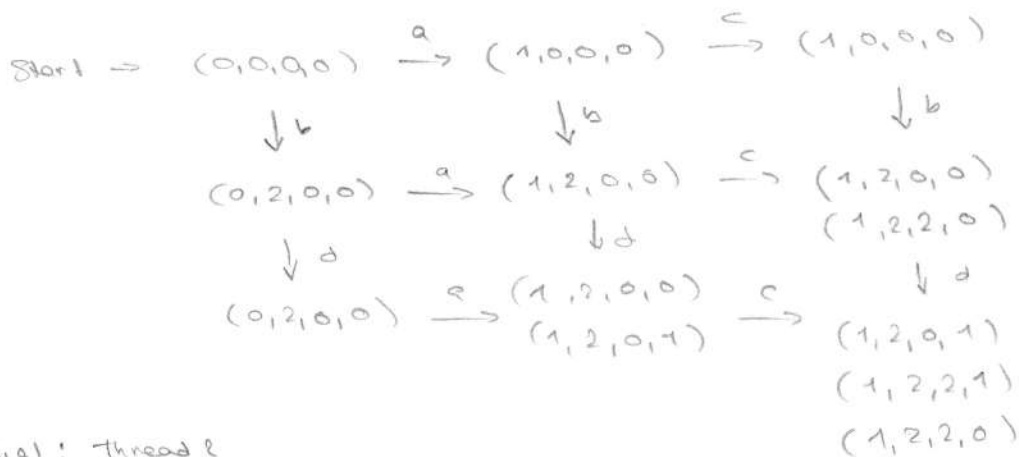
## Interleavingsgraph

(A, B, C, D) = (0, 0, 0, 0)

Thread 1:  
a: A := 1  
c: C := B

Thread 2:  
b: B := 2  
d: D := A

↓  
Verschiedene Scheduling-Strategien  
(Hardware abhängig)  
alle möglich



↓ Vertikal: Thread 2  
→ Horizontal: Thread 1

Woran erkennt man dass wir in einer sequentially consistent world sind?  
Sequentielle Ausführung inhaltslos eines Threads?

Beispiel:

Output kann nie (1, 2, 0, \*) sein, da a immer vor d ausgeführt werden muss.  
D

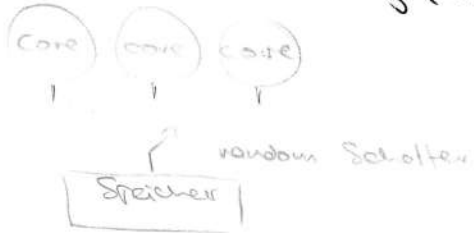
→ SEQUENTIALLY CONSISTENT MEMORY MODEL (SC)

## Race Condition (im SC Modell)

Konstellation, in der Output von der zeitlich versetzten Ausführung Anweisungen abhängt. (bei mehreren Threads)

↳ Setzt voraus, dass -Zuweisungen atomar sind

## Vereinfachte Erklärung für SC



- Ein einziger globaler Speicher, jeder Core erzeugt Speicheroperationen in Programm-Ordnung
- Speicher-Ordnung sequentiell weil Schalter zufällig einen Core access erlaubt zum Speicher

→ Alle Operationen atomar, weil jeweils global nur eine Operation ausgeführt werden kann

## Zusammenfassung:

Das SC Modell setzt voraus:

- atomare Wertzuweisungen
- Ausführung von Thread internem Code nach Programm-Ordnung (in modernen Rechnern nicht)

Man kann mit Hardware Hilfsmittel SC gewährleisten.

- + einfach zu programmieren / debuggen
- + einfach Korrektheit zu beweisen
- langsam und nutzt Multi-Core Potential nicht aus.

