

Literatur

1. David Money Harris & Sarah L. Harris: "Digital Design and Computer Architecture", Morgan Kaufmann 2007

- 2. John F. Wakerly: "Digital Design, Principles & Practices", Prentice-Hall 2001
- 3. David A. Patterson, John L. Hennessy: "Computer Organization & Design, The Hardware/Software Interface", Morgan Kaufmann 1998 (git's and als PDF in Netz)
- 4. Peter J. Ashenden: "The Designer's Guide to VHDL", Morgan Kaufmann 1996
- 5. Donald E. Thomas, Philip R. Moorby: "The Verilog Hardware Description Language", Kluwer Academic Publishers 2002

Übersicht

- 1. Hardware - Entwurf
 - Kombinatorische Digitalschaltungen ("Schaltetze")
 - Sequentielle Digitalschaltungen ("Schaltwerke")
 - Moore- und Mealy-Automaten
 - Spider und programmierbare Logikbausteine
- 2. Rechnerarchitektur
 - Performance-Messung
 - Instruction-Set-Architecture

Arithmetik

Datenpfad und Steuerung

Pipelining

Caches

~~Virtueller Speicher~~

Eingabe/Ausgabe

~~3. Eingebettete Systeme und ihre Umgebung~~

~~Meßprinzipien~~

~~Sensoren~~

~~Wandler (ADC/DAC)~~

~~Aktoren~~

~~Kommunikation~~

4. Werkzeuge (werden im Praktikum behandelt)

Hardware-Beschreibungssprachen

Simulatoren

Syntheswerkzeuge

3 Testfragen:

- 1) Was ist der Unterschied zwischen einem Schaltplan und einem Schaltwerk?
- 2) Was zeichnet eine synchrone Schaltung aus?
- 3) Was ist "Metastabilität"? Muß man darauf überhaupt achten?

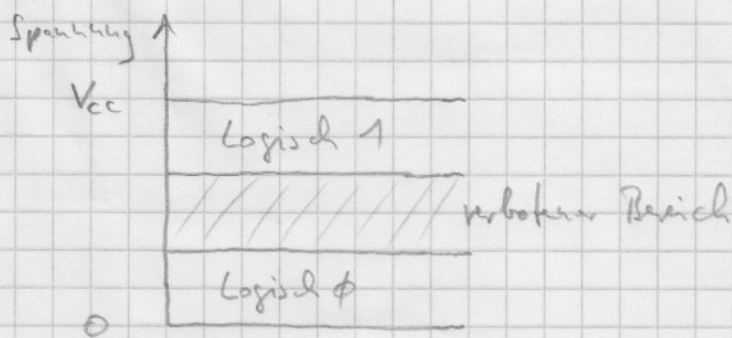
1. Hardware - Entwurf

1.1. Analog vs. Digital

Analogrechner: Die Rechengrößen werden durch eine physikal. Größe (z.B. Spannung) repräsentiert, die jeden Wert zwischen zwei Grenzen (z.B. $\pm 10V$) annehmen kann und die der Rechengröße bis auf eine Proportionalitätskonstante gleich ist.

Digitalrechner: Die Rechengrößen werden durch eine gewisse Zahl (z.B. 32) von Elementen mit diskreten Zuständen (meistens zwei) dargestellt. Die zwei Zustände heißen FALSE und TRUE (oder kürzer = 0 und 1).

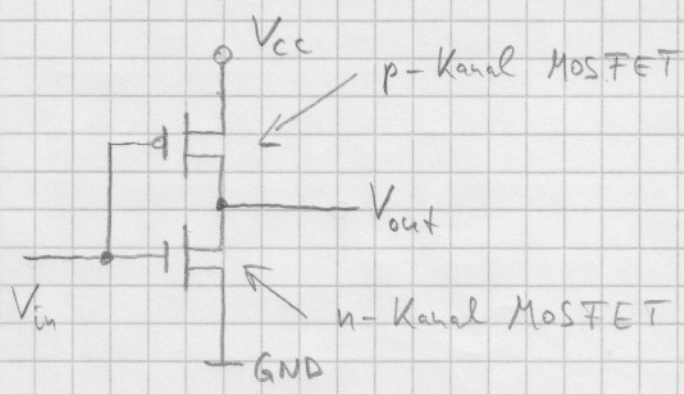
Heißt das, wir finden Nullen und Einsen in unseren Digitalschaltungen? Nein! Zwei Spannungsbereiche dienen zur Darstellung der 0 und der 1:



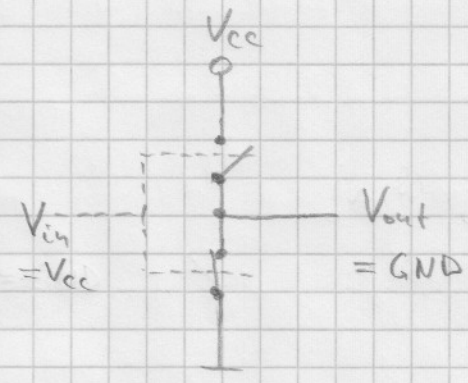
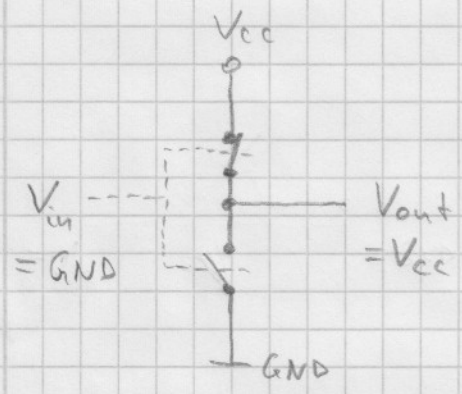
1.2. Gatter

3 Gattertypen reichen aus, um jede Boole'sche Funktion $(x_1, x_2, \dots, x_n) \mapsto f(x_1, x_2, \dots, x_n)$ zu realisieren.

a) NICHT - Gatter ("Inverter")



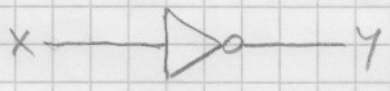
Schaltung



Verhaltensweise

x	y
0	1
1	0

Wahrheitstabelle



Schaltzeichen

Algebraische Schreibweise: $y = \bar{x}$
 $y = x'$

Verilog:
 a) not $g1(y, x);$
 b) assign $y = \sim x;$
 "Schalterschreibweise"
 "Verhaltensbeschreibung"

b) UND-Gatter

x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1

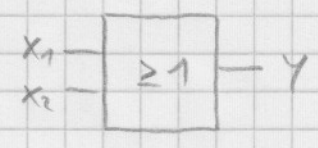
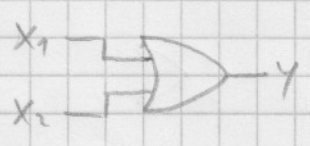
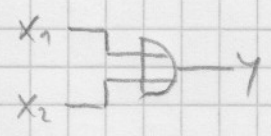


$y = x_1 \cdot x_2$

and $g1(y, x_1, x_2);$
 assign $y = x_1 \& x_2;$

c) ODER - Gatter

X_1	X_2	Y
0	0	0
0	1	1
1	0	1
1	1	1

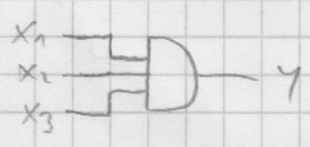


$Y = X_1 + X_2$ or $g_1(Y, X_1, X_2)$;

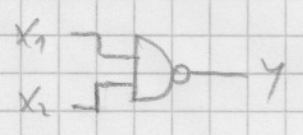
assign $Y = X_1 | X_2$;

Bem.:

a) Mehr als zwei Eingänge ^{sind} möglich; z.B. $Y = X_1 \cdot X_2 \cdot X_3$

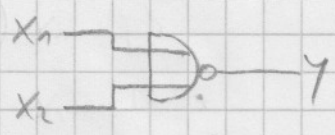


b) NAND \equiv "not and" = $Y = \overline{X_1 \cdot X_2}$



X_1	X_2	Y
0	0	1
0	1	1
1	0	1
1	1	0

NOR \equiv "not or" = $Y = \overline{X_1 + X_2}$



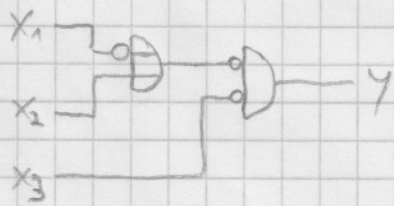
X_1	X_2	Y
0	0	1
0	1	0
1	0	0
1	1	0

XOR \equiv "exclusive or" : $y = x_1 \cdot \bar{x}_2 + \bar{x}_1 \cdot x_2$ (6)



x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

c) Der "Invertierklingel" kann auch an die Eingänge von Gattern geschrieben werden; z.B.:



$$y = \overline{(\bar{x}_1 + x_2)} \cdot \bar{x}_3$$

d) Die analoge Wirklichkeit schlägt zurück:

- Jeder Eingang belastet den treibenden Ausgang mit einem kleinen Strom \Rightarrow Man kann nicht bel. viele Eingänge an einen Ausgang hängen ("Fan-out" ≈ 20)
- Ein Signalwechsel am Eingang eines Gatters überträgt sich nur mit einer gewissen Verzögerung an seinen Ausgang ("Propagation Delay" \approx nsec)

1.3. Kombinatorische vs. sequentielle Logik

7

Kombinatorische Logik ("Schaltetze"):

- Werte der Ausgänge hängen nur von den momentanen Werten der Eingänge ab, hat also kein "Gedächtnis"
- Bel. viele Gatter und Signalverzweigungen, aber keine Rückkopplungsschleifen (= Pfade, die vom Ausgang eines Gatters direkt oder indirekt über andere Gatter wieder zu einem Eingang dieses Gatters zurückführen)

Sequentielle Logik ("Schaltwerke"):

- Werte der Ausgänge hängen außer von den momentanen Werten der Eingänge auch von bel. vielen vorhergehenden Werten ab, d.h. die Schaltung ist "zustandsbehaftet" (hat also "Gedächtnis")
- Dieses Verhalten wird durch "Rückkopplung" von Ausgängen einer Schaltung auf Eingänge derselben Schaltung erreicht. Genauereres später!

1.4. Kombinatorische Logik

(8)

1.4.1. Boole'sche Algebra (Schaltalgebra)

$$x + \phi = x$$

$$x \cdot 1 = x$$

$$x + 1 = 1$$

$$x \cdot \phi = \phi$$

$$x + x = x$$

$$x \cdot x = x$$

$$\overline{\overline{x}} = x$$

$$x + \overline{x} = 1$$

$$x \cdot \overline{x} = \phi$$

$$x + y = y + x$$

$$x \cdot y = y \cdot x$$

$$(x + y) + z = x + (y + z)$$

$$(x \cdot y) \cdot z = x \cdot (y \cdot z)$$

$$x \cdot (y + z) = x \cdot y + x \cdot z$$

$$x + (y \cdot z) = (x + y) \cdot (x + z)$$

$$x + x \cdot y = x$$

$$x \cdot (x + y) = x$$

$$x \cdot y + x \cdot \overline{y} = x$$

$$(x + y) \cdot (x + \overline{y}) = x$$

$$\overline{x + y} = \overline{x} \cdot \overline{y}$$

$$\overline{x \cdot y} = \overline{x} + \overline{y}$$

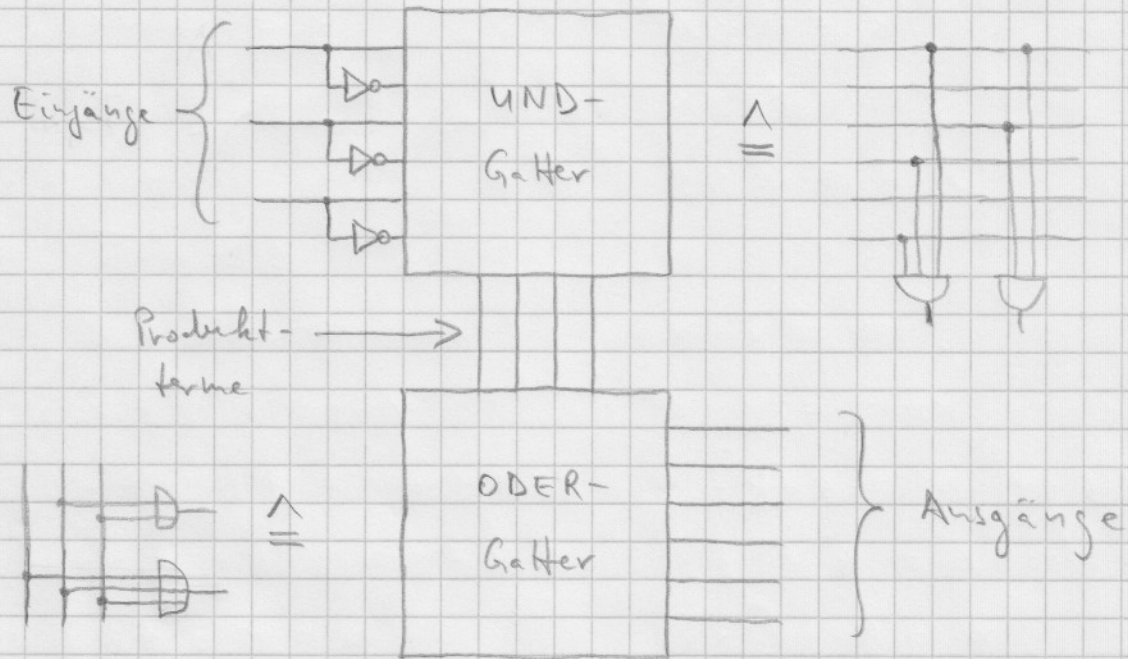
1.4.2. Normalformen, Minimierung, PLAs

Jede Boole'sche Funktion lässt sich in den beiden "Normalformen" DNF (disjunktive Normalform, "sum-of-products") und KNF (konjunktive Normalform, "product-of-sums") darstellen. Es gibt Verfahren (Karnaugh-Diagramme, Quine-McCluskey-Algorithmus, Espresso-II), die eine minimale (oder angenähert minimale) Repräsentation berechnen.

$$\text{Bsp.: } AB\overline{c}\overline{d} + \overline{A}B\overline{c}d + A\overline{B}\overline{c}d + \overline{A}\overline{B}\overline{c}\overline{d} + A\overline{B}c\overline{d} = B\overline{c} + A\overline{B}d$$

↑ ↗
"Produktterme"

Der Aufbau eines PLA ("programmable logic array") spiegelt diese zweistufige Logik wieder: (9)



Bem.: Ein ROM lässt sich auffassen als ein vollständig dekodiertes PLA: Alle 2^n Eingangskombinationen sind als Produktterme vorhanden. Die Programmierung des ROMs geschieht nur mittels der ODER-Gatter.

