

Levels of abstraction

1. Physics: Motion of electrons
2. Devices: Transistors and Terminals (measure voltage/current)
3. Analog Circuits: In which devices are assembled (amplifiers)
4. Digital Circuits: Logic gates, restrict voltage
5. Logic: More complex structures (memory)
6. Microarchitecture: Describes hardware for programmer

Gates

-OR x v y x+y	\Rightarrow	x y	-NOR $\neg(x \vee y)$ $\overline{x+y}$	\Rightarrow	D $\neg(a \vee b)$
-AND x \wedge y $x \cdot y$	\Rightarrow	x \wedge y	-NAND $\neg(x \wedge y)$ $\overline{x \cdot y}$	\Rightarrow	D $\neg(a \wedge b)$
-NOT $\neg x$ \overline{x}	\Rightarrow	Buffer x			
-XOR $x \neq y$ $x \oplus y$	\Rightarrow	x \wedge y	{ 00 \rightarrow 0 10 \rightarrow 1 01 \rightarrow 1 11 \rightarrow 0 }	-XNOR $x = y$ $\overline{x \oplus y}$	\Rightarrow D $\neg(x \neq y)$

Number System

- Binary: 2^N (0, 1, 2, ..., 2^{N-1}) | $2^4 \Rightarrow 0, 1, 2, \dots, 15 \Leftrightarrow 0000, 0001, 0010, \dots, 1111$
- Hexadecimal: 16^N (0, 1, 2, ..., 9, A, B, ..., F) | $16^1 = 2^4$
- Bits Bytes (= Unsigned)
- 1 Byte $\hat{=} 8$ Bits $\hat{=} 2^8 \hat{=} 256$ possibilities (Unsigned: $[0/2^{N-1}]$)
- 1 Word (Handeled by Microprocessors) | 32/64 Bit words
- Bit in the 1's column = least significant bit > most significant bit

z. Darstellungen für:

- Sign/magnitude: most significant bit ± 0 : $-S_2: 1101 / +S_2: 0101 | [-2^{N-1} + 1 / 2^{N-1} - 1]$
- Two's complement: $[-2^{N-1} / 2^{N-1} - 1] | \text{zero} = 0 \dots 0 / \text{kleinste Zahl: } 1000 \dots 0 / \text{größte: } 0111 \dots 1$
 - Most significant bit ist sign-bit: 1 = negativ / 0 = positiv
 - Sign reversed: 1. invert all bits 2. add 1 to Lsb: $+2_{10} = 0010_2$
 - Addition/Subtraction: $\pm \text{Zahl} + \pm \text{Zahl}$ (immer addieren) $-2_{10} = 1110_2$
 - Wenn number is extended to more bits, sign bit to msb

Logic Levels: 0 - V_{OL} - V_{IL} - Forbidden zone - V_{IH} - V_{OH} - V_{OD} (V_{IL}/V_{IH} for noise margins)

Combinational Logic Blocks (have no memory / building blocks)

- Output depends only from the input \rightarrow has no memory
- Every circuit block itself is combinational and has no cyclic paths

Boolean equation variables are called literals

Minterm: Formel, wenn wir überall true einsetzen (And/V/I/.)

A	B	y	
0	1	1	$\bar{A} \cdot B \text{ Min}$
0	0	0	$A + B \text{ Max}$

Maxterm: Formel, wenn wir überall false einsetzen (OR/V/I/+)

Sum of product: Summing up each minterm, where output/y = true / y = $\bar{A}B + \dots + \dots$

Products of sum: Producing up each maxterm, where output/y = false / y = $(A+B) \cdot (B+C) \dots$

De Morgan: $y = \bar{A} \cdot B = \bar{A} + \bar{B} / y = \bar{A} + \bar{B} = \bar{A} \cdot \bar{B} \Rightarrow \text{D} \Leftrightarrow \text{D} \quad (= \text{Bubble pushing})$

From logic to gates: $\text{D} \Leftrightarrow \bar{A} \bar{B} \bar{C} + \text{No connection} \quad \text{D} \Leftrightarrow \text{D} \quad \text{D} \Leftrightarrow \text{D} \quad \text{D} \Leftrightarrow \text{D}$

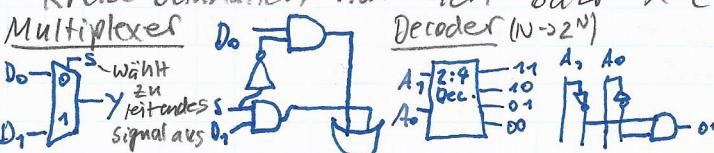
Illegal Values: X: unknown value / Z: Entweder 0/1

