

	Answers
	<ul style="list-style-type: none"> - Prozessor <ul style="list-style-type: none"> - Operation unit - Controll Unit - Main Memory - I/O devices - Single-Bus-System
Von-Neumann-Architektur	- keine Aufteilung von Daten und Befehlen
Von-Neimann-Flaschenhals	Es kann zum Rückstau im Bussystem kommen, da alles über 1 Bussystem läuft.
Harvard Architektur	Trennung von Datenspeicher und Befehlsspeicher (2 Bussysteme dafür)
Aufbau eines Microcomputers	<ul style="list-style-type: none"> - Prozessor <ul style="list-style-type: none"> - Operation unit - Controll Unit - Main Memory - I/O devices
Aufgabe eines Prozessors	Kontrolliert den Ablauf von Instruktionen. Besteht aus der Controll-Unit und ALU
Aufgabe der Controll-Unit	Interpretiert die geladenen Befehle und erstellt alle Controllbefehle weiterer Einheiten (übersetzt Befehl in microinstructions die der Rechner ausführen kann)
Aufgabe der Arithmetic-Logic-Unit (ALU)	Erladigt alle Befehle, die ihr erteilt werden. Zentrum aller Berechnungen und Logischen Aufgaben
Aufgabe der Mainmemory	Speichert Befehle und Daten in Speicherzellen mit fester größe, die über Adressen angesprochen werden können
Aufgabe des IO	Akzeptiert Inputs/Gibt Outputs (Kommunikation des Computers mit der Außenwelt)
Mooreschen Law	Alle 18 Monate verdoppelt sich die Rechenleistung der Computersysteme
Pro/Contra Von-Neumann-Architektur	<p>Pro:</p> <ul style="list-style-type: none"> - minimale Hardware bedingungen - minimaler Speicherbedarf <p>Contra:</p> <ul style="list-style-type: none"> - Von-Neumann-Flaschenhals - Keine trennung von Daten und Befehlen - Kein Schutz vom Speicher (alle programme können überall zugreifen/schreiben und können sich so gegenseitig zerstören)
Definitionen von Computer-Architektur	<ul style="list-style-type: none"> - Amdahl, Blaauw and Brooks 1967 - Wie der Prozessor aufgebaut ist, was sein Instruction-Set ist und welche Register und Adressen angesprochen werden können - Processor architecture - Instruction-Set und Prozessoraufbau gesehen als Blackbox => wie sich ein Prozessor verhält und nicht wie der Physikalische Aufbau ist
Klassification von Computersystemen	<p>Klassifiziert werden Systeme über ihre Parallelität</p> <p>5 Level:</p> <ul style="list-style-type: none"> - Programm Ebene - Prozess Ebene <ul style="list-style-type: none"> - Tasks - Block Ebene - Threads - Instruktionen Ebene <ul style="list-style-type: none"> - Basis Instruktionen - Sub-Operationen Ebene - Aufgeteilt in Micro-operationen
Darstellungsformen von Negativen-Zahlen	<ul style="list-style-type: none"> - Betrag + Vorzeichen - Einerkomplement - Zweierkomplement - Offsetdual-Darstellung
Betrag + Vorzeichen	<ul style="list-style-type: none"> - Carrybit gibt an, ob die Zahl Positiv oder Negativ ist <ul style="list-style-type: none"> - 0 => positiv - 1 => negativ - $2^{(n-1)}$ positive als auch negative Zahlen - 2 Darstellungen für die 0!!! 1000 = 0000 (-0=+0) - $[-2^{(n-1)}-1 : 2^{(n-2)}-1]$ - Kippen der Bits - Carrybit gibt wieder an, ob die Zahl positiv oder negativ ist
Einerkomplement	<ul style="list-style-type: none"> - $2^{(n-1)}$ positive als auch negative Zahlen - 2 Darstellungen für die 0!!! - $[-2^{(n-1)}-1 : 2^{(n-1)}-1]$
Zweierkomplement	<ul style="list-style-type: none"> - Kippen der Bits + 1 - Macht man, um bei der 0 einen überlauf zu erzeugen (-0 im EK = 1111 => 0 im ZK => 0000) - Carrybit codiert die Zahl und das Vorzeichen - $2^{(n-1)}$ negative Zahlen und $2^{(n-1)}-1$ positive Zahlen, da die 0 im negativen bereich nicht existiert - $[-2^{(n-1)} : 2^{(n-1)}-1]$
Offsetdual-Darstellung	<ul style="list-style-type: none"> - Verschiebe den Zahlenraum um n positionen (schiebe die Darstellung der Negativen Zahlen in den Positiven Bereich) => genau gleichviele Positive wie negative Zahlen => Offset darf frei gewählt werden, doch nach konvention ist er für Zahlen $2^{(n-1)}$ Addiere auf alle Zahlen diesen Offset - $2^{(n-1)}$ negative Zaheln und $2^{(n-1)}-1$ positive Zahlen, da die 0 im positiven Bereich liegt - $[-2^{(n-1)} : 2^{(n-1)}-1]$
Reellezahlen	<ul style="list-style-type: none"> - Fix-Point-Numbers <ul style="list-style-type: none"> - festes Komma und damit eine feste genauigkeit nach dem Komma - Floating-Point-Number <ul style="list-style-type: none"> - Komma nicht fest und kann damit fließen - Zahlen sind ungleich auf dem Zahlenraum verteilt: <ul style="list-style-type: none"> - Höhere genauigkeit nahe 0 doch großer Zahlenraum (je größer die Zaheln, desto geringer ist die Genauigkeit nach dem Komma)