# Atomics

Beispiel: 16-Bit Wert Zuweisung erfolgt nicht atomar: je 8 Bit

1. Thread 1 überträgt 1. Hälfte auf Variable store\_here:

(-1) 1111 1111 1111 1111  $\rightarrow$  store\_here = 1111 1111 0000 0000

2. Thread 2 liest store\_here  $\rightarrow$  Falscher Wert! (-128)

3. Thread 1 überträgt den Rest.

### Wie kann man Ausführungen atomar machen?

Nur durch Hardware: atomare Operationen als Instruktionen

Programmierbar durch atomare Typen und Operationen auf atomaren Variablen in dieser Programmiersprache

### Synchronisation & Kommunikation zwischen Threads (Einfachste Art)

Beispiel: Data von Thread 1 auf Thread 2 ohne Unterbrechung  $\rightarrow$  mit flag

$F=1 \rightarrow$  lesen  
gefolterlos  
fertig übertragen!

Anwendung von Atomaren Operationen:  
"Producer - Consumer - Pattern"

(F ist atomar, weil es nur ein Bit hat)

$(D, F, X) = (0, 0, 0)$

Thread 1:

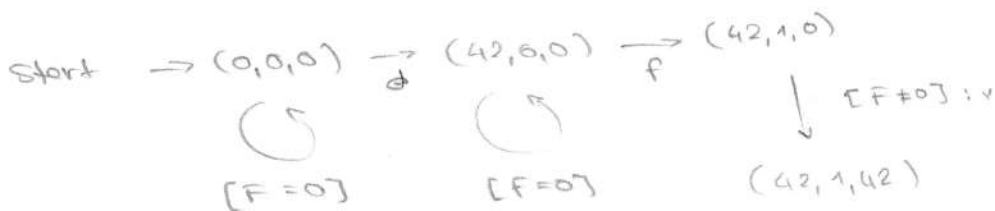
d: D=42

f: F=1

Thread 2:

if: if F=0 then goto if (wait)

x: X=D



hier keine operation für  $F \neq 0$  existiert.

### Großer Nachteil:

Thread 2 muss so lange warten bis Thread 1 fertig ist.

Wir nutzen also die Vorteile von Multithreading bei diesem Lösungssatz nicht!

(mehr dazu folgt...)

Alles bisher war in der sequentially consistent world mit dem SC-Memory-Mod.

Dabei wurde vorangeseht dass:

- ① - Atomare Wertzuweisungen
- ② - Programmordnung = Ausführungsreihenfolge in einem Thread

Aus dem Interleavings-Graphen konnte man sehen, dass:

Mit einem Core / Thread, kann man Ausführungsreihenfolge nicht umordnen.  
Mit mehreren Cores / Threads, kann Ergebnis unterschiedlich sein

Optimierung im Hintergrund:

Programmordnung  $PO \neq$  Exekutionsordnung  $EO$   
(Ausführungsreihenfolge)

- Programme wird von Compiler und CPU umgeordnet.  
Es soll aber nicht zu einem anderen Ergebnis führen (bei einem Thread / Core)

Bei Multicore-Rechnern aber:  
nicht mehr sequentially consistent!

↳ also immer noch sequentially consistent obwohl 2. Bedingung verletzt wurde

Beispiel für SC-Verletzung bei Multithreading wobei  $EO \neq PO$ :

Einfaches Beispiel:

Thread 1:	Thread 2:
a: A = 1	b: B = 2
c: C = B	d: D = A

Ohne Einhaltung der PO innerhalb der Threads:

Thread 1:	Thread 2:
a: A = 1	d: D = A
c: C = B	b: B = 2

SC ist verletzt, dadurch ist (1, 2, 0, 0) als Endzustand möglich

obwohl es kontra intuitiv ist und zeitlich relativ verletzt, wobei Kausalität

→ RELAXED MEMORY MODEL

Im relaxed memory model funktioniert "Producer - Consumer - Pattern" nicht.