Karnaugh-Maps: (K-Maps) Eliminate complementary forms of variables, größer, besser

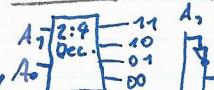
- Weniger Kreise, desto besser / Jeder Kreis $1/2/4/8/2^N$ lang/breit/berges

- Kreise beinhalten nur 1en oder X (= Don't care)

Multiplexer



Decoder ($N \rightarrow 2^N$)



A_1, A_0

S

D_0, D_1, D_2, D_3

$D_{0:3}$	00	01	10	11	$\bar{D}_2 \bar{D}_1 \bar{D}_0$
00	1	0	0	0	$\bar{D}_2 \bar{D}_1 \bar{D}_0$
01	0	1	0	0	$\bar{D}_3 \bar{D}_2 \bar{D}_1$
11	1	1	0	0	$\bar{D}_3 \bar{D}_2 \bar{D}_0$
10	1	1	1	0	$\bar{D}_3 \bar{D}_1 \bar{D}_0$

größer,
besser

größer,
besser

- Low \rightarrow High: Rising Edge / High \rightarrow Low: Falling Edge

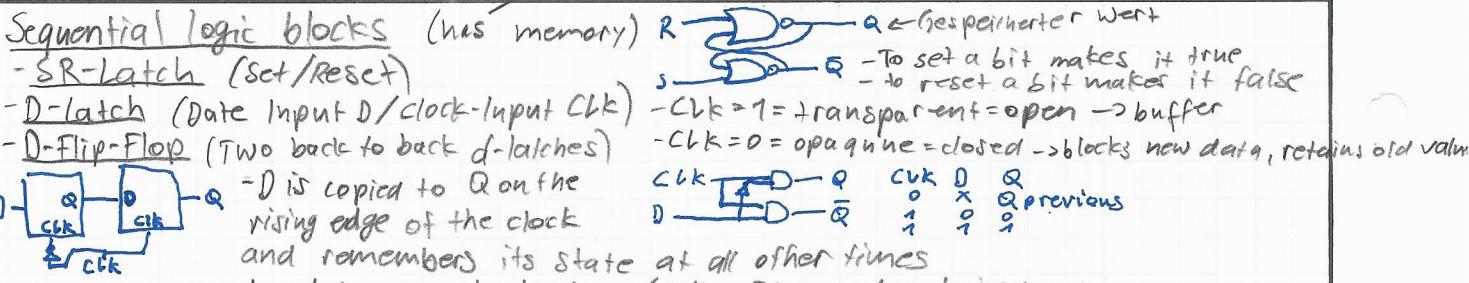
- Propagation delay (t_{pd}): Max time from when an input changes, until the output changes

- Contamination delay (t_{cd}): minimum time from when an input changes until any output starts to change its value (value muss nicht final sein)

Critical path (longest, and therefore slowest path)

2-Level-Logic: Only and/or

→ And cyclic paths



- Latch: Not edge-triggered, bistable / Flip-Flop: edge-triggered
- Enabled flip-flop: Has another input called EN (EN=True: Normal flip-flop)
- Resettable flip-flop: Reset=True → Output=0 (EN=false: ignores clock)

Finite-State Machine: - Circuit with k registers has 2^k unique states

- Next-state-logic: holds current state and on each clock-cycle it changes to next-state
- State-register: stores the different states
- Output-logic: gives out the outputs
- Moore: Only output only depends on current state
- Mealy: Output depends on current state and input
- One-hot-encoding: a separate bit for each state: S1: 001 / S2: 010 / S3: 100 (simpler for next)
- Binary-encoding: Not as simple to encode next-state, aber weniger Flip-Flops
- Moore has more states than Mealy



Timing: element has a propagation time around clock edge, during which the aperture time

↳ for input must be stable for the flip-flop (tsetup-rising edge+hold)

↳ sequential ↳ The clock period has to be long enough for all signals to settle

↳ veraltet ↳ For asynchronous inputs from outside device (button) we need synchronizer flip-flop ↳ Dynamic discipline: inputs of synchronous circuit must be stable during tsetup+hold

Frequenz: T_c between rising edge to falling edge / Frequency: $f_c = 1/T_c$

↳ $T_c \geq t_{propagation} + t_{setup} + t_{flip-flop\ clock\ to\ Q\ propagation\ delay}$

Metastability: Wenn ich Knopf bei rising-edge drücke, weiß Flip-Flop nicht den Value