Floating-Point-Number	<ul style="list-style-type: none"> - Besteht aus Vorzeichen, Charakteristik (Exponent in Offsetdual darstellung) und Mantisse (im ZK) - Stelle des Kommas der Mantisse ist festgelegt (meistens vor dem most-significant-bit) - Normalisiert heißt ein Float, wenn die Mantisse: $1/b \leq \text{mantisse} < 1$ ist mit $b = \text{basis des floats (im Rechner basis = 2)}$ - Bei der Normalisierung trifft man häufig die annahme, dass die Zahle mit 1.... anfängt und spart sich damit die 1 vor dem Komma - Erhöht die größe der Mantisse um eine Stelle ohne einen größeren Speicher zu benötigen - Länge von Charakteristik gibt den Darstellbaren Zahlenraum an - Länge der Mantisse gibt die genauigkeit des Floats an
Floating-Point-Number Overflow/Underflow	Vom Overflow spricht man, wenn man versucht eine zu große Zahl darzustellen.
Floating-Point-Number max real	Vom Underflow spricht man, wenn man versucht eine Zahl nahe 0 darzustellen, die nicht darstellbar ist
Floating-Point-Number min real	Maximal Positive Zahl, die mit einem Float darstellbar ist (Bsp: Charakteristik alles 1 und Mantisse alles 1)
Floating-Point-Number small real	Minimal Positive Zahl, die mit einem Float darstellbar ist ($\neq 0$)
IEEE Rundung	<ul style="list-style-type: none"> - Minimaler Float, den man auf 1.0 Addieren kann und die Zahl verädert $1 + (\text{small real}/2) + (\text{small real}/2) = 1$ und nicht $1 + ((\text{small real}/2) + (\text{small real}/2)) = 1 + \text{small real}$ Das ist der Fall, weil wir mit begrenzent resourcen arbeiten - Exates rechnen mit folgender Rundung muss das selbe Ergebnis liefern, als würden wir erst exact rechnen und dann runden - Einführung der g(uard) und r(ound) stelle + s(ticky bit) - sticky bit == 1, wenn durch das shiften der bits der Mantisse irgendwann eine 1 rausgeschiftet wurde - $r = 0 \rightarrow$ g abrunden (nicht modifizieren) - $r = 1 \rightarrow$ anschauen des sticky bits - sticky bit = 1 \rightarrow aufrunden - " = 0 \rightarrow abrunden

Addition und Subtaktion	<ul style="list-style-type: none"> - Eine Addition kann über einen halb/volladdierer passieren - Eine Subtraktion ist das gleiche wie die Addition der Negativen Zahl
	<ul style="list-style-type: none"> - Nimmt a und b - $a \text{ xor } b = s$ - a und b = ü <p>a und b sind bits s ist das summen bit ü ist das Übertragsbit</p> <ul style="list-style-type: none"> - Reicht nicht für mehrstellige Addition
Halbaddierer	siehe 2.76
	<p>Reihe von 2 halbaddierer</p> <p>nimmt a, b und ü $a \text{ xor } b \text{ xor } ü = s$ (a xor b) und ü oder a und b</p>
Volladdierer	siehe 2.78
Carry-Ripple-Addierer	<ul style="list-style-type: none"> - Aneinandereiung von Volladdierer - Übertrag muss durchgereicht werden - sequenzielle Rechnung - dauer der übertragsleitung gibt max. taktung an
Carry-Lookahead-Addierer	<ul style="list-style-type: none"> - Löst Rechenzeit durch rekursive Berechnung des Übertrages - Übertrag wird mehrfach Berechnet => kein weiterleiten des Übertrages => konstante Laufzeit
Überlaufserkennung bei der Addition	<ul style="list-style-type: none"> - Der Überlauf vom vorletzten bit und des letzten bits muss gleich sein - ist das nicht der Fall gab es einen Überlauf, der nicht gewollt ist
	<ul style="list-style-type: none"> - Registersatz - 2 Eingabe Ports - Multiplexer-Schaltnetz <ul style="list-style-type: none"> - entscheidet, welche und wie die Inputs gelesen werden - ALU-Schaltnetz <ul style="list-style-type: none"> - 2 Eingaben und 1 Ausgabe - Zuständig für alle Arithmetisch und Logische Funktionen - Schiebeneetz <ul style="list-style-type: none"> - Zuständig für das Verschieben der Bits <p>- Alle Teile werden über eine Folge von Bits angesprochen</p>
Komponenten einer einfachen ALU	siehe 2.101
Aufbau eines Micropozessors	siehe 3.4
	<ul style="list-style-type: none"> - Taktgenerator <ul style="list-style-type: none"> - gibt die Geschwindigkeit der CPU vor - Befehlsregister <ul style="list-style-type: none"> - speichert gefeatchte Befehle - Decoder <ul style="list-style-type: none"> - wandelt die Befehle in microprogramme um (falls nötoig) - Steuerregister <ul style="list-style-type: none"> - Übermittelt Aufgaben an die ALU