Initialisierung

$$P = F = x = 0$$

Thread 1

$$d: D = 42$$

$$f: F = 1$$

Thread 2:

if: if  $F = 0$  then goto .if

$$x: x = D$$

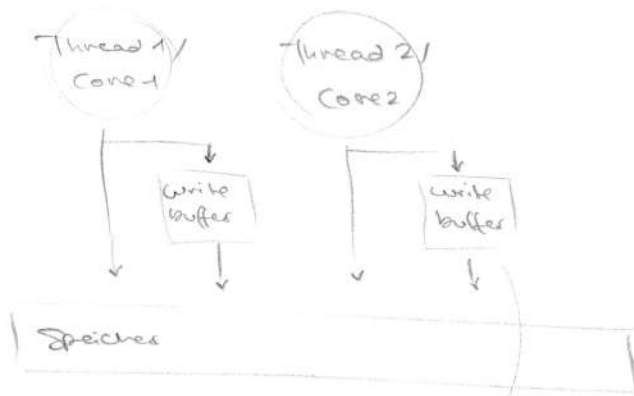
Reihenfolge:

$$f, \text{if}, x, d \longrightarrow (42, 1, 0)$$

nicht erfolgreicher Synchronisierest!

Beispiel für SC-Verletzung bei Multithreading wobei  $EO \neq PO$ ;

Architektur ohne Caches  $\rightarrow$  Producer - Consumer Pattern



Initial:

Flag 1 = 0

Flag 2 = 0

Thread 1:

Flag 1 = 1

if Flag 2 = 0 then

... critical (no interruption)

Thread 2:

Flag 2 = 1

if Flag 1 = 0 then

... critical

wenn Buffer voll,  
dann alles mit Speicher  
abgleichen und Buffer  
ausleeren

Problem:

Lesen-Operationen (wie die in if Abfrage) überholen write Buffer!

Möglicher Ablauf:

Flag 1 = 1 (im write-buffer)

if (Flag 1 = 0) then ...  $\rightarrow$  wird ausgeführt obwohl es nicht sein darf!  
ERROR

Conclusion:

SC ist verletzt, weil durch Cache die Wertzuweisungen nicht mehr atomar sind.

Definition von Atomizität mit mehreren Caches und Cores:

SCHREIBEN:

Operation darf nicht unterbrochen werden

Alle Cores müssen gleichzeitig Zugriff auf den neu geschriebenen Wert haben, wenn es ein Core hat.

LESEN:

So lange aufschreiben bis alle Cores gleiche Cache-Kopien haben ab letzter Schreiboperation



Weiteres Beispiel:

(Demonstration der Verletzung der Zeitlichen Relativität)

|||  
○  
Beobachter 1:  
"a war vor b"

Thread  
a: ...  
b: ...

|||  
○  
Beobachter 2:  
"a war nach b"

Beispiel:

T1      T2      T3      T4  
a: A=1    b: B=1    c: C=A    e: E=A  
            d: D=B    f: F=B

Initialisierung:  
A=B=C=D=E=F=0

Da alle Anweisungen in beliebiger Reihenfolge ausgeführt werden dürfen:

(A, B, C, D, E, F) = (0, 0, 0, 0, 0, 0)

↓ d: D=B

(0, 0, 0, 0, 0, 0)

↓ b: B=1

(0, 1, 0, 0, 0, 0)

↓ f: F=B

(0, 1, 0, 0, 0, 1)

↓ e: E=A

(0, 1, 0, 0, 0, 1)

↓ a: A=1

(1, 1, 0, 0, 0, 1)

↓ c: C=A

(1, 1, 1, 0, 0, 1)

( A B C D E F )  
  1 1 1 0 0 1

Aus der Sicht von Thread 3:

C=1, D=0 ⇒ a vor b ausgeführt (also T2 nach T3)

Aus der Sicht von Thread 4:

E=0, F=1 ⇒ b vor a ausgeführt (also T1 nach T4)

Das würde bedeuten:

T2 < T3

T1 < T4

↪ nicht möglich mit diesem Output

## Allgemein: Speichermodelle

- Hardware muss wissen: welche Art von Instruction Scheduling
- Programmiersprache sagt Hardware wo fences gelegt werden müssen, beschreibt Thread Interaktion mit Speicher
- Fences können nicht automatisch gesetzt werden, (und benötigen zusätzliche Hardware)

## Speichermodell-Operationen

store : atomic\_store (atomic\_var, memory\_order)  
load : atomic\_load (atomic\_var, memory\_order)  
exchange : atomic\_exchange (atomic\_var1, atomic\_var2, memory\_order)

## Auswahl unter Sprachen:

JAVA: SC

C++ : SC, RA, Relaxed ...