1.4.3. Don't Cares

Manchmal möchte man den Wert eines Digitalsignals nicht spezifizieren: ein sog. "don't care", bezeichnet mit "x".

Don't Cares kommen an zwei Stellen vor, an denen sie verschiedene Dinge bedeuten.

Bsp.:

A	B	C	Out
0	x	0	1
0	0	1	0
0	1	1	0
1	x	0	x
1	0	1	x
1	1	1	1

a) Don't Cares in den Eingängen einer Funktion:

(10)

Die Zeile $0 \times 0 \mid 1$ ist in Wirklichkeit eine Kurzschreibweise für die beiden Zeilen

$$\begin{array}{ccc|c} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \end{array}$$

und bedeutet also, daß die Funktion von der betreffenden Variablen (hier B) an dieser Stelle nicht abhängt.

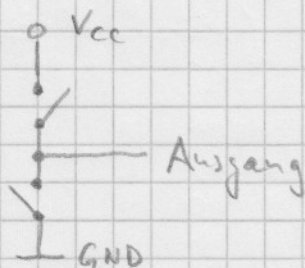
b) Don't Cares in den Ausgängen einer Funktion:

Die Zeile $101 \mid x$ bedeutet, daß der Wert der Funktion für die Eingangskombination 101 keine Rolle spielt und deshalb beliebig gewählt werden kann (z.B. so, daß die Funktion einfach realisiert werden kann).

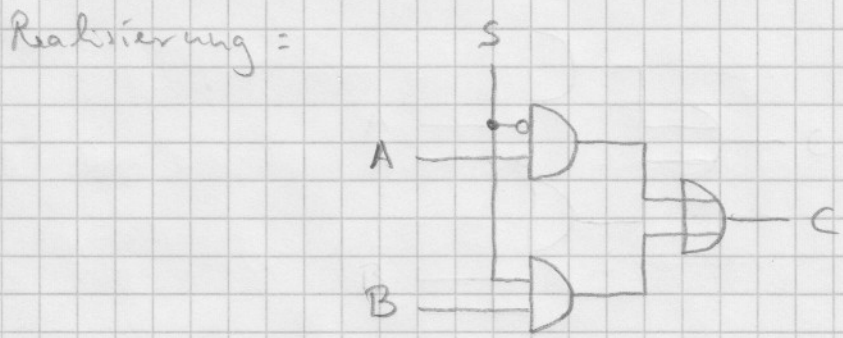
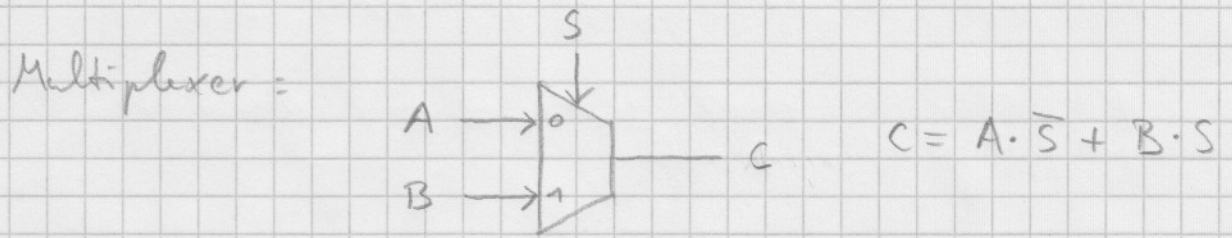
Kombinationen aus a) und b) sind möglich.

1.4.4. Tri-State-Gatter

Man kann Gatter mit einem zusätzlichen Steuereingang ("output enable") ausrüsten. In der Stellung "enabled" erlaubt er, daß das Gatter eine 0 oder eine 1 am Ausgang erzeugt. In der Stellung "disabled" wird durch Sperren beider Transistoren (öffnen beider Schalter) das Gatter effektiv vom Ausgang abgetrennt ("hochohmig gemacht", Ausgangswert "Hi-Z" oder kurz z).

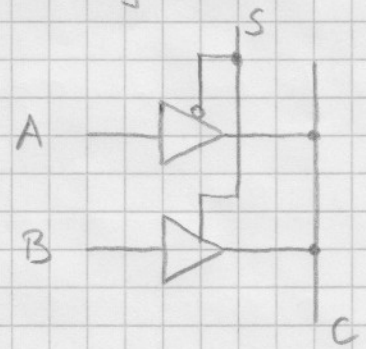


Verwendung = "Verteilter Multiplexer"



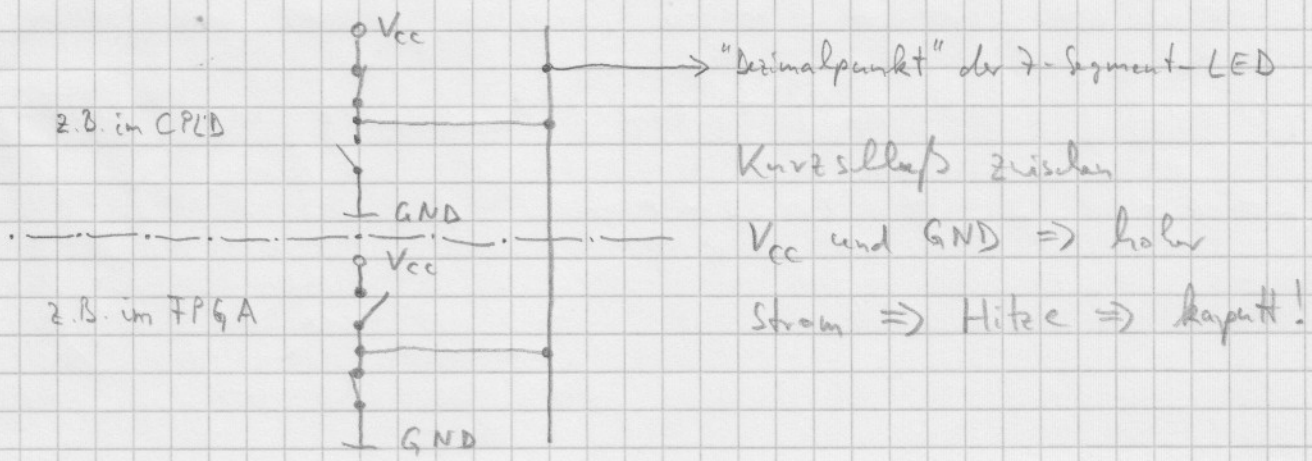
Bei 2^n Eingängen n Steuerleitungen nötig. Bei großem n hoher Schaltungsaufwand, viele Leitungen.

Realisierung mit Tri-State-Gattern:



Vorteile: Leicht erweiterbar, Logik dezentral

Nachteile: Langsamer (wg. kapazitiver Last am Ausgang), Selbstzerstörung bei falscher Ansteuerung



1.4.5. Vorsicht: Hazards!

(12)

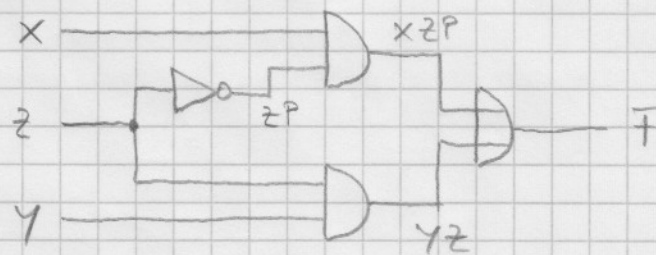
Hazard: Kurzer Puls ("Glitch") auf einem Ausgang, hervorgerufen durch Laufzeiten in der Schaltung

Hazards gibt's bei kombinatorischen Schaltungen in 3 Varianten:

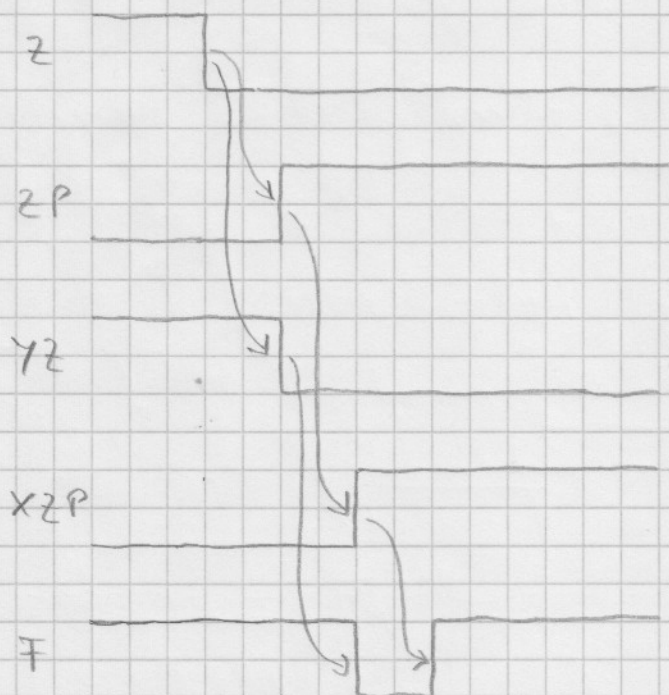
a) Statischer 1-Hazard

Ein Paar von Eingangswert-Kombinationen, die sich nur in einer Eingangsvariablen unterscheiden und beide 1 am Ausgang erzeugen, so daß am Ausgang ein kurzer ϕ -Impuls erscheint.

Bsp.:



$$X=1, Y=1, Z=1 \rightarrow \phi$$



b) Statistischer ϕ -Hazard

Das ist das duale Gegenstück zu a).

c) Dynamischer Hazard

Ein Ausgang ändert sich mehr als einmal, obwohl sich nur ein Eingang einmal geändert hat.

Das kann passieren, wenn es mehrere Pfade vom Eingang zum Ausgang gibt, mit unterschiedlichen Verzögerungszeiten.

Wahrheiten über Hazards:

1. Hazards in ad-hoc-Schaltungen sind schwer zu finden.

2. Zweistufige Logik weist keine dyn. Hazards auf.

3. Zweistufige UND-ODER-Logik hat keine stat. ϕ -Hazards, kann aber stat. 1-Hazards haben.

4. Zweistufige ODER-UND-Logik hat keine stat. 1-Hazards, kann aber stat. ϕ -Hazards haben.

5. Stat. Hazards lassen sich mit Karnaugh-Diagrammen finden und beheben. Unser Bsp. von oben: $F = x\bar{z} + yz$

Vor:
Nimmals
liegt x
und \bar{x} an
gleichen
Gatter an!

xy z	00	01	11	10
0			1	1
1		1	1	

Hinzunahme des "Konsens-Terms"
 xy behebt den Hazard;
 $F = x\bar{z} + yz + xy$ ist frei
von Hazards.

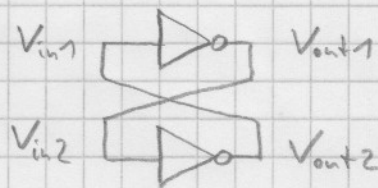
6. Hazards sind extrem wichtig beim Entwurf von nichtsynchrone sequentiellen Schaltungen. Sie spielen keine Rolle in synchronen Schaltwerken: Die Ausgänge von Schaltnetzen werden nur an Taktflanken ausgewertet, d.h. wenn alle Hazards abgeklungen sind.