Parallelism: - spatial parallelism ↳ violates dynamic discipline

↳ multiple copies of hardware so that multiple tasks can be done

- temporal parallelism / pipelining: No additional hardware, just one register per stage

- Latency: L / Throughput: (no parallelism) $1/L$ / (spatial) $1/L$ / Pipelining $1/\text{longest stage}$

- Throughput: Amount of work in a given time: $\frac{1}{\text{longest stage}}$

- latency: Dauer eines ganzen Durchlaufs: # stages · longest stage

Verilog: combinational logic

↳ Multiplexer 4:1

- module sillyfunktion (input a,b,c, output y); - module mux4(input [3:0] d0,d1,d2,d3, input [1:0] s, assign y = a & b & c; output [3:0] y);

endmodule

↳ 9 Bit Bus assign y = s[1]? (s[0]? d3:d2):(s[0]? d1:d0);

- module inv(input [3:0] a, output [3:0] y); - wire p; ↳ internal wire

assign y = ~a; ↳ Little Endian ↳ // % (Mult/Div/Mod) / <<, >> (Left, Right Shift)

endmodule ↳ Multiplexer 2:1 ↳ ==, != (Comparison) - N'Bvalue (Number/Base/Value) z.B. 3'b101

- module mux2(input [3:0] d0,d1, input s, output [3:0] y); - mux2 test(.d0(A), .d1(B), .s(c), .output(y));

assign y = s? d1: d0; // s=1 => y=d1 - case: always @(*) case (data)

endmodule ↳ if/else ↳ s=0 => y=d2 0: segments = 1'b0,

In Verilog, always-statement signals keep old value, until 1: segments = 1'b1;

sensitivity list takes place; Always with q=<= d default: segments = 1'b1; endcase

module flop(input clk, input [3:0] d, output reg [3:0] q); - if: always @(*) if (a[3]) y = 2'b01; always @ (posedge clk) - else if (a[2]) y = 2'b00;

↳ q <= d Must be a reg → left always @ (posedge clk, posedge reset)

endmodule

if (reset) state <= S0; else state <= nextstate;

↳ moore=only states

= Blocking: Values are assigned immediately/++: //next state logic

↳ output logic

always @(*)

- Non-Blocking: At the end of block/++ case (state)

assign y = state

↳ mealy=also input

S0: nextstate=S1;

S1: nextstate=S2;

S2: nextstate=S0;

endcase

↳ parameter S0 = 2'b00;

parameter S1 = 2'b01;

parameter S2 = 2'b10;

Arithmetic Circuits

- Half adder: $A/B/S/C_{out}$ / Full adder: $A/B/S/C_{out}/C_{in}$

- Magnitude Comparator: Usually done by computing $A-B$ and looking at the sign

- Arithmetic Logic Unit (ALU) combines a variety of math/logic operations

- SLT = Set less than: $A < B \Rightarrow Y=1/1, S=A-B/2$. zero-extend unit produces N-Bit 0 Array

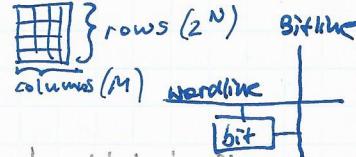
- Shifters: $\cdot 2/S=0101 \cdot 2 \rightarrow 1010$ 3. Sign bit is input

- Multiplication: $0100 \times 0011 = 010000$ Memory is organized in an array

$\begin{array}{r} 0011 \\ \times 0100 \\ \hline 010000 \end{array}$

- N Bit Address

- Each column is called a word



Memory (DRAM/SRAM/ROM) - Ram braucht immer Strom / Rom braucht kein Strom

- Flip-flops: Very fast, parallel accessible, expensive other address in block

- SRAM: Relatively fast, only one data word at a time, less expensive ✓ capacitor leak

- DRAM: Slower, reading destroys content, must be rewritten, cheap, refreshed every $\times M s$.