C : SC, RA, Relaxed ...

} (schwaches Speichermodell)

Auswahl bei memory\_order:  
SC, relaxed, release und/oder acquire

## Auswahl unter Architekturen:

Intel x86/64: SC

ARM: schwaches Speichermodell

RISC-V: schwaches Speichermodell

## Performance:

SC < RA < Relaxed

↑  
Sinnvoll wenn keine negativen Konsequenzen durch Anwendung.

# → RELEASE / ACQUIRE MEMORY MODEL

Benötigt Memory fences / Memory Barriers (die nur load und store operationen betreffen)

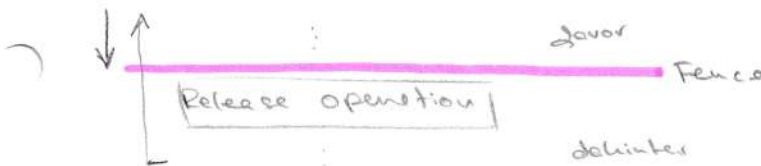
- Versteht Einseitig des Umwandeln über Zaun
- Angesprochen mit Release / Acquire Commands

↳ Nur Operation auf atomare Variablen (≠ gewöhnlich) möglich

## Relaxed

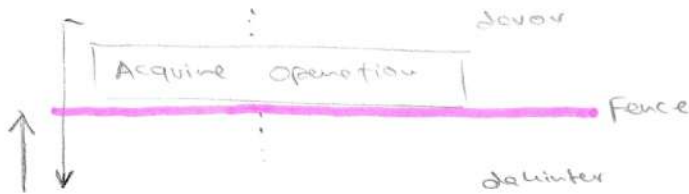
keine Einschränkung

### Release (store, speichern)



operationen dahinter dürfen Fence überschreiten

### Acquire (load, lesen)



operationen davor dürfen Fence überschreiten

Alles unter der Annahme, dass Befehl selbst NIE verschieben wird

Wichtig: Datenabhängigkeiten per default immer berücksichtigt:

```
z.B. atomic(counter++)  
if (atomic(counter) > 0)  
    ...
```

← schreiben auf counter  
← lesen auf counter

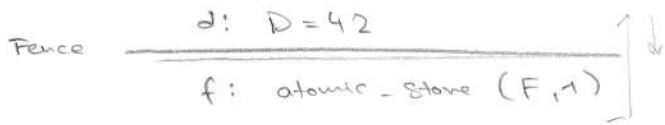
Darf relaxed werden

# Release - Acquire Beispiels:

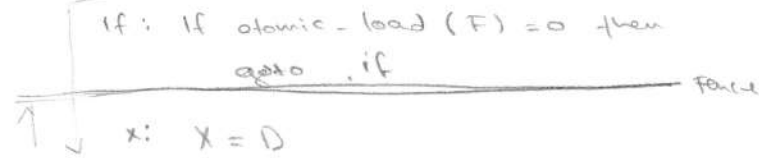
Producer - Consumer - Pattern

$$D = F = X = 0$$

Thread 1



Thread 2

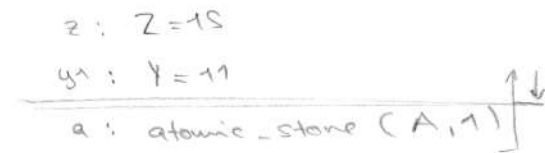


Dadurch sieht man:

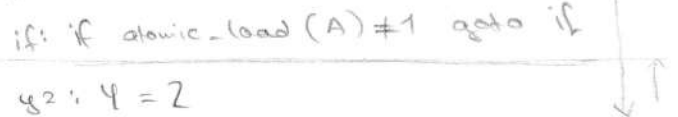
Synchronisation und Kommunikation mit diesem Pattern funktioniert mittels RA - Pattern

Wann ist es sinnvoll aus RA relaxed zu machen?