1.5. Sequentielle Logik

(14)

1.5.1. Bistabile Elemente, Metastabilität

Schaltung:

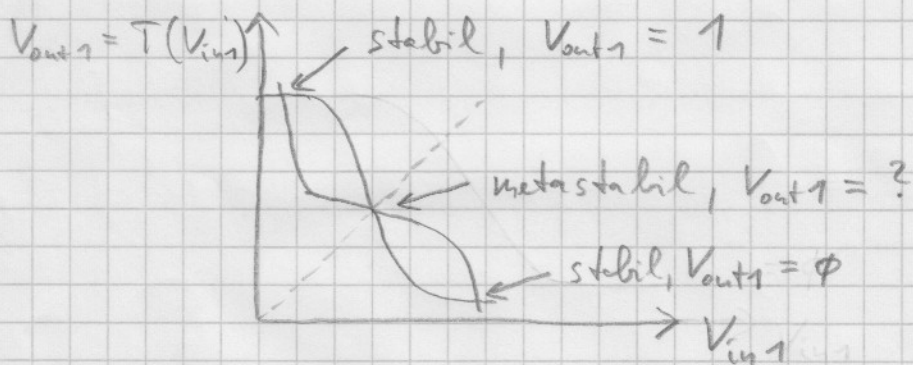


Analyse:

$$V_{out1} = T(V_{in1}), \quad V_{out2} = T(V_{in2})$$

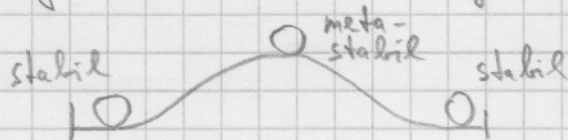
$$\Rightarrow V_{in2} = V_{out1} = T(V_{in1}) = T(T(V_{in1}))$$

$$\Rightarrow T^{-1}(V_{in1}) = T(V_{in1})$$



Bem.:

1. Jede sequentielle Schaltung kann in metastabile Zustände gebracht werden — nicht nur beim Einschalten der Betriebsspannung, sondern auch beim Verletzen der Timing-Bedingungen (mehr in Kürze). Verweildauer: unbegrenzt!
2. Verlassen des metastabilen Zustandes durch Rauschen: ist der Gleichgewichtspunkt erst mal ein wenig verlassen, wirkt die Rückkopplung (exponentiell) verstärkend.
3. Analogon: "Ball und Hügel"

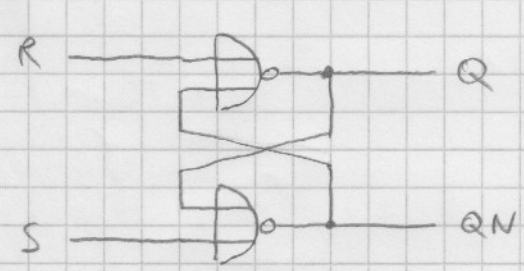


1.5.2. Latches und Flip-Flops

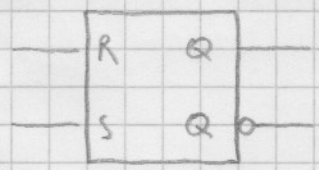
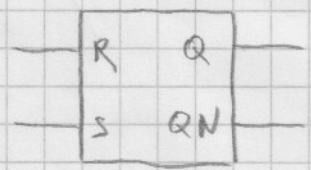
Flip-Flop: sequentielle Schaltung, die ihre Ausgänge nur zu best. Zeiten ändert, gegeben durch "Takt"

Latch: sequentielle Schaltung, die kontinuierlich auf Eingangssignaländerungen reagiert

1.5.2.1. S-R-Latch



S	R	Q	Q _N
0	0	letztes Q	letztes Q _N
0	1	0	1
1	0	1	0
1	1	0	0



↑ nicht ganz korrekt f. S=R=1

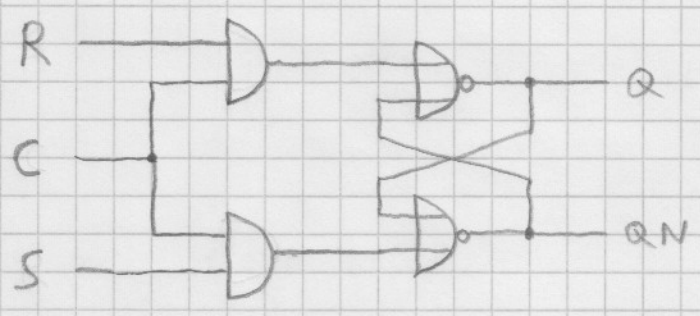
Das S-R-Latch kann metastabil werden, wenn die Pulsdauer an R oder S zu kurz ist, oder wenn bei S=R=1 die Übergänge R: 1 → 0 und S: 1 → 0 zeitlich zu nahe beieinander liegen. Dann ist sogar "Schwingen" der Ausgänge möglich:

Q: 0 → 1 → 0 → 1 ...

Q_N: 0 → 1 → 0 → 1 ...

1.5.2.2. S-R-Latch mit Enable

"Enable" steuert, ob das Latch auf seine Eingänge reagiert.



S	R	C	Q	QN
0	0	1	letztes Q	letztes QN
0	1	1	0	1
1	0	1	1	0
1	1	1	0	0
x	x	0	letztes Q	letztes QN

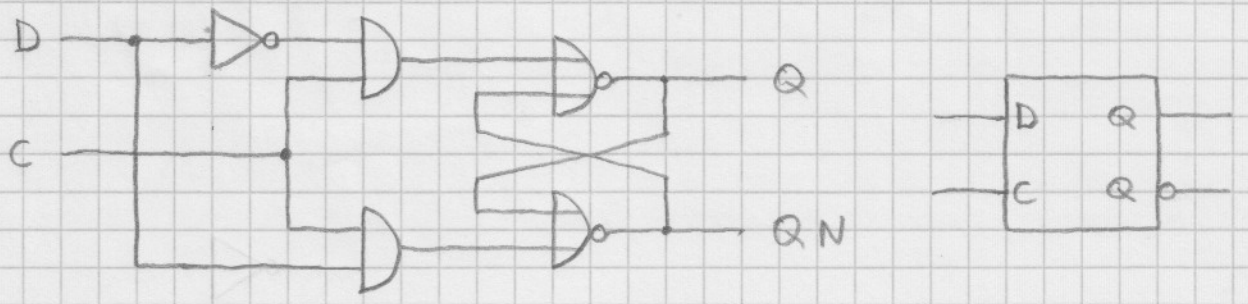
Metastabilität: bei $S=R=C=1$ $C=1 \rightarrow 0$

oder: bei $S=R=C=1$ $S=1 \rightarrow 0$ und $R=1 \rightarrow 0$

oder: zu kurze Pulsdauer bei S, R oder C

1.5.2.3. D-Latch

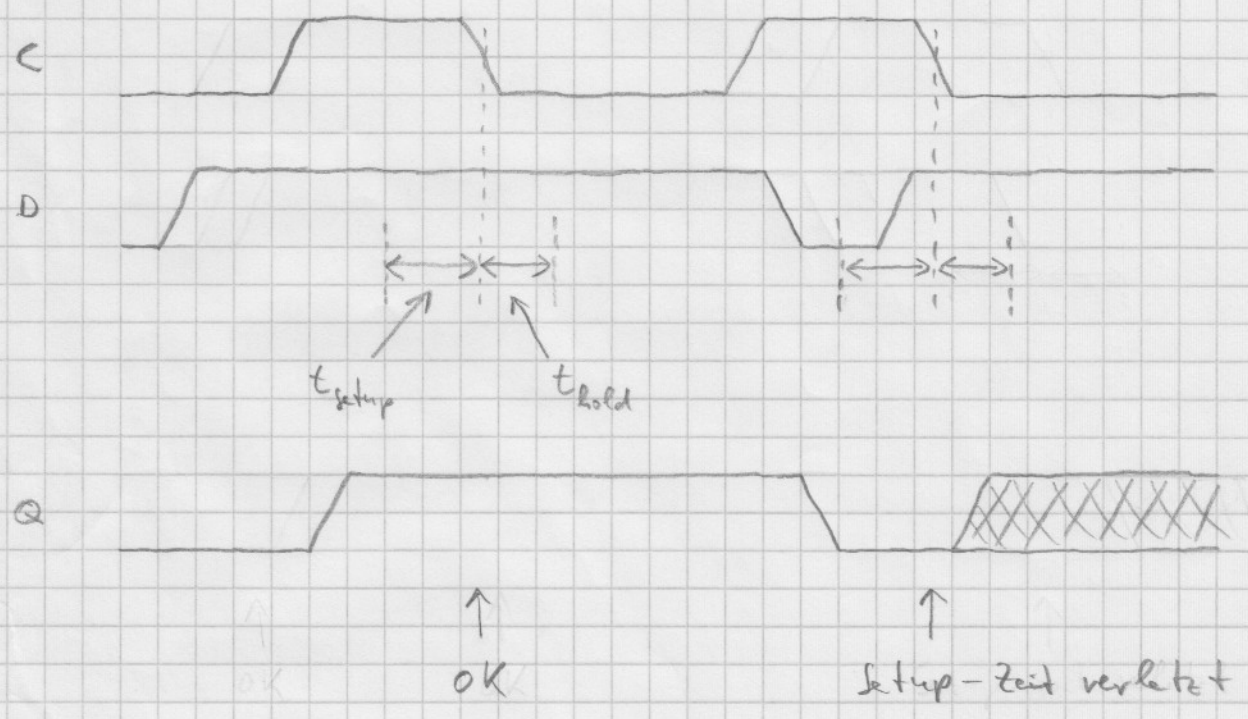
Speichert ein Datum (1 Bit). Dazu wird $S = \bar{R} = D$ gesetzt:



C	D	Q	QN
1	0	0	1
1	1	1	0
0	x	letztes Q	letztes QN

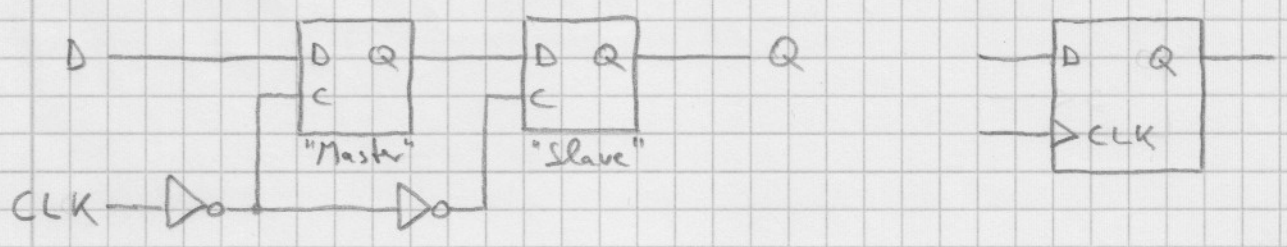
} Latch ist offen ("transparent")
 } Latch ist geschlossen ("hält")

Metastabilität: kritisch ist die fallende Flanke von C
 In einem best. "Fenster" um diese Flanke herum darf sich der Eingang D nicht ändern, sonst droht Metastabilität!



1.5.2.4. Taktflankengesteuertes ("edge triggered") D-Flip-Flop

Schaltung, die den D-Eingang nur auf einer Flanke des Takt-Eingangs speichert. Besteht aus zwei D-Latches, die mit gegenphasigen Takt angesteuert werden:



Wirkungsweise: Während $CLK = \phi$ ist der Master transparent; sein Ausgang folgt dem Daten Eingang. Der Slave hält aber den vorigen Wert. Wird $CLK = 1$, dann schließt der Master und der Slave übernimmt diesen Wert. Während der ganzen Zeit mit $CLK = 1$ folgt der Slave dem Master; dessen Wert ändert sich aber nicht mehr, weil er geschlossen ist.

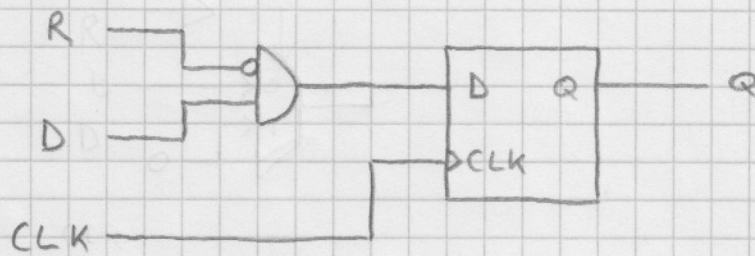
D	CLK	Q	
0		0	<u>In Verilog:</u> always @(posedge CLK) begin Q <= D; end
1		1	
x	0	letztes Q	
x	1	letztes Q	

Metastabilität: Genau wie beim D-Latch sind Setup- und Holdzeiten in Bezug auf die aktive Flanke des Taktsignals () zu beachten.

Bem.: D-Flip-Flops können mit zusätzlichen asynchronen Setz- und Rücksetzeingängen ausgestattet sein, die es erlauben, das Flip-Flop unabhängig vom Takt auf 1 oder 0 zu setzen ("preset" und "clear"). Diese sollten in synchronen Schaltungen nicht verwendet werden (außer bestenfalls zum Herstellen eines bekannten Startzustandes beim Reset — selbst da ist Vorsicht geboten).