- Volatile: loses data, when power is off (RAM) / Non-Volatile: keeps data when power off (SSD)

- Cache (Static RAM) / Main Memory (Dynamic RAM) / Virtual Memory (Harddisk/SSD)

Cache miss: compulsory (first load) / capacity / conflict:

- Temporal locality: if you used a book recently, you are likely to use it again

- Spatial locality: if you use a particular book, you are likely to use books in the same shelf

- Cache hit: When CPU requests data, that is already in the cache \rightarrow quickly

- Cache miss: Otherwise CPU gets data from main memory

- Miss-rate: $\frac{\# \text{Misses}}{\# \text{Total}} = 1 - \text{Hit Rate} / \text{Hit-Rate} = \frac{\# \text{Hits}}{\# \text{Total}} = 1 - \text{Miss-Rate}$

- Average memory access time: AMAT = cache + Miss Rate \times (t_{cache} + t_{Main Memory} + Miss Rate \times t_{Virtual Memory})

- Cache exploits temporal + spatial locality to achieve a low miss rate

- Cache block: The fetched word + surrounding words for spatial locality

- Block size: number of words in a cache block

- Capacity C: contains $B = \frac{C}{S}$ blocks each memory location is exactly mapped to one set

- Sets S: A cache is organized into S sets, each of which holds one or more blocks

- Direct mapped cache: each set contains only one block, so cache has $S=B$ sets

- N-way set associative: each set contains $\frac{B}{S}$ blocks, so cache has $S=B/N$

- Fully associative: Only one set

- Each set has a memory address: Tag (which of the many possible addresses is held in that set) +

- Ein Set ist eine Art Vorfilter, dass Set (Set Number) + B-bit Offset

nicht ganzer Cache angeschaut werden muss

- What data is replaced? direct-mapped: the new block for the set / fully associative: LRU

- There are virtual addresses, that are either located in cache or in storage

↳ CPU translates them into physical addresses

Assembly MIPS-language (has 32 registers, that hold commonly used operands) \$0 = 0

- add a, b, c // $a=b+c$ / addi \$50, \$50, 4 R-Type: op, rs, rt, rd, shamt, funct (Registers)

- sub a, b, c // $a=b-c$ offset / addi \$50, \$50, -4 ↳ add / sub

- lw \$53, 1(\$0) // load memory word 1 into \$53 I-Type: op, source dest, imm (immediate)

- sw \$57, 5(\$0) // write \$57 to memory word 5 ↳ addi / lw/sw

- sll \$70, \$1, 4 // $2^4 \cdot 2^{16} = 2^{20}$ ↳ J-Type: op, address (jump)

- srl \$70, \$1, 4 // $2^{16} \cdot 2^4 = 2^{20}$ \$51

- Branch: beq \$50, \$51, target // $$50 == $51 \rightarrow$ target

bne \$50, \$51, target // $$50 \neq $51 \rightarrow$ target

- Jump: j target

- For-loop: addi \$70, \$50, 128

while: if $$70 == 128 \rightarrow$ done

beq \$50, \$70, done

addi \$50, \$50, 1

j while

done:

- Procedures: call: jump and link (=via)

method: jump register (=jr)

arguments: \$90 - \$93

return value: \$v0

- Stacks: Wenn während Method-call keine Daten

zu überschreiten (\$sp: Stack pointer)

stores addi \$sp, \$sp, -4 (Stack geht nach unten)

sw \$R16, \$sp

retrieves lw \$R16, \$sp

addi \$sp, \$sp, 4

Microarchitecture Execution time = (# instructions) · (Cycles) · (seconds) / instructions / cycle

Single Cycle

- Each instruction is executed in a single cycle
- Clock cycle ist abhängig von der langsamsten Operation (Tcritical path)
- 1. Read one instruction / 2. Decode what it means / 3. Find operands from register or memory
- 4. Perform operation / 5. Write back data / 6. Go to next instruction
- Program counter (PC) tells us, where to find the next instruction
 - ↳ Needs to be incremented +4 during each cycle
- Register file: read write from register
- Arithmetic logic unit (ALU): performs operation
- Control unit: setzt alle Funktionssignale: MemtoReg / MemWrite / Branch / ALuSRC / RegDst / RegWrite
- CPI = Cycles per Instruction = 1 (single cycle) + AluOP
- Disadvantages: Clock Cycle: slowest instruction (lw) / 3 Adders / Separate instruction / data memory

The von Neumann Model

- Memory is unified between instructions / data
- Instructions stored in a linear memory array → darum immer PC+4

Multicycle Processor

1. Fetch: fetch instruction from memory
2. Decode: the instruction is being encoded
3. Execute: ALU
4. Memory access
5. Write Back results

} 1 phase per cycle
↳ finite state machine

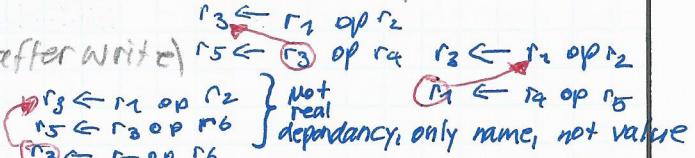
- Execution time = #instructions · CPI · Tc
- What if memory takes longer than one cycle to complete
 - Asynchronous: we wait for ready bit from memory-controller
 - Synchronous: we wait exactly x-cycles, because we know, how long a mem access needs