Aufbau des Steuerwerks	
Phasen der Befehlsausführung	<ul style="list-style-type: none"> - Instruction Featch - Instruction Decode - Execute
Register im Steuerwerk	<ul style="list-style-type: none"> - Befehlsregister <ul style="list-style-type: none"> - multiple Register damit mehrere Befehle gleichzeitig vom Memory gefeatched werden können - Steuerregister <ul style="list-style-type: none"> - enthält Bitfolge zum ansprechen der ALU/des Rechenwerks
Register im Rechenwerk	<ul style="list-style-type: none"> - Statusregister <ul style="list-style-type: none"> - z.B. für Flags - Hilfsregister (für zwischenergebnisse?)
Verbindung des Steuerwerks und Rechenwerk	- Microbefehle aus dem Stuerwerk spricht ALU an
Prozessorstatusword (PSW)	<ul style="list-style-type: none"> - Besteht aus status- und steuerregister (im Bezug des Contextwechsel auch die regulären Register) - Bei Prozesswechsel wird PSW und PC gespeichert
Stackframe	<ul style="list-style-type: none"> - Begrenzt Bereich im Speicher, der Dynamisch vom Programm angelegt werden kann - Benötigt das Speichern der Basis des Stackframs und des Index, wo sich der Pointer befindet - Ermöglicht effizienteren Zugriff auf den Speicher und effizientere Speicherfreigabe
Automatisch modification von Registern	<ul style="list-style-type: none"> - Es existieren Register, die automatisch operationen auf sich selbst ausführen können - e.g <ul style="list-style-type: none"> - Increment - Decrement - Scaling
System-Bus-Schnittstelle	<ul style="list-style-type: none"> - Verbindet CPU mit Hauptspeicher - Rechnet Spannungen um - Datenbus gröÙe mindestens == Adressbus gröÙe
Leistungssteigernde Möglichkeiten	<ul style="list-style-type: none"> - Technologische Maßnahmen; <ul style="list-style-type: none"> - Transistoren anstatt vakuuum röhren (historisch) - Strukturelle Maßnahmen; <ul style="list-style-type: none"> - Anzahl der Transistoren erhöhen - Paralelisierung (ganz wichtig!) - Multicore
Rechnerarchitektur nach Flynn	<ul style="list-style-type: none"> - SISD (Single Instruction Single Data) <ul style="list-style-type: none"> - klassischer von-Neumann-Rechner - SIMD (Single Instruction Multiple Data) <ul style="list-style-type: none"> - bsp. GPU - MIMD (Multiple Instruction Multiple Data) <ul style="list-style-type: none"> - Klassische mehrprozessoren Systeme - MISD (Multiple Instruction Single Data) <ul style="list-style-type: none"> - Pipeline als weit entferntes Modell
Pipelining	<ul style="list-style-type: none"> - Paralelisierung von Aufgaben (wie am Fließband) - Während 1. in Instruction Decode ist kann 2. schon in Instruction Featch - befehle müssen unabhängig von einander sein - Anzahl der Teilschritte == Pipeline Stages (k) - IF - ID - OF - EX - WB = 5 Pipeline Stages (k) - Taktung der Pipeline ist durch den Langsamsten Teilschritt vorgegeben - Throughput = number of instruction leaving pipeline per clock cycle
Effizienz von Pipelins	<p>Ohne Pipeline:</p> <ul style="list-style-type: none"> - $k * n = \text{cycles until all instructions executed}$ ($n \Rightarrow \text{instructions}$) <p>Mit Pipeline:</p> <ul style="list-style-type: none"> - $k + (n-1) = \text{cycles until all instructions executed}$ ($n \Rightarrow \text{instructions}$)
Pipeline Speedup s in perfect world	$k / ((k/n) + 1 + (1/n))$; k = pipeline stages and n = number of instructions
Problems of Pipelines in real world	- Abhängigkeiten von Befehlen

Pipeline Hazards	<ul style="list-style-type: none"> - Data Hazard <ul style="list-style-type: none"> - value needed, that is currently in calculation - Structural hazard <ul style="list-style-type: none"> - needed resources, that are not available - Control hazard <ul style="list-style-type: none"> - Jumps and Conditional jumps = non-linear program
Data Hazards	<ul style="list-style-type: none"> - RAW (Read after write) <ul style="list-style-type: none"> - Instruction 2 wants to read operand before instruction 1 wrote it - WAR (Write after read) <ul style="list-style-type: none"> - Instruction 2 writes operand before instruction 1 read it - only in superscalar pipelines - WAW (Write after Write) <ul style="list-style-type: none"> - Instruction 2 wants to write before instruction 1 wrote it - only in superscalar pipelines