1.5.2.5. D-Flip-Flop mit synchronem Reset

(19)

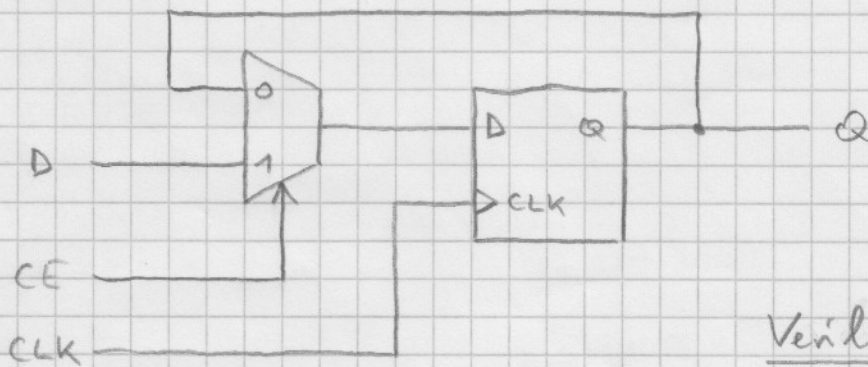


Verilog: siehe ①

1.5.2.6. D-Flip-Flop mit Taktfreigabe ("clock enable")

WICHTIG: Entgegen dem deskriptiven Namen "Taktfreigabe"

hat das mit dem Takt nicht das geringste zu tun !!



Verilog: siehe ②

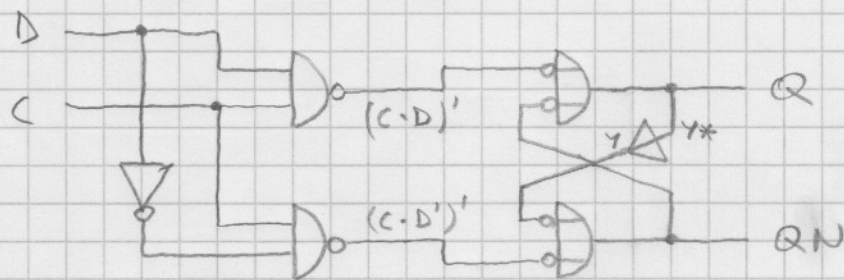
```
① always @(posedge CLK) begin
    if (R)
        Q <=  $\phi$ ;
    else
        Q <= D;
    end
```

```
② always @(posedge CLK) begin
    if (CE)
        Q <= D;
    end
```

1.5.3. Analyse sequentieller Schaltungen

(20)

Ist i. A. extrem anfrendig. Wir behandeln ein Bsp.:



1. Schritt: Aufbrechen der Rückkopplungen durch fiktive Puffer (hier ist ein Puffer erforderlich). Die Ausgänge der Puffer sind die "Zustandsvariablen" (hier: y). Die "nächsten" Werte der Zustandsvariablen werden mit einem "*" versehen.

Dann läßt sich der nächste Zustand aus dem jetzigen Zustand und den Eingangsvariablen ausdrücken ("Anregungsfunktion"):

$$\begin{aligned} y* &= (C \cdot D) + (y' + C \cdot D)' \\ &= C \cdot D + y \cdot C' + y \cdot D \end{aligned}$$

2. Schritt: Aufstellen der Übergangstabelle für y^* :

$y \backslash CD$	00	01	11	10
0	0	0	1	0
1	1	1	1	0

Im Allgemeinen hat eine solche Tabelle 2^n Zeilen für n Rückkopplungsschleifen und 2^m Spalten für m Eingänge.
Zunächst Annahme: Schaltung wird im "Grundmodus" betrieben,

d.h. es wird zu einem gegebenen Zeitpunkt nur ein (21)

Eingang geändert und die Schaltung hat danach genügend Zeit, in einen stabilen Zustand überzugehen.

"Totaler Zustand": Kombination aus "internem Zustand"

(= Werte in den Rückkopplungsschleifen) und "Eingangs Zustand"

(= Werte der Eingangsvariablen). "Stabiler totaler Zustand":

Kombination so, dass nächster interner Zustand gleich jetzigem internem Zustand ist.

3. Schritt: Aufstellen der Zustands- und Ausgangstabelle

Wir benennen die internen Zustände (hier: S_0, S_1), markieren die stabilen totalen Zustände und schreiben die

Werte der Ausgangsgrößen mit in eine Tabelle:

$$Q = C \cdot D + C' \cdot Y + D \cdot Y$$

$$QN = C \cdot D' + Y'$$

CD	00	01	11	10
S_0	($S_0, 01$)	($S_0, 01$)	$S_1, 11$	($S_0, 01$)
S_1	($S_1, 10$)	($S_1, 10$)	($S_1, 10$)	$S_0, 01$

Nun kann man Übergänge untersuchen, z.B. $S_0/00, D \rightarrow 1$.

Dann gelangt man zum stabilen totalen Zustand $S_0/01$, d.h.

interner Zustand (und Ausgänge) bleiben gleich. Wechselt jetzt

$C \rightarrow 1$, kommt man in den instabilen totalen Zustand $S_0/11$

und von dort in den stabilen totalen Zustand $S_1/11$ (Ausgänge: 10).

4. Schritt: Analyse des Verhaltens bei "gleichzeitiger" (22)

Veränderung an den Eingängen. Wirklich gleichzeitig gibt's nicht - wir müssen bei m "gleichzeitig" wechselnden Eingängen jede der $m!$ verschiedenen Reihenfolgen gesondert betrachten! Bsp.: Start in $S1/11$, $CD \rightarrow 00$

$S \backslash CD$	00	01	11	10
S0	(S0, 01)	(S0, 01)	S1, 11	(S0, 01)
S1	(S1, 10)	(S1, 10)	(S1, 10)	S0, 01

Wenn C zuerst wechselt: $S1/11 \rightarrow S1/01 \rightarrow S1/00$

Wenn D zuerst wechselt: $S1/11 \rightarrow S1/10 \rightarrow S0/10$
 $\rightarrow S0/00$

\Rightarrow Zustand unbekannt!

Aber: Nicht jeder "gleichzeitige" Wechsel von Eingangssignalen muss unvorhersagbares Verhalten zur Folge haben!

Bsp.: $S0/00$, $CD \rightarrow 11$

Beide Übergangssequenzen enden bei $S1/11$.

Fazit: Beliebige sequentielle Schaltungen sind nur mit enormem Aufwand zu analysieren, und noch schwieriger zu konstruieren! Ausweg: Synchroner Schaltwerke, in denen die sequentiellen Schaltungen auf Speicherglieder beschränkt werden, die gut verstanden sind (D-Latches bzw. D-Flip-Flops).

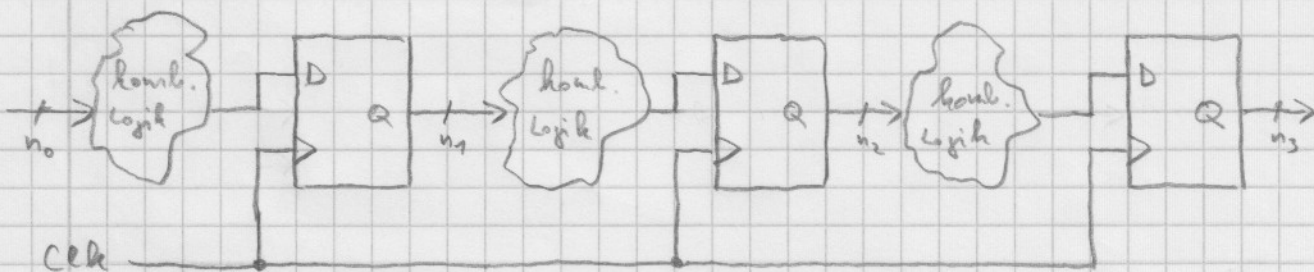
1.6. Synchroner Schaltwerke und Automaten

(23)

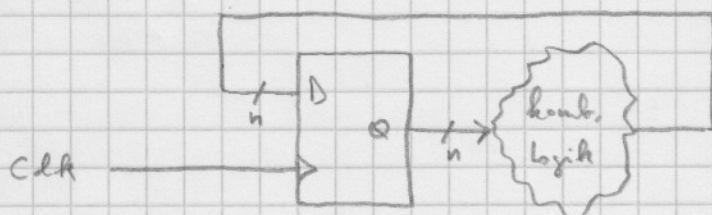
1.6.1. Prinzip und Realisierung

Prinzip: Sequentielle Schaltkreise sind nur in speichernden Elementen zu finden. Diese sind alle von derselben Sorte (z.B. positiv-flankengetriggerte D-Flip-Flops) und werden alle vom selben Takt angesteuert.

Vorteil: Das Einhalten der Setup-Zeiten kann immer durch Senken der Taktfrequenz erreicht werden; Hazards und Durchlaufzeiten der kombinatorischen Schaltnetze wirken sich höchstens begrenzend auf die Taktfrequenz aus, machen die Schaltung aber nicht funktionsunfähig!



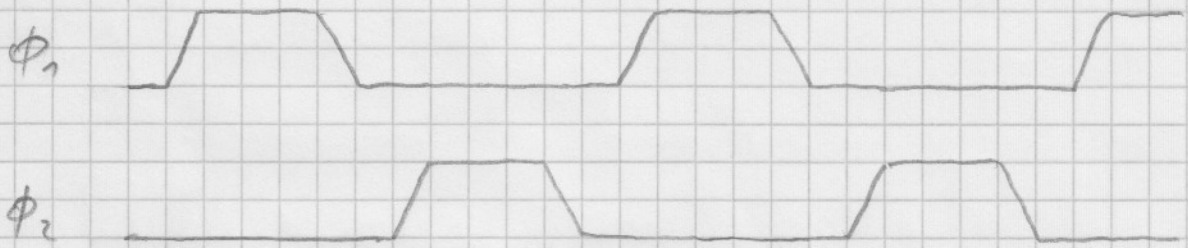
Beacht: Jedes Speicherelement wird gleichzeitig gelesen und (mit der Taktflanke) geschrieben! Also ist auch folgende Anordnung zulässig ("synchroner Automat"):



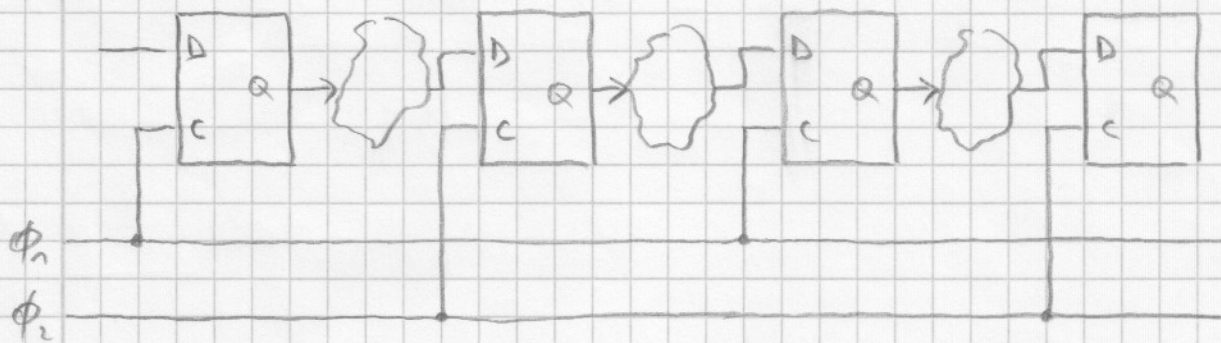
Der Wert von Q heißt "Zustand", der Wert, der gleichzeitig an D anliegt "Folgezustand" des Automaten.

Die komb. Logik berechnet also den Folgezustand aus (24) dem Zustand.

Bem.: Man kann anstelle von D-Flip-Flops auch D-Latches benutzen, muß allerdings die oben gezeigte Schaltung modifizieren, da sonst bei $CLK = 1$ (Latches sind transparent) kein Zustand gehalten werden könnte. Abhilfe schafft ein "nichtüberlappendes Zweiphasen-Takt":



Dann müssen die speichernden Elemente abwechselnd mit ϕ_1 und ϕ_2 getaktet werden:



Beachte: Wenn man jedes zweite Schaltelement weglässt und die beiden Latches zusammenfaßt, erhält man wieder die D-Flip-Flop-Variante.

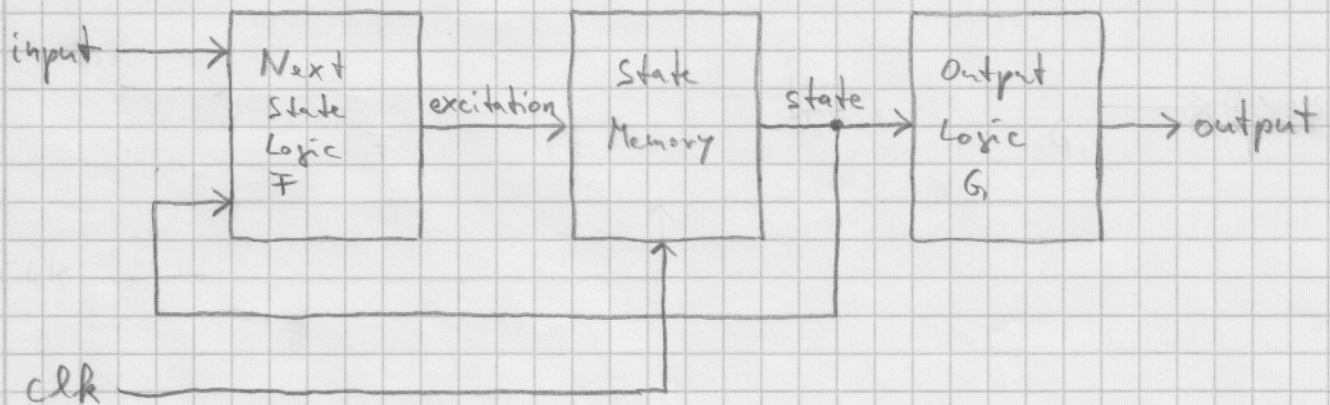
1.6.2. Moore-Automat

(25)

Im Allgemeinen wird ein Automat auch Eingänge besitzen, mit denen er etwas über seine Umgebung erfährt, und Ausgänge, mit denen er seine Umgebung beeinflusst.