Thread 1:



Thread 2:

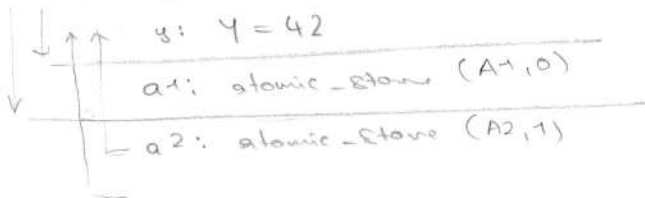


→ Lösung:

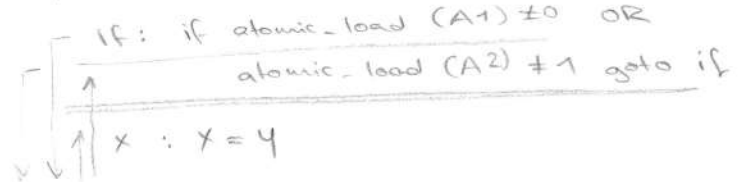
hier gar nicht:

1. z darf nicht hinter a kommen
2. y2 darf nicht vor if kommen
3. y1 darf nicht hinter a kommen, sonst könnte nach y2 nochmal y1 y überschreiben

Thread 1:



Thread 2:



→ Lösung:

if Schleife wird verlassen, wenn  $(A1 = 0 \wedge A2 = 1)$

a1 darf relaxed werden, weil es egal ist ob zuerst a1 oder a2 gespeichert wird, aber a1 UND a2 dürfen nicht gleichzeitig relaxed werden

a1 vor a2 darf relaxed werden

related werden darf

Wenn Compiler if - booleans von links nach rechts bewertet, beeinflusst das (neben der Tatsache dass man a1 oder a2 relaxed hat) welche if Abfragt

Das relaxen der ersten (links) if-Abfrage, wäre ausreichend dafür, dass nichts noch vorne wandern kann.

Somit hätten wir am Ende 2 Memory-Fences weniger!

Das relaxen der ersten (links) if-Abfrage, wäre ausreichend dafür, dass nichts noch vorne wandern kann.

Somit hätten wir am Ende 2 Memory-Fences weniger!

Fortsetzung des Beispiels:

Angenommen

Thread 1

y: y = 42

a1: atomic\_store(A1, 0)

a2: atomic\_store(A2, 1)

Thread 2:

if atomic\_load(A1) != 0 OR  
atomic\_load(A2) != 0 goto if

x = y

Case 1: relaxed a1

Mögliche Anordnungen:

Thread 1

y: y = 42

a1: A1 = 0

a2: atomic\_store(A2, 1)

möglich

$$\begin{pmatrix} y \\ a1 \\ a2 \end{pmatrix} \begin{pmatrix} a1 \\ y \\ a2 \end{pmatrix}$$

✓ gültig

Case 2: relaxed a2

Mögliche Anordnungen

Thread 1

y: y = 42

a1: atomic\_store(A1, 0)

a2: A2 = 1

möglich

$$\begin{pmatrix} y \\ a2 \\ a1 \end{pmatrix} \begin{pmatrix} a2 \\ y \\ a1 \end{pmatrix}$$

✓ gültig

Der einzig ungültige Fall den  
man vermeiden muss:

$$\begin{pmatrix} a1 \\ a2 \\ y \end{pmatrix} \begin{pmatrix} a2 \\ a1 \\ y \end{pmatrix}$$

## Weiteres Beispiel:

Thread 1  
↓ ↑  
y: y = 42  
a: atomic\_store(A, 1)

Thread 2  
↑ ↓  
x: x = atomic\_load(A)  
if: if x ≠ 1 then goto v  
z: z = x + 1

### → Lösung:

- a kann relaxed werden, weil y irrelevant ist
- x kann relaxed werden, weil if und z über x Datenabhängig sind und by default nur danach stehen können (nur z darf mit if vertauscht werden)
- theoretisch muss A nicht atomic sein.