Solutions to Data hazards	<p>Software solutions:</p> <ul style="list-style-type: none"> - No-Operations (NOP's) - Reordering of commands <p>Hardware solutions:</p> <ul style="list-style-type: none"> - Pipeline Stall / Bubbles - Forwarding of informations <ul style="list-style-type: none"> - Result forwarding: Exe -> Exe - Forwarding vom MEM -> Exe <p>(In MIPS-Architecture)</p>
Structural Hazards	<ul style="list-style-type: none"> - Resource Problems - Stalling or resource upgrade is the only solution
Control Hazard	<ul style="list-style-type: none"> - Stalling of Pipeline till conditional jump is evaluated - Branch Prediction to have a "linear"-programm - Worst-Case: Pipelineflush!
Branch Prediction	<ul style="list-style-type: none"> - Deciding, if a conditional jump should be taken or not (as a prediction) - One Solution is the Branch-Delay-Slot: Loading a command into the Pipeline that can be executed in both scenarios - Static Predictions: <ul style="list-style-type: none"> - always take/not take jump - compile decides if jump is always taken/not taken - Dynamic Predictions: <ul style="list-style-type: none"> - evaluate the past predictions and therefore adapt behaviour - one-bit predictor and two-bit predictor
Branch-Target-Buffer	<ul style="list-style-type: none"> - Table of known jumps + it's locations to jump to - In complex versions it also contains the state of the Branch-Prediction-Automat
One-Bit-Predictor	<p>2 states:</p> <ul style="list-style-type: none"> - Predicate Take - Predicate Not Take <p>Miss judgement will change the state for next predictions</p>
Two-Bit-Predictor	<p>4 states:</p> <ul style="list-style-type: none"> - Strongly Take - Weakly Take - Strongly Not Take - Weakly Not Take
Hysteresis Two-Bit-Predictor	See: 3.130
Saturation-Counter Two-Bit-Predictor	See: 3.131
Hysteresis vs Saturation-Counter	Hysteresis has the advantage not to flutter between Weakly Not Taken and Weakly Taken because there are no direct connections between the two
Predicated Instructions	Predicated Instructions are always executed, but the result has only an impact, if the condition is true
Pro/Contra Predicated Instructions	<p>Pro:</p> <ul style="list-style-type: none"> - eliminated the need for a branch prediction - keeps the pipeline busy <p>Contra:</p> <ul style="list-style-type: none"> - Unnecessary code might be executed - Might be executing too much code that is unnecessary, if the execution of the predicated instructions takes longer than a flush would take
Vektor-Pipelining	<p>Ausführung des selben commands auf einen Vektor von Daten. Bsp. Array von n Zahlen in einem großen register und auf jede Zahl 3 Addieren schneller als n Register mit je 1 Addition</p> <p>SIMD Prinzip</p>
Chaining	Eine Instruction führt mehrere Vektor-Instructionen hintereinander aus
Superscalare Prozesse	<ul style="list-style-type: none"> - Eine Pipeline hat mehrere Einheiten von verschiedenen Operationen (Multiple Addition, Multiplikations, etc. Einheiten) - Langsamster Takt gibt nicht die gesamte Taktung vor - echte Gleichzeitigkeit! <ul style="list-style-type: none"> - mehrere Befehle können pro Takt ausgeführt werden - CPU (Cycles per Instruction) << 1.0 möglich! - Befehle können einander überholen (da verschiedene Operationen verschiedene Zeiten haben) - in-order -> out-of-order -> in-order <ul style="list-style-type: none"> - IF ID -> OF EXE -> Retire/WB - Siehe 3.150 für Darstellung
Issue-Unit	Verteilt in Superscalaren Prozessen die Befehle auf die verschiedenen Recheneinheiten
Reservation- Station	Entweder in-/out-of-order (je nach implementation)
Completion of Superscalar Processes	Vor der Executestation; hier warten die instructionen auf ihre Daten
Instruction-Set-Architectur (ISA)	<p>Happens in-order</p> <p>Has to wait for operations to finish in order before it can write them back</p> <p>RISC (Reduced-Instruction-Set)</p> <ul style="list-style-type: none"> - directly executable without microprograms <p>CISC (Complex-Instruction-Set)</p> <ul style="list-style-type: none"> - execution of microprograms
CISC (Complex-Instruction-Set)	<ul style="list-style-type: none"> - ISA soll mächtige Befehle enthalten, sodass selten auf den Speicher zugegriffen werden muss - 1 CISC-Befehle => mehrere Befehle im Microprogramm - Contra: <ul style="list-style-type: none"> - Benötigt sehr großes Schaltnetz - nur wenig Befehle werden oft genutzt (lohnt sich also der Aufwand?)