Je nachdem, was die Eingänge steuern bzw. wie die Ausgänge berechnet werden, unterscheidet man verschiedene Automaten.

Beim Moore-Automaten hängt die Ausgabe nur vom Zustand ab, nicht von der Eingabe:



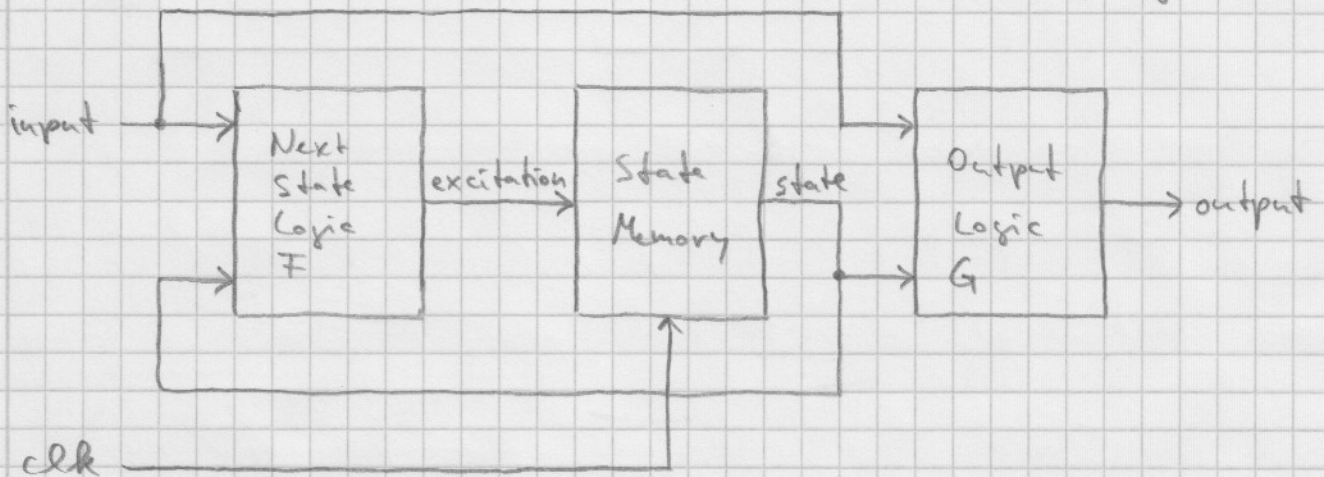
In Verilog:

```
assign output = G(state);
always @(posedge clk) begin
    case (state)
        :
        n: begin
            if (input == ...) begin
                state <= m; // excitation = F(state, input)
            end else
                :
            end
        end
    endcase
end
```

1.6.3. Mealy-Automat

(26)

Hier hängt die Ausgabe (asynchron) auch von der Eingabe ab:

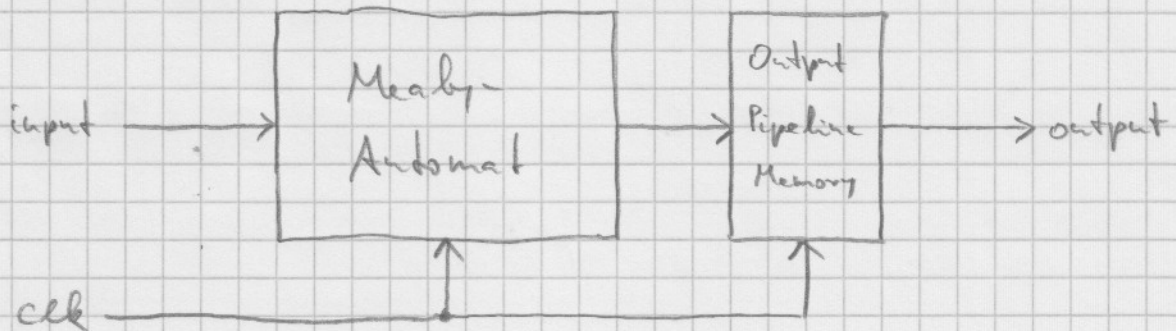


↳ Verilog: Wie beim Moore-Automaten, aber mit
 $assign\ output = G(state, input);$

Vergleich Moore/Mealy: Moore-Automaten sind potentiell schneller als Mealy-Automaten, haben aber möglicherweise mehr Zustände als diese.

1.6.4. Mealy-Automat mit synchronem ("Pipeline"-) Ausgang

Dieser Automat entsteht aus dem Mealy-Automaten durch "Anhängen" einer Speicherstufe an die Ausgänge:



In Verilog:

```

always @(posedge clk) begin
  case (state)
  :
  n: begin
    if (input == ...) begin
      state <= m;
      output <= x;
    end else
    :
  end
  :
endcase
end

```

Bem.: Diese Variante ist sehr schnell, man muß aber die Ausgangswerte schon einen Taktzyklus früher berechnen können.

(Persönliche Anmerkung: Achtung, Hirnsausen!)

① →

1.6.6. Asynchrone Eingänge und Synchronisierer

Nicht alle Teilsysteme eines Rechners werden mit dem gleichen Takt getaktet (wie drückt man synchron zu 16Hz eine Taste?).

Man benötigt "Synchronisierer", die ein asynchrones Eingangssignal synchron zu einem gegeb. Takt abtasten. Es ist normalerweise nicht schlimm, wenn der komplette Wert erst einen Takt später übernommen wird, aber es gibt zwei Fehler, die man hier leicht macht:

1. Synchronisation des gleichen Signals an mehreren Stellen

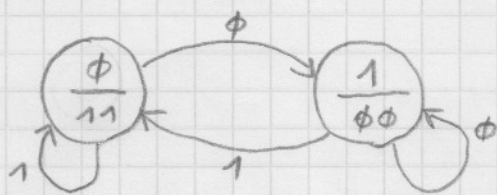
1

1.6.5. Zur Konstruktion von Automaten ("FSMs")

27a

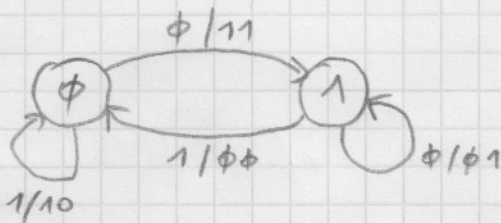
1. Stellen die Ausgaben an den Zuständen: Moore-Automat

Bsp.:



2. Stellen die Ausgaben an den Übergängen: Mealy-Automat

Bsp.



3. Häufig legt man beide Sorten Ausgänge gleichzeitig (formal ist das dann ein Mealy-Automat).

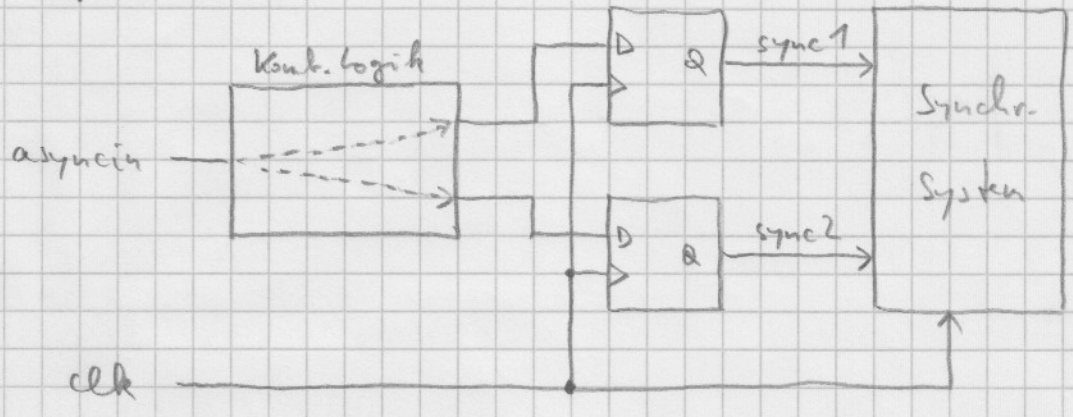
4. Minimierung der Zustände: ist machbar, aber ein vernünftiges Design macht das i.d.R. überflüssig.

5. Optimale Kodierung der Zustände: ist abhängig von der verwendeten Technologie \rightarrow die Synthesewerkzeuge können das erledigen. Tatsächlich wird häufig eine "One-Hot-Kodierung" gewählt: n Zustände werden durch n Bits (nicht $\log n$ Bits!) repräsentiert, von denen genau eines 1, der Rest ϕ ist.

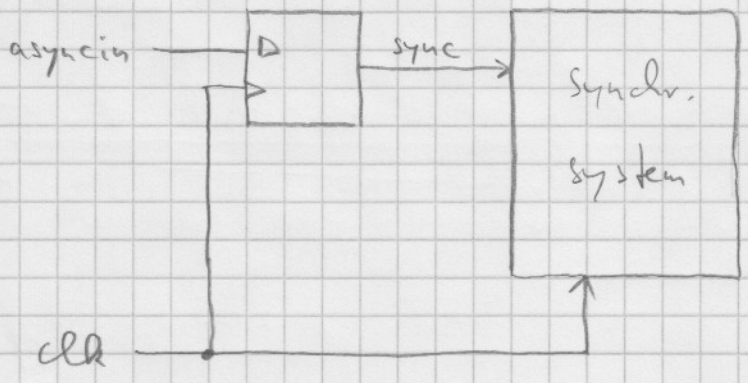
6. Niemals vergessen: FSMs müssen einen Reset-Eingang haben (sonst funktioniert z.B. One-Hot-Kodierung nicht)!

7. Ein gutes Design sieht Übergänge von "unermüdlichen" Zuständen zu einem Fehlerzustand oder zum Resetzustand vor!

Das Problem tritt insbesondere auf, wenn die Mehrfachsynchroisation durch kombinatorische Logik verdeckt wird:



Hier ist es sehr wahrscheinlich, dass bei einer Änderung von asyncein das eine Flip-Flop den alten Wert aus der komb. Logik, das andere aber den bereits geänderten Wert erkennt.
 Richtige Vorgehensweise: zuerst (vor irgendwelcher Logik) und nur an einer Stelle synchronisieren!



2. Unterschätzen der Gefahr von Metastabilität

Beim Erhöhen der Taktrate steigt die Gefahr von Metastabilität exponentiell an! Deshalb sollten zwei (in Spezialfällen drei) Synchronisierungsstufen vorgesehen werden:



2. Rechnerarchitektur

Die Rechnerarchitektur versucht, Gesetze zu finden und Strukturen zu beschreiben, die zum Bau von Rechnern mit hoher Leistung ("Performance") führen.

2.1. Performance

Zwei gebräuchliche Metriken:

Als einzelner Benutzer ist man interessiert an der "Antwortzeit".

Als Computer-Center-Manager ist man am "Durchsatz" interessiert.

Verkürzung der Antwortzeit verbessert i. A. auch den Durchsatz.

Im Folgenden wird die Antwortzeit ("Execution time") zur Performance in Bezug gesetzt:

$$\text{Performance} = \frac{1}{\text{Execution time}}$$

~~Beim Messen der Ausführungszeit ist zu unterscheiden zwischen "User CPU time", "System CPU time" und "Elapsed time". Hier geht's um die "User CPU time".~~

2.1.1. Die Performance-Gleichung

$$\text{Execution time} = \frac{\text{Clock cycles}}{\text{Program}} \cdot \frac{\text{Seconds}}{\text{clock cycle}}$$

$$= \frac{\text{Instructions}}{\text{Program}} \cdot \frac{\text{clock cycles}}{\text{Instruction}} \cdot \frac{\text{Seconds}}{\text{clock cycle}}$$

↑
"CPI", dient zum Vergleich verschiedener Implementierungen des gleichen Replissates

Der CPI ist ein geschicktes Mittel:

$$CPI = \sum_{i \in \text{Instrs}} CPI_i \cdot p_i$$

(p_i : rel. dyn. Häufigkeit der Instr. i)

Bsp.: Maschine A: clock cycle time = 1 nsec, CPI = 2.0

Maschine B: " " " = 2 nsec, CPI = 1.2

Wie ist das Performance-Verhältnis der beiden Maschinen?

Die Zahl der im Progr. ausgeführten Instruktionen sei n .

$$\text{Execution time}_A = n \cdot 2.0 \cdot 1 \text{ nsec} = 2n \text{ nsec}$$

$$\text{Execution time}_B = n \cdot 1.2 \cdot 2 \text{ nsec} = 2.4n \text{ nsec}$$

$$\Rightarrow \frac{\text{Performance}_A}{\text{Performance}_B} = \frac{\text{Ex. time}_B}{\text{Ex. time}_A} = \frac{2.4n}{2n} = 1.2$$

Bem.: Komplexere Instruktionen führen zwar zu einer kleineren Anzahl von ausgeführten Instruktionen, erhöhen aber i.d. Regel die clock cycle time oder den CPI.