## Blockierendes Verhalten

bisher in „Kommunikation & Synchronisation“: 1 Thread wartet in Endlos-Schleife auf Flag

## Alternative: Semaphore

Eine Semaphore ist eine atomare Variable die wie ein Zähler funktioniert und am Anfang auf eine Zahl  $\geq 0$  gesetzt wird (max # der Threads die parallel laufen dürfen)

Semaphore = Zähler = 1 initialisiert

↓  
(Siehe Wikipedia: Bäckerei Beispiel mit Zimmern)

(wait()) Lock()

Zähler --

Wenn Zähler  $\geq 0$ , darf aufrufender Thread weiterexecutieren, wenn aber  $< 0$  muss Thread in Queue warten (exekution stoppen)

(post()) Unlock()

Zähler ++

Wenn Zähler  $\leq 0$  zu 1 wird, Thread aus Queue holen, sonst nichts

## Race Condition

Ausführungsreihenfolge beeinflusst Ergebnis wenn nicht alles atomar abläuft, memory fences deshalb notwendig!

Lösung:

Read-Modify-Write Operationen (in Hardware  $\rightarrow$  Prozessorarchitektur)

versichert Atomicität, für Semaphore geeignet

abhängige Instruktion

## Vor und Nachteile von Semaphoren

- Korrektheit schwer prüfbar → für jedes Lock braucht es ein Unlock



Deshalb leichter handhabbare Alternativen in höheren Sprachen:

- Monitore
- Synchronized Objects, Java
- Protected Object, Ada

+ leicht verständlich (siehe Beispiel)

+ Performant

+ Sequentially Consistent wenn Zähler = 1

- Wenn ein Thread während kritischer Operation zwischen lock() und unlock abstürzt, können andere keinen Fortschritt machen

- Thread dispatching!

(Stoppen und Starten von Threads) sehr langsam

## Beispiel für Semaphoren:

globale Variable Z soll nur gleichzeitig von einem Thread benutzt werden.

Global:

Z = ...

S = 1

Thread:

lock(S)

Z-local = Z

Z-local = Z-local + ...;

Z = Z-local

Unlock(S)

## Nicht Blockierendes Lesen

Statt Semaphoren direkt Read-Modify-Write-Operation von Hardware benutzen.

"Optimistisches Vorgehen":

- zB
1. Probieren Änderung auf globale Variable Z durchzuführen
  2. Falls nicht unterbrochen Ende  
sonst bei 1. fortsetzen

Funktion: RMV (V, alter\_wert, neuer\_wert)

returned true: Atomare Variable V hatte immer noch alten\_wert und Wertzuweisung erfolgreich.

returned false: sonst

Beispiel:

Z ist atomar

Thread:

1. Z\_local = Z

2. if not RMV (Z, alter\_wert, neuer\_wert) then goto 1

↗ false  
↘ true

Auswirkung:

Endlosschleife, so lange bis  $Z = Z\_local$ .

Dann bekommt  $Z = Z\_local + \dots$

und Programm endet in diesem Thread

## Vor und Nachteile

- über 100% schneller →
- + Effizient bei wenigen Threads, da sie nicht gestoppt / gestartet werden müssen.
  - + Wenn ein Thread abstürzt können andere Fortschritte machen
  - + Sequentially consistent und Release and Acquire Speichermodell möglich
  - Schwer verständlich → vor allem mit RA-Speichermodell