RISC (Reduced-Instruction-Set)	<ul style="list-style-type: none"> - Reduzierte Befehlsanzahl <ul style="list-style-type: none"> - <= 128 - Reduzierte Formate <ul style="list-style-type: none"> - <= 4 - Reduzierte adressierungs modelle <ul style="list-style-type: none"> - <= 4 - Pro: <ul style="list-style-type: none"> - Creates space on chip - memory access only allowed for load and store (so scheduling of instructions is easier because these two are the "heaviest" for the cpu) - Pipelining developed favoured because of this - Contra: <ul style="list-style-type: none"> - von-Neumann-Bottleneck
Procedures, Traps and Interrupts	<ul style="list-style-type: none"> - Programme werden unterbrochen durch: <ul style="list-style-type: none"> - Procedures (functions Aufrufe von z.B. OS funktionen), - Multithreading (OS wechselt auf einen anderen prozess) und - Hardware bzw. Software Interrupts

	<ul style="list-style-type: none"> - Programme können unterbrochen werden, wenn etwas unmittelbar unterbrochen werden muss - Der Prozess des Interrupts wird Interrupt Service Routine (ISR) genannt <ul style="list-style-type: none"> - sollte relativ schnell sein - Interrupt-Vektor-Table (IVT) <ul style="list-style-type: none"> - Mapping von der Service routine zu dem jeweiligen Interrupt - Die Adresse des Vektor-Tables ist fixed, dass heißt das Betriebssystem muss diesen Table füllen mit ISA - Bevor die Interrupt-Service-Routine ausgeführt werden kann muss der PC und Context des Programms gesichert werden <ul style="list-style-type: none"> - Während der Contextsicherung darf der Interrupt nicht Unterbrochen werden! Sonst kommt man in einen nicht definierten Programmstatus - Interrupts können einander unterbrechen, wenn dies gewünscht ist. Regel ist: Interrupt eines Interrupts nur wenn die Priorität des Interrupts höher ist
Interrupts	<ul style="list-style-type: none"> - Software Interrupts => Process internal (synchron ausführung) - Hardware Interrupts => Process external (asynchron ausführung)
Ablauf einer Interrupt Service Routine (ISR)	<ul style="list-style-type: none"> - Aktivierung des Interrupts - Beenden des gerade in Ausführung befindlichen Befehls - Feststellen ob es sich um ein SW oder HW interrupt handelt - Feststellen, ob Interrupts erlaubt sind (Interrupt Enable bit == true) - Falls HW-Interrupt: Quelle feststellen und Interrupt Acknowledge schicken - Interrupt Enable Bit = 0 - Save PSW und PC - Startadresse des ISR ermitteln und starten - ISR ausführen <ul style="list-style-type: none"> - Context sicherung - Interrupt Enable Bit = 1 - Routine ausführen - IRET-Befehl ausführen (Ende der ISR) - Context, PSW und PC wieder herstellen
Interrupt-Controller	<ul style="list-style-type: none"> - Managed einkommende Interrupt anfragen
Methoden von Interrupt Handling	<ul style="list-style-type: none"> - Polling: <ul style="list-style-type: none"> - CPU checkt periodisch Interrupt-Controller, ob ein Gerät einen Interrupt angefordert hat - Polling 1: <ul style="list-style-type: none"> - cyclische abfrage der Geräte, ob sie ein Interrupt haben. Wenn ja und der bearbeitet wurde, dann chekt die CPU das nächste - Polling 2: <ul style="list-style-type: none"> - cyclische abfrage der Geräte, startet nach bearbeitung eines Interrupts wieder bei 1 - Interrupt: <ul style="list-style-type: none"> - Interrupt-Controller meldet sich bei der CPU sobald ein Gerät ein Interrupt meldet
Speicherhierarchie	<ul style="list-style-type: none"> - Durch den effekt: klein und schnell oder groß und langsam wird der Speicher in Systemen hierarchisch angeordnet - Kommunikation zwischen den verschiedenen Speicherschichten ist nur unmittelbar erlaubt (see 4.8) <ul style="list-style-type: none"> - um Zugriffszeiten zu minimieren (es wird immer im nächst langsameren Medium nach dem Datum gesucht) - Speicher soll so wirken als sei er beides: groß und schnell - Speicherverschiebung wird genutzt um zu gewährleisten, dass benötigte Daten so schnell wie möglich gelesen/bearbeitet werden können
Lokalitätsprinzipien des Speichers	<p>Zeitliche lokalität:</p> <ul style="list-style-type: none"> - was benötigt wurde wird schnell wieder benötigt <p>Örtliche lokalität:</p> <ul style="list-style-type: none"> - was nebeneinander steht wird oft gemeinsam/nacheinander benötigt - bsp. Arrays
Zugriffszeit und Zykluszeit	<p>Zugriffszeit:</p> <ul style="list-style-type: none"> - Definiert als Zeit die benötigt wird um ein Datum aus dem Speichermedium zu lesen <p>Zykluszeit:</p> <ul style="list-style-type: none"> - Definiert als Zeit die benötigt wird um das Datum auszulesen und sich der Speicher regeneriert hat um den nächsten Zugriff zu erlauben
Cache-Speicher	<ul style="list-style-type: none"> - Wird genutzt um die verhältnismäßig langsame Geschwindigkeit des RAMs zu kompensieren - Wichtige Informationen aus dem DRAM werden als Blöcke in den Cache gelagert, in der Hoffnungs, dass sie bald wieder benötigt werden - Darstellung des Cache-Speichers 4.26
Was Passiert bei einer Abfrage eines Datums	<ul style="list-style-type: none"> - Adresse wird über den Adressbus an den Cache und Speicher weitergeleitet <ul style="list-style-type: none"> - Steuersignal wird zunächst nur an den Cache geschickt - ist das Datum im Cache (Cache-Hit) wird das Datum auf den Datenbus gelegt - ist das Datum im Cache (Cache-Miss) so wird das Steuersignal an den Hauptspeicher weitergeleitet, der dann nach dem Datum sucht <ul style="list-style-type: none"> - Datum wird auf den Datenbus gelegt - Cache speichert Datenblock ab mit dem benötigtem Datum
Cache Performance	<ul style="list-style-type: none"> - Cachegeschwindigkeit ist am höchsten, wenn die beiden Lokalitätsprinzipien eingehalten werden - Hitrate = Anzahl der Treffer Pro Zugriff - $t_{access} = (Hitrate) * t_{hit} + (1-Hitrate) * t_{miss}$ <ul style="list-style-type: none"> - t_{hit} = time needed for hit - t_{miss} = time needed for miss
Cache-Schreibzugriffe	<ul style="list-style-type: none"> - Bei Cache-Miss: <ul style="list-style-type: none"> - Änderung des Datums im Speicher und Schreiben des neuen Datums in den Cache - Bei Cache-Hit <ul style="list-style-type: none"> - Veränderung des Datums im Cache und je nach implementation direct Weiterleitung des Datums an den Hauptspeicher
Cache write through	<p>Änderung des Datums direct im Cache und weiterleitung an den Arbeitsspeicher</p> <p>Pro:</p> <ul style="list-style-type: none"> - garantierte konsistenz zwischen Cache und Speicher <p>Contra:</p> <ul style="list-style-type: none"> - Schreibzugriffe dauern lange
Cache buffered write through	<p>Änderung des Datums im Cache und weiterleitung über einen Buffer an den Arbeitsspeicher um die Schreibgeschwindigkeit des Speichers zu kompensieren</p> <p>Pro:</p> <ul style="list-style-type: none"> - Prozessor muss nicht auf beendigung des Schreibbefehls warten <p>Contra:</p> <ul style="list-style-type: none"> - Temporäre Inkonsistenz - mächtiger Controller wird benötigt um doppelte Schreibzugriffe abzufangen (neuer Schreibbefehl obwohl vorheriger noch läuft)
Cache write back	<p>Änderung des Datums im Cache, doch keine weiterleitung des Datums an den Speicher. Erst bei Auslagerung des Datums wird geschaut, ob das Dirty-bit gesetzt ist (ob das Datum verändert wurde) und wenn ja dann wird es im Speicher geändert</p> <p>Pro:</p> <ul style="list-style-type: none"> - Schreibbefehle können in schneller Cache geschwindigkeit geschehen <p>Contra:</p> <ul style="list-style-type: none"> - Inkonsistenz zwischen Cache und Speicher (großes problem für multi-core prozessoren)
Aufbau eines Cache-Speichers	<p>See 4.48</p> <ul style="list-style-type: none"> - Aufgeteilt in Adress und Datenspeicher - Komparator checkt, ob Adresse im Cache ist - Cache-Lines werden Datenspeicher abgespeichert <ul style="list-style-type: none"> - mehrere bytes auf einmal einladen - Örtliche lokalität
Begriffe zum Cache	<p>Block:</p> <ul style="list-style-type: none"> - Datenportion die zu einem Block zusammen gefasst ist <p>Blockrahmen:</p> <ul style="list-style-type: none"> - Eine Reihe im Cache mit der Basis Adresse des jeweiligen Blocks, Statusbit (ob Cacheline valide ist) und Daten des Blocks (und vielleicht Dirty-bit) <p>Blocklänge:</p> <ul style="list-style-type: none"> - Anzahl der Speicherplätze im Blockrahmen <p>Assoziativität:</p> <ul style="list-style-type: none"> - Anzahl der Blockrahmen, die zu einem Satz zusammengefasst wurden
Cachespeicher Assoziativitätsvarianten	<p>Voll-Assoziativ:</p> <ul style="list-style-type: none"> - Cache besteht nur aus einem Satz (alle Blöcke konkurrieren um alle Plätze, dafür kann jeder Block überall im Cache liegen) <p>Direct-Abgebildet:</p> <ul style="list-style-type: none"> - Jeder Satz enthält nur einen Blockrahmen - (Ein Block kann immer nur an einer Stelle im Cache stehen und konkurriert mit allen anderen die diesen Platz brauchen) <p>n-fach-Assoziativ:</p> <ul style="list-style-type: none"> - Sätze bestehen aus n Blöcke die nur miteinander konkurrieren, wenn sie auf den selben Satz mappen

Verdrängungsstrategien	<ul style="list-style-type: none"> - Verdrängungsstrategien werden genutzt um Platz für Blöcke zu schaffen der bereit belegt ist - Nur verwendet bei Voll-Assoziativ und n-fach-Assziativ <ul style="list-style-type: none"> - Bei Direct-Abgebildet ist das nicht notwendig, dass der Block nur an einer stelle stehen kann
Voll-assoziativer Cache	<ul style="list-style-type: none"> - Alle Blöcke können überall im Cache stehen - Echtparallel aller Adressen im Cache <p>Pro:</p> <ul style="list-style-type: none"> - Ein Block kann überall im Cache stehen - Optimale Voraussetzungen um am effektivsten zu verdrängen - Optimale Ausnutzung des Caches <p>Contra:</p> <ul style="list-style-type: none"> - Sehr hoher Hardwareaufwand <ul style="list-style-type: none"> - nur für sehr kleine Caches verwirklichtbar - Verdrängungsstrategie braucht komplexe Logik in der Hardware
Direct-mapped-Cache	<ul style="list-style-type: none"> - Ein Datum Mapped auf nur einen Platz im Cache <p>Pro:</p> <ul style="list-style-type: none"> - Geringer Hardwareaufwand - Keine Logik benötigt für die Verdrängung <p>Contra:</p> <ul style="list-style-type: none"> - Ständiges Konkurrieren der Blöcke (Flattern kann passieren, da viele Blöcke um genau den selben Platz konkurrieren)
n-Way Set Associative Cache	<ul style="list-style-type: none"> - Sätze enthalten mehre Blöcke, dass heißt weniger Konkurrenz zwischen verschiedenen Blöcken, da sie auf mehr Plätze pro Satz mappen - Bei Suche nach Datum müssen alle n Tags des Satzes überprüft werden - Ist ein kompromiss zwischen Voll-assoziativer- und Direct-mapped-Cache
Virtueller-Speicher	<p>Je größer der Satz, desto näher n-Way an den Pro und Contras des Voll-assoziativen und umso kleiner desto näher an Direct-mapped Pro und Contra</p> <ul style="list-style-type: none"> - Jeder Prozess bekommt seinen eigenen Logischen-Adressraum <ul style="list-style-type: none"> - Die Adressen des logischen-Adressraumes müssen auf die realen-Adressen gemappt werden - Bei Programmstart wird entschieden, wo das Programm im realen-Adressraum abgelegt wird - Betriebssystem kümmert sich um die Verwaltung des Virtuellen-/Logischen-Adressraums <ul style="list-style-type: none"> - Hardwaretechnisch von der MMU (Memory-Management-Unit) unterstützt <ul style="list-style-type: none"> - Enthält eine Übersetzungstabellen, die vom Betriebssystem verwaltet wird - Wird genutzt um nicht benötigte Daten aus dem Hauptspeicher in den Harddrive auszulagern
Aufgaben der MMU (Memory-Management-Unit)	<ul style="list-style-type: none"> - Sinnvoll für Multitasking/Multithreading/Multiuser - Umrechnung der Logischen-Adressen in Reale-Adressen - Schutz vor Falschzugriff auf nicht berechnigte Real-Adressen - Datenschutz