\Rightarrow Beim Vergleich von Rechnern alle drei Faktoren bedenken!

2.1.2. Programme, um Performance zu messen

Sog. "Benchmarks" dienen zum Vergleich von Rechnern.

Fazit aus vielen Jahren mit fehlerhaft entworfenen, geschicht ausgeklügelten, etc, etc, Benchmarks lassen nur einen Schluss zu: Die zu erwartende Performance muß an realen Programmen, am besten der eigenen Application, gemessen werden.

2.1.3. Amdahl's Gesetz

Trugschluß: Man erwartet, daß eine Beschleunigung eines Aspekts einer Maschine die Performance darüber um einen Betrag proportional zur Beschleunigung steigert.

Richtig ist: Execution time (mit Beschleunigung)

$$= \left(\frac{\text{Execution time (betroffen von Beschleunigung)}}{\text{Faktor der Beschleunigung}} \right)$$

+ Execution time (nicht betroffen von Beschleunigung)

Bsp.: Ein Programm läuft 100 sec auf einer best. Maschine, davon werden 80 sec für Multiplikationen verbraucht.

Um welchen Faktor müssen Multiplikationen beschleunigt werden, so daß das Programm danach 5 mal so schnell läuft?

$$\text{Exec. time (mit Besch.)} = \frac{80 \text{ sec}}{n} + 20 \text{ sec} \stackrel{!}{=} \frac{100 \text{ sec}}{5}$$

$$\Rightarrow \frac{80 \text{ sec}}{n} = 0 \Rightarrow \text{keine (endl.) Beschleunigung der}$$

Multiplikationen kann eine 5-fache Performance der Maschine bewirken, wenn Multiplikationen nur 80% der ausgeführten Instruktionen ausmachen!

$$\Rightarrow \text{Amdahl's Gesetz: "Speedup"} = \frac{\text{Performance mit Besch.}}{\text{Performance ohne Besch.}}$$

$$= \frac{\text{Exec. time (ohne)}}{\text{Exec. time (mit)}} = \frac{1}{(1-p) + \frac{p}{s}}$$

$$\text{mit } p = \frac{\text{Exec. time (betroffen v. Besch.)}}{\text{Exec. time (ohne Besch.)}}, \quad s = \text{Faktor der Besch.}$$

2.2. Instruktionssatz-Architektur

32

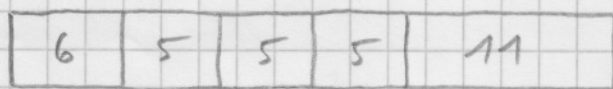
Das ist eine der wichtigsten Abstraktionen im Rechner: das Interface zwischen Hardware und Software. Für ein- und dieselbe ISA sind viele Implementierungen möglich, mit sehr verschiedener Performance (und Kosten).

Bei uns: "ECO32e", eine Untermenge des MIPS-Befehlssatzes (mit anderer Kodierung), 34 Befehle, jeder belegt genau 1 Wort.

2.2.1. Befehlsformate

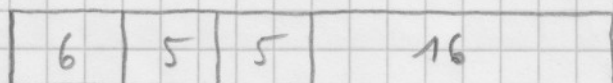
Beschreiben die Einteilung der Befehle in Klassen und die entspr. Zuordnung von Bits in den Befehlen zu "Feldern", die zur Kodierung der Befehlsbestandteile dienen.

R-Format:



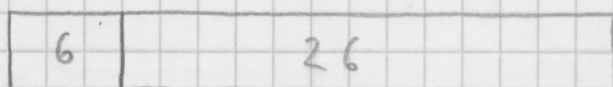
opcode rs1 rs2 rd -unbenutzt-

I-Format



opcode rs1 rd immediate

J-Format



opcode immediate

Festlegung des Formats durch den Opcode!

Bem.: Es werden verschiedene kleinere Variationen

(33)

der Hauptformate benutzt, z.B.:

- bei Arithmetik-Befehlen wird der Immediate-Operand vorzeichenrichtig aufgeweitet ("sign-extended"), bei Logik-Befehlen aber mit Nullen aufgefüllt ("zero-extended")
- zum Berechnen einer Datenadresse wird der Immediate-Operand so verwendet, wie er vorliegt, zum Berechnen einer Instruktionsadresse (bei Sprüngen) aber erst um zwei Bits nach links geschoben (da Instruktionen nur an durch 4 teilbaren Adressen liegen)

2.2.2. Befehlsatz

Im E032e haben wir Befehle für

- Arithmetik (+, -)
- Logik (&, |, xor, xnor)
- bedingte Sprünge ("branches")
- unbedingte Sprünge ("jumps")
- Unterprogramm Sprünge ("jump and link")
- Datentransfer vom ("load") und zum ("store") Speicher in den drei Transferbreiten Word, Halfword, Byte

Austauschen "E032e Architecture"

2.3. Implementierung des Befehlssatzes

(34)

2.3.1. Strategien

Es gibt drei Strategien, um einen gegebenen Befehlssatz zu realisieren:

a) Single-Cycle: jede Instruktion benötigt genau einen Taktzyklus zur Ausführung

Nachteile: Die Zyklusdauer muss sich nach dem kompliziertesten (längsten) Befehl richten; Hardware (wie z.B. Addierer) kann nicht für verschiedene Zwecke ($PC \leftarrow PC+4$, $RES \leftarrow OP1+OP2$) innerhalb der gleichen Instruktion wiederverwendet werden, sondern muss mehrfach vorhanden sein; getrennte Speicher für Instruktionen und Daten erforderlich.

Abklärung des längsten Pfades (komplizierteste Instruktion ist ldw):

- Zugriff auf Instruktionsspeicher
- Zugriff auf Register-File (Lesen)
- ALU (zur Berechnung der effektiven Adresse)
- Zugriff auf Datenspeicher
- Zugriff auf Register-File (Schreiben)

⇒ Obwohl ~~HW~~ $CPI = 1$ ist, ist die Performance schlecht, da die Länge, wenn ldw notwendige Taktzeit von vielen Befehlen nicht gebraucht wird!

Beim nächsten Ziel muss mit dem nächsten Befehl...

b) Multi-Cycle: jede Instruktion wird in eine unterschiedliche Anzahl etwa gleich langer Teilaufgaben aufgespalten. Die längste der Teilaufgaben bestimmt die Zyklusdauer; der CPI ist für versch. Befehle verschieden.
 Vorteile: Ein gemeinsamer Speicher für Instruktionen und Daten; Wiederverwendung von Hardware für versch. Teilaufgaben innerhalb einer Instruktion; gute Anpassung der Instruktionsdauer an die wirklich benötigte Zeit. Zusätzlich braucht man:

2) Ein paar wenige Multiplexer, die Eingangssignale an Baugruppen auswählen.

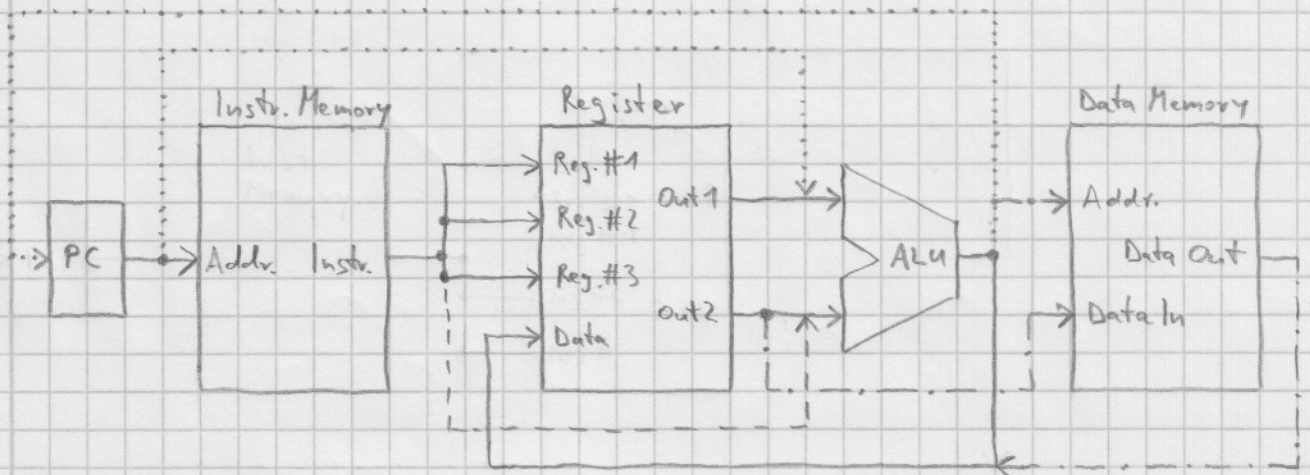
- 1) Ein paar wenige Register, die die Ergebnisse eines Teilschrittes einer Instruktion an die nächsten Teilschritte weitergeben.
- 3) Einen Automaten ("Steuerwerk"), der die Einzelschritte für jede Instruktion ablaufen lässt und dabei Steuersignale für den "Datenpfad" erzeugt.

c) Pipeline: Das ist eine Kombination aus a) und b). Jede Instruktion wird in eine Anzahl gleich langer Teilaufgaben zerlegt, aber jede Teilaufgabe erhält eigene Hardware. Also können bei n Teilschritten n verschiedene Teile der Hardware ("Pipeline-Stufen") parallel arbeiten - allerdings auf n unterschiedlichen Instruktionen. Damit wird zwar die Ausführungszeit einer Instruktion nicht verringert, aber die "Durchsatz" von Instruktionen (= ausgeführte Instruktionen pro Zeiteinheit) auf das n -fache erhöht (besterfalls).

2.3.2. Eine Multi-Cycle-Implementierung

36

2.3.2.1. Abstraktes Blockschaltbild

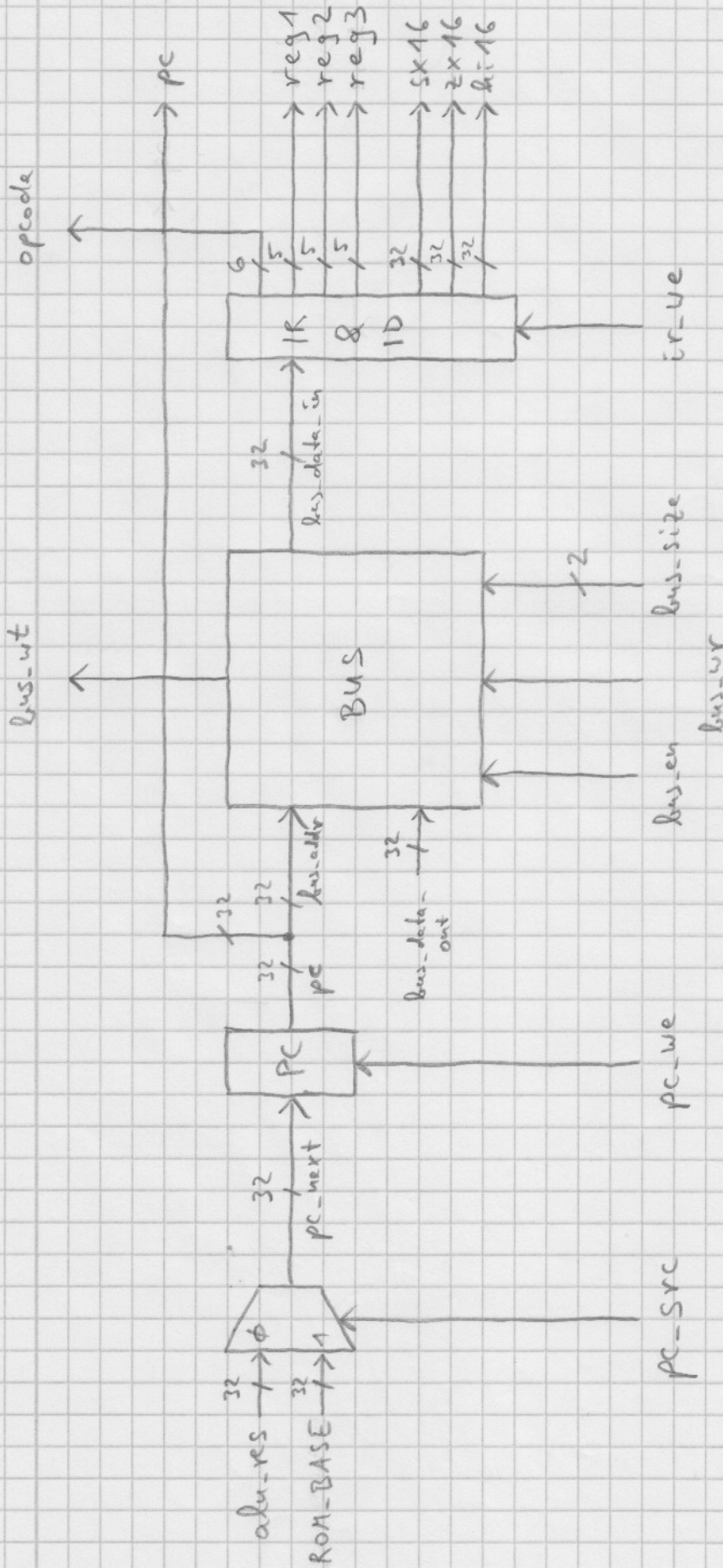


- Register-Instructionen
- - - Immediate-Instructionen
- Sprünge
- - - Load/Store-Instructionen