Pipelining "Assembly line production" same pipeline stages

- Ideal: identical operations / uniform suboperations / independent operations

Data dependences

- Flow dependence (true dependence, read after write) $r_3 \leftarrow r_1 \text{ op } r_2$
- Output dependence (write after write) $r_5 \leftarrow r_3 \text{ op } r_4$
- Anti dependence (write after read) $r_3 \leftarrow r_5 \text{ op } r_6$



How to handle data dependencies?

1. detect and wait until value is available (stalling)
2. detect and forward / bypass data to dependent instructions
 - ↳ Compiler ändert Reihenfolge, indem er unabhängige Instructions zwischen Abhängige gibt
3. Detect and eliminate dependency at the software level
4. fine grained multithreading: switch to another thread every cycle → no dependencies
- Detection of dependencies: Scoreboarding: jedes Register hat eigenes valid bit, das gesetzt wird, wenn ein Thread ins Register schreibt will, spätere Instructions wissen das dann
5. data forwarding: sobald Ergebnis berechnet wurde, wird es direkt an abhängige, spätere Instruktion weitergeleitet, und erst dann im Register gespeichert

Control dependencies / Branch prediction

- Branch is not determined until 4. stage of pipeline

1. Always predict, the branch is not taken
2. Backward branches are usually taken (loops)
3. Look at history, whether branch was previously taken

destination in register
write back
register

Hardware/Software Scheduling

- Software: static scheduling, fixed, known before run-time
 - compiler orders, hardware executes them in order
 - software doesn't know source-registers
 - for example, multiplication with 0, software doesn't know, but hardware does, an copy directly 0
 - Everything, that is determined at run-time is bad
 - Profiling: Compiler testet mit mehreren Inputs
- Hardware: dynamic-scheduling, out-of-order

Exceptions vs. Interrupts

- Exception: caused by internal to running thread
- Interrupt: caused by external to running thread (mouse/keyboard)
- Handle exceptions at the same time
- Handle interrupts, when all instructions are finished in pipeline
- Precise Exception: wenn sie geworfen wird, müssen alle vorherigen Instruktionen abgeschlossen sein, und keine mehr angefangen werden
- When Exception is detected:
 1. All previous instructions are finished
 2. Flush all younger exceptions
 3. Save PC and registers
 4. Redirect fetch address to exception-handling-address
 5. Exception is handled
 6. Resumes previous PC-address

↳ ohne precise exceptions ist Debugging unmöglich

Reorder buffer

- complete instructions out-of-order, but reorder them before making results visible
- Value gets written in reorder-buffer, and he reorders results and writes them back
- Dependencies: Remapping: Wenn ich irgendwo hinschreiben muss, mache ich Pointer von Memory to reorder-buffer

Superscalar-Execution

- Fetch, decode, execute, retire multiple instructions per cycle
- Need to add extra hardware

Combinations

- Scalar > C Superscalar + - Inorder > C out-of-order

Very-long-instruction-word

- consists of multiple, independant instructions packed together in one instruction
- SIMD-Processing (=single instruction operates on multiple data)
- Eine Implementation von SIMD sind Vector-Prozessoren (GPU) (=Arrayprocessor)
- Beispiel: Matrixmultiplikation, same operation performed on many data elements
- Advantages:
 - weniger branches (loops fallen weg)
 - weniger fetch, denn eine Instruktion für viele Daten
 - keine Abhängigkeiten
- Eine Operation muss SIMD-Voraussetzungen erfüllen
- Memory of Vector-GPU is divided into banks, that can be accessed independently

Vector-Chaining: The same as forwarding, simply for SIMD-instructions

GPU

- Each thread executes the same code, but on a different data piece
- Warp: a set of threads that execute the same instruction at the same PC-Count
- Jeder Thread im Warp durchgeht alle 6-Stages
 - ↳ Also "Single-instruction-on-single-data" in jeder Einheit im Warp, aber zusammen dann SIMD



Systolic arrays

- Replace a single processing element with an array of processing-elements
- Difference from pipelining
 - These are individual things to be done on the data
 - Array structure can be non-linear and multi-dimensional
- Advantage: Specialized, improved efficiency, simple design, high concurrency
- Disadvantage: not general \rightarrow specialized
- Taken further: pipeline parallelism and stream processing