Virtuelle-Speicher-Verwaltung	<ul style="list-style-type: none"> - Paging <ul style="list-style-type: none"> - Daten werden in Seiten mit fest bytelänge unterteilt - "Feste Segmentlänge" - Segmentierung <ul style="list-style-type: none"> - Programme können in verschiedene Segmente unterteilt werden - jedes Segement verhält sich wie ein eigener Speicherraum
Speicherverwaltung Paging	<ul style="list-style-type: none"> - Daten werden in Seiten mit fester länge unterteilt - Seiten sind reativ klein (max. 4kbyte) - Ein Prozess wird auf viele Seiten verteilt - Im Virtuellen Speicherraum sind die Prozesse zusammenhängend; im Realen-Speicherraum können die Pages verteilt sein - Page-Table speichert die Abbildung des Virtuellen-Speicherraums auf den Reallen-Speicherraums <ul style="list-style-type: none"> - Übersetzung der Virtuellen-Adressen in Real-Adressen <p>Problem:</p> <ul style="list-style-type: none"> - Dynamische Speicherstrukturen können nur schlecht dargestellt werden, da Pages eine feste Länge besitzen
Speicherverwaltung Segmentierung	<ul style="list-style-type: none"> - Virtueller-Adressraum wird in Segmente verschiedener Länge unterteilt - Ein Programm besteht aus verschiedent vielen Segmenten - Jedes Segment ist in sich logisch zusammenhängend <p>Pro:</p> <ul style="list-style-type: none"> - Segmentierung spiegelt logische Programmstruktur wieder - Große Segmente => seltener Datentransfer - Segmentspezifische Zugangsrechte <ul style="list-style-type: none"> - Read-Only oder Read-Write Zugriffe <p>Contra:</p> <ul style="list-style-type: none"> - Wenn Segmente transferiert werden müssen, dann sind es große Transfere - Keine einlagerung von Teilsegmenten
Virtueller Cache	<ul style="list-style-type: none"> - Adressen werden im Cache Virtuell abgespeichert - Keine Umrechnung der MMU benötigt um auf die Adressen zu zugreifen <p>Pro:</p> <ul style="list-style-type: none"> - bei Cache-Hit keine Umrechnung der MMU <p>Contra:</p> <ul style="list-style-type: none"> - Dadurch, dass Virtuelle Adressen öfters vergeben werden können kann der Cache nicht zwischen den Virtuellen Adressen verschiedener Prozesse Unterscheiden
Physikalischer Cache	<ul style="list-style-type: none"> - Wahre Adressen werden im Cache abgespeichert - MMU muss Adressen zunächst umrechnen bevor der Cache abgefragt werden kann <p>Pro:</p> <ul style="list-style-type: none"> - Keine Duplizierten Adressen <p>Contra:</p> <ul style="list-style-type: none"> - Umrechnung benötigt von der MMU - Potentiell kleinerer Adressraum
Schutzmechanismenen des Hauptspeichers	<ul style="list-style-type: none"> - Schutz der Realenspeicherräume ist Hardwaretechnisch unterstützt - MMU
Cache-Verwaltung bei Multicoresystemen	<ul style="list-style-type: none"> - Speichergekoppelt (geteilter Speicher) <ul style="list-style-type: none"> - Besitzen einen geteilten Adressraum pro Prozess - Shared Memory (beide Prozessoren greifen auf den Physikalisch seleben Speicher zu) - Distributed Memory (beide Prozessoren besitzen eigenen Physikalischen Speicher (Daten sind auf die verschiedenen Speicher aufgeteilt)) - Nachrichten basiert (Send/Recieve von Infos) <ul style="list-style-type: none"> - Physikalisch zwei Speicher die über Send und Recieve miteinander Kommunizieren
Probleme von Multicoreprozessoren	<ul style="list-style-type: none"> - Höherer Verwaltungsaufwand - Systemverklemmung möglich (A wartet auf B wartet auf C wartet auf A) - Bottleneck wenn Prozessoren miteinander Kommunizieren müssen - Abstimmung des Cache (MESI)
Cache-Kohärenzproblem	<ul style="list-style-type: none"> - Kohärenz -> Caches können unterschiedliche Daten beinhalten (== nicht konsistent) doch die korrekte Funktionsweise ist gewährleistet - Konsistenz -> Cache ist bei allen Prozessoren auf dem selben Stand - Cache