2.3.2.2. Wo werden "Auffangregister" platziert?

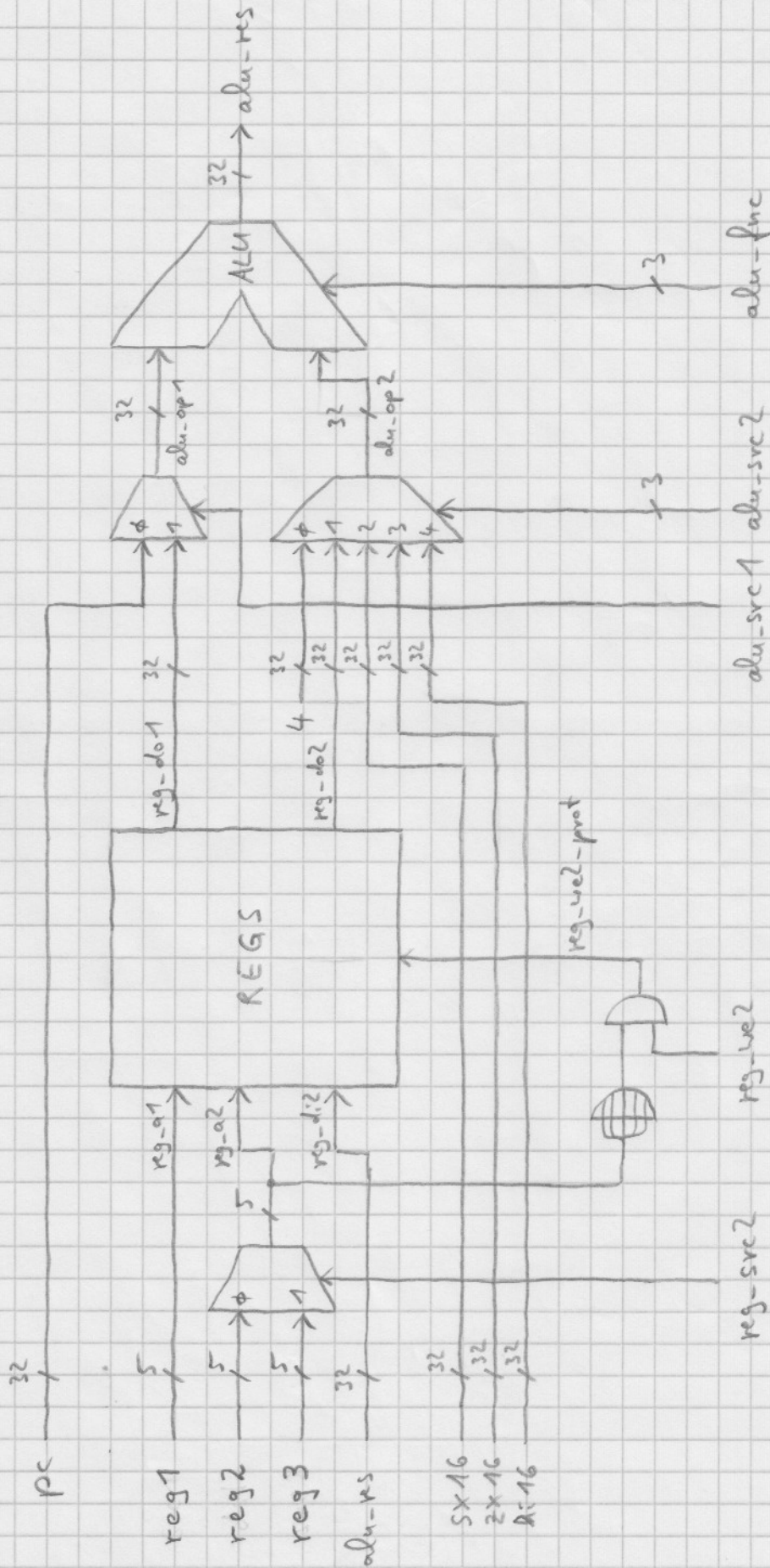
Das richtet sich ganz danach, in welche Teilschritte man die Instructionen zerlegt. Dies wiederum hängt von der Durchlaufzeit der Signale durch die Teilerke und von der beabsichtigten Taktrate ab. Eine grobe Abschätzung für unseren FPGA ergibt: Bei 50 MHz hat man recht viel Zeit, kann also viel in einem Taktzyklus erledigen und braucht nur wenige Register. Das wird unterstützt durch die Tatsache, daß das Block-RAM zur Realisierung des Register-Files getaktet ist. Eine Ausnahme bildet der Bus - der kann langsam sein. Also werden hier Auffangregister benötigt.

2.3.2.3. Datenpfad für Instruction-Fetch/Decode



- SX 16 = sign-extend (16-bit-immediate)
- EX 16 = zero-extend (16-bit-immediate)
- hi 16 = { 16-bit-immediate, 16'0000 }

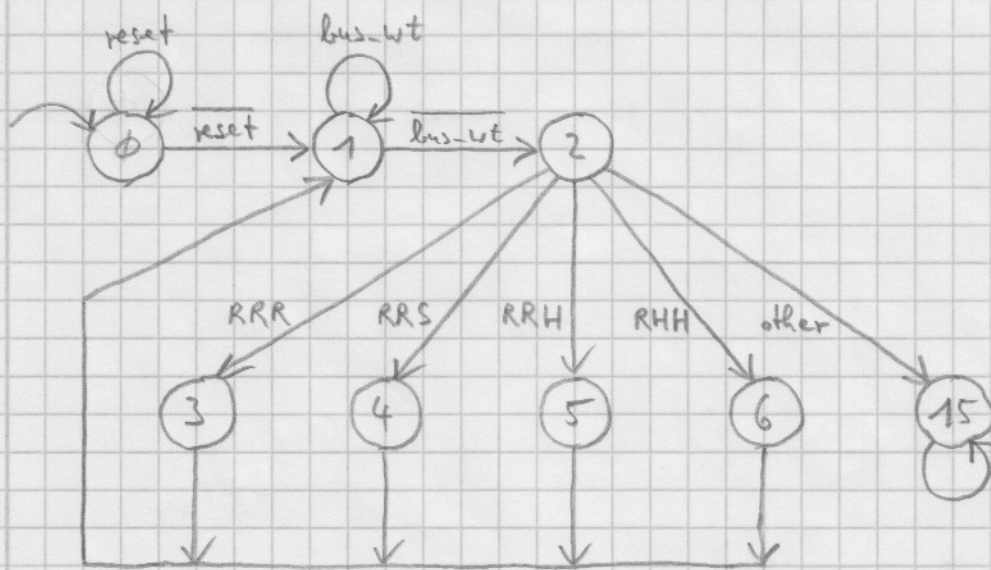
2.3.2.4. Datenpfad für Register- und Immediate-Befehle



2.3.2.5. Steuerwerk für Register- und Immediate-Befehle (39)

Synchroner Automat mit den Zuständen:

- ϕ Reset
- 1 Fetch Instruction
- 2 Decode Instruction & Increment PC
- 3 Execute RRR-Inst (z.B. add \$5, \$6, \$7)
- 4 Execute RRS-Inst (z.B. add \$5, \$6, 12)
- 5 Execute RRH-Inst (z.B. and \$5, \$6, 12)
- 6 Execute RHH-Inst (z.B. lddi \$5, $\phi \times 8 \phi \phi \phi \phi \phi \phi \phi \phi$)
- 15 Error

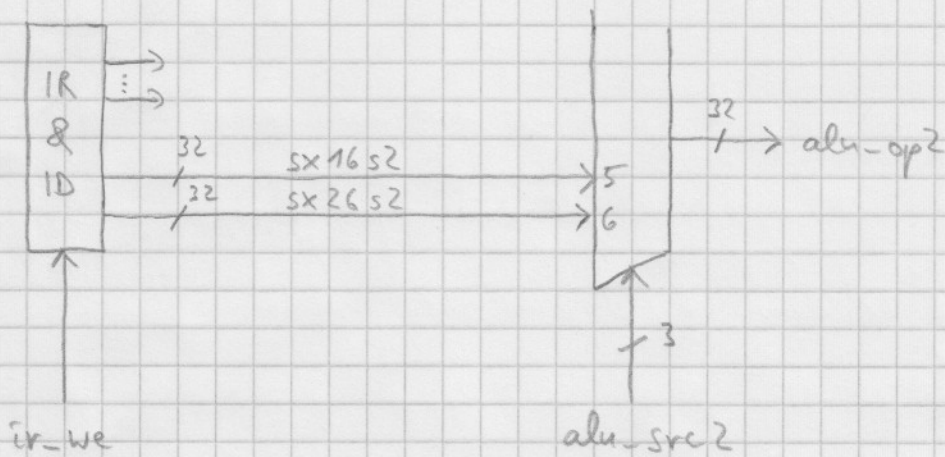


Steuerwerk ist Moore-Automat: Kontrollsignale für den Datenpfad werden ^(kur) aus dem Zustand erzeugt.

Dabei rechnet man formal das Instruktionsregister mit zum Zustand des Automaten.

2.3.2.6. Datenpfad für Sprünge

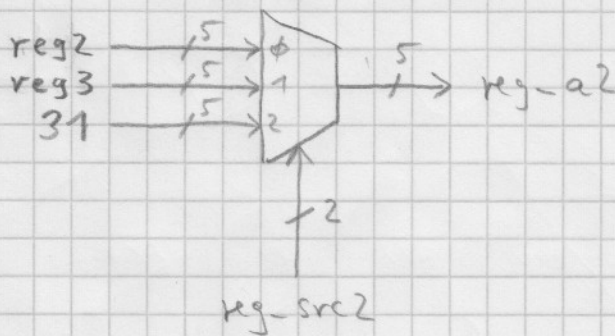
a) Ergänzungen am Instruktions-Dekoder / ALU-Operand2-Selektor



$sx\ 16\ s2 = \{ \text{sign-extend (16-bit-immediate)}, 2'b\ \phi\phi \}$

$sx\ 26\ s2 = \{ \text{sign-extend (26-bit-immediate)}, 2'b\ \phi\phi \}$

b) Ergänzung zum Speichern der Return-Adresse



c) Berechnung der Entscheidungsgrundlage für bedingte Sprünge

Benötigte Information: $a=b$, $a <_u b$ ($<_u$: vorzeichenloser Vergleich)

Alle anderen Sprünge lassen sich daraus ableiten. (*)

1) $a=b \Leftrightarrow a-b = \phi$: ALU subtrahiert, dann NOR(result)

2) $a <_u b \Leftrightarrow a-b <_u \phi$ (???) \downarrow

Kleiner Trick: Fasse die vorzeichenlosen 32-Bit-Zahlen auf als positive vorzeichenbehaftete 33-Bit-Zahlen!

(*) Aber Achtung: $a > b \Leftrightarrow b < a \Leftrightarrow b \leq a \Leftrightarrow r(b > a)$

$a' = \phi a_{31} \dots a_0$, $b' = \phi b_{31} \dots b_0$

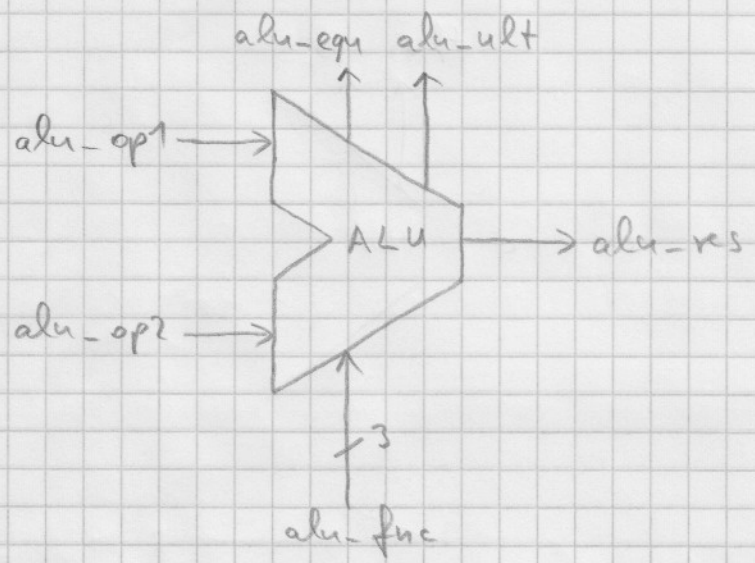
$a <_u b \Leftrightarrow a' < b' \Leftrightarrow a' - b' < \phi$

Die Differenz ist in 33 Bits darstellbar:

$0 \leq a \leq 2^{32} - 1$, $0 \leq b \leq 2^{32} - 1 \Rightarrow$

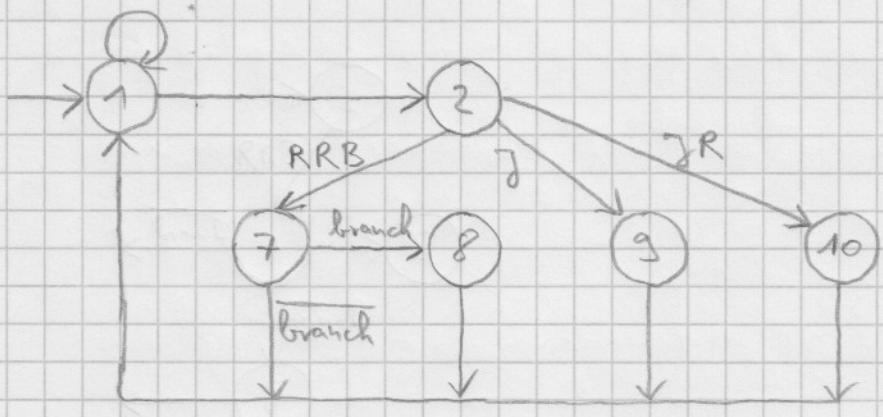
$-(2^{32} - 1) \leq a - b \leq 2^{32} - 1$

Also: $a <_u b \Leftrightarrow \text{MSB}(a' - b') = 1$



2.3.2.7. Steuerwerk für Sprünge

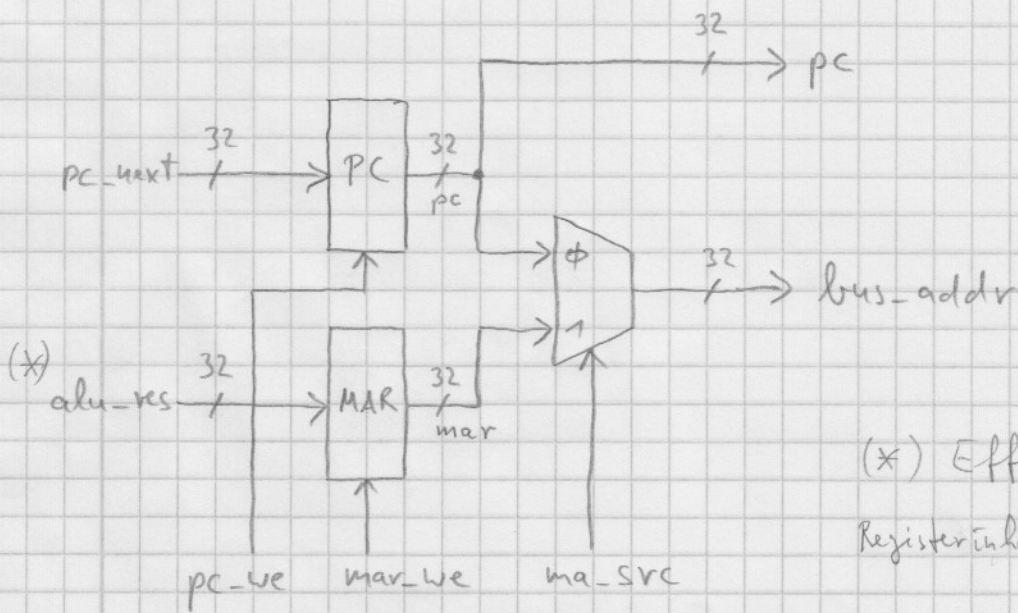
- Zustände: 2 Spricht ggf. PC+4 in §31
- 7 Execute RRB-Inst (Subtraktion)
- 8 Execute RRB-Inst (neuer PC)
- 9 Execute J-Inst
- 10 Execute JR-Inst



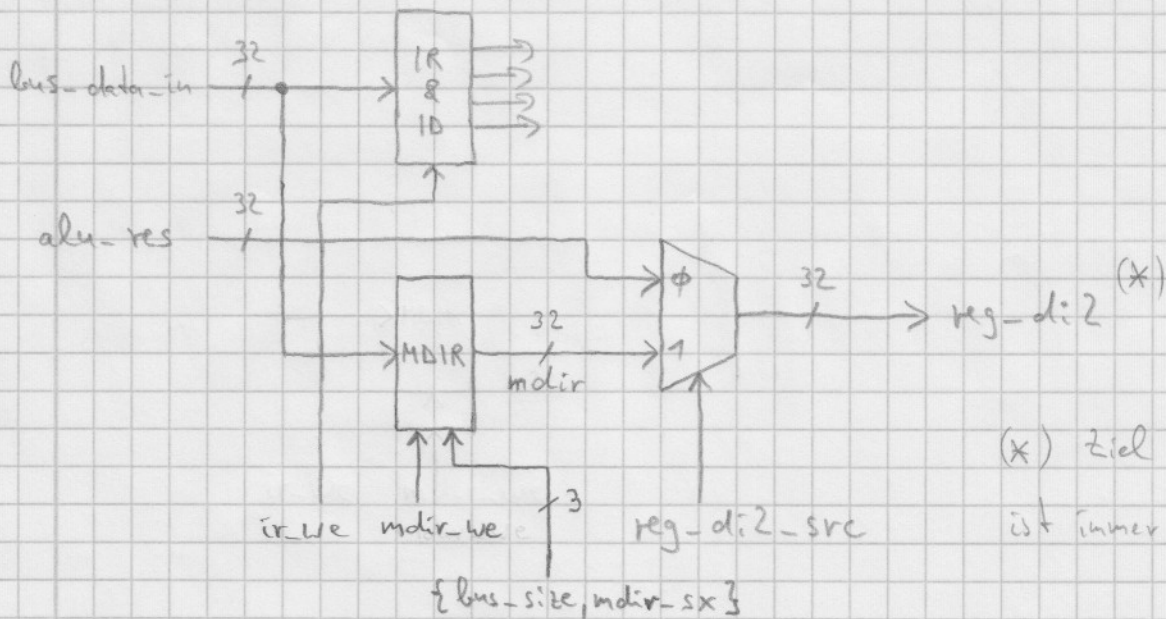
2.3.2.8. Datenpfad für Load/Store-Operationen

(42)

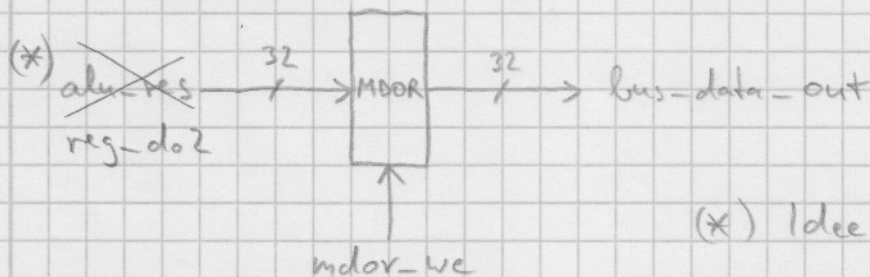
a) Bus-Adresse



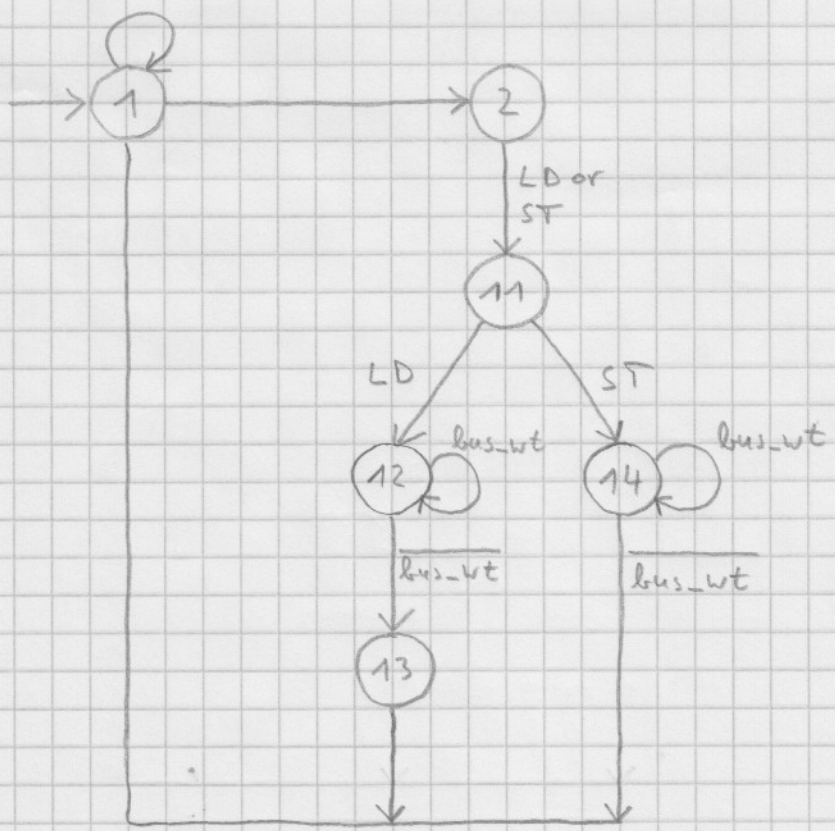
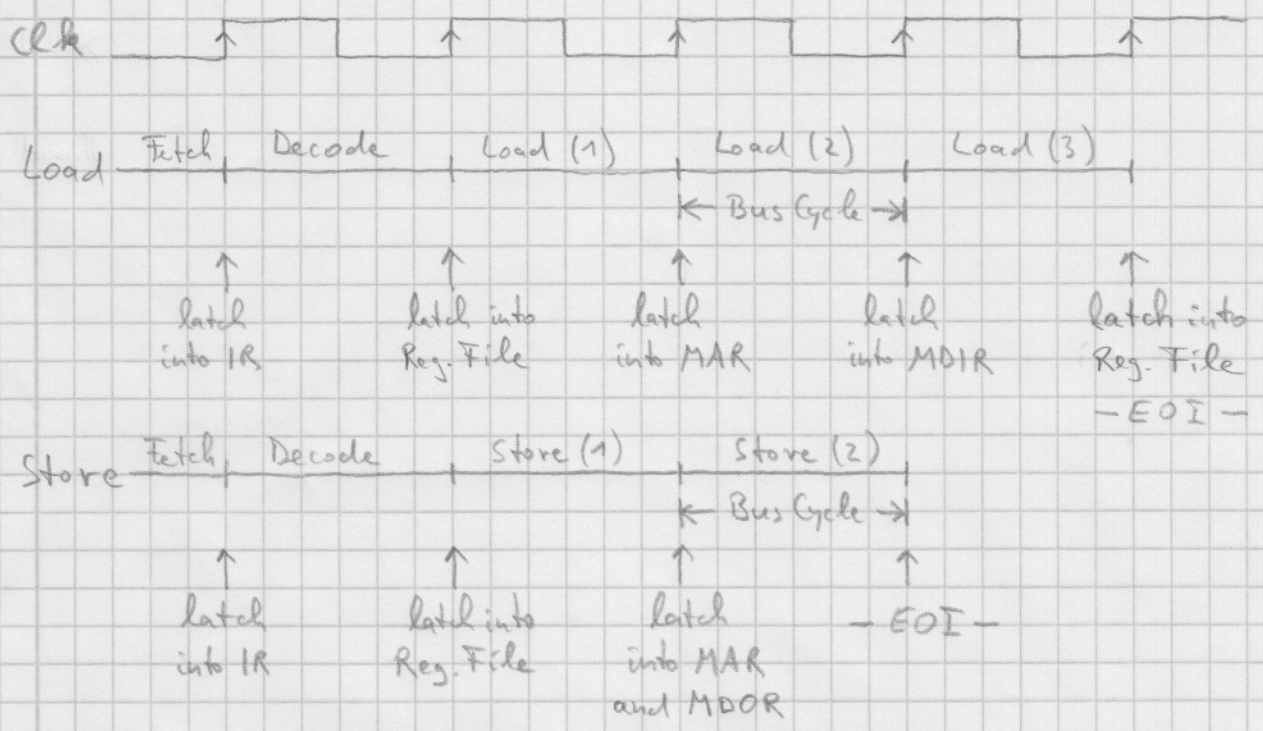
b) Daten: Bus → CPU



c) Daten: CPU → Bus



2.3.2.9. Steuerwerk für Load/Store-Operationen



2.3.3. Interrupts und Exceptions

(44)

Interrupt: Unterbrechen des Kontrollflusses durch CPU-externes Ereignis; konzeptionell "zwischen" zwei Instruktionen (Bsp.: Plattenoperation beendet)

Exception: Unterbrechen des Kontrollflusses durch die momentan ausgeführte Instruktion (Bsp.: Division durch 0, illegaler Opcode, Zugriff auf geschützte Seite, TCB-Miss, etc). Achtung: die Instruktion darf nicht weiter ausgeführt werden!

Software zur Behandlung von Interrupts und Exceptions benötigt mindestens:

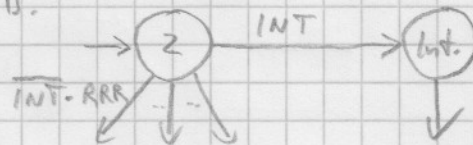
- die Adresse der betroffenen Instruktion
- die Ursache der Unterbrechung

→ 2 Spezial-Register, die von Software gelesen werden können

Handhabung der Unterbrechung in der Multi-Cycle-Impl.:

- Zykl. Abfragen des Interrupt-Eingangs im Steuerwerk,

z.B.



- Ergänzung des Steuerwerks um Zustände / Übergänge für Exceptions, z.B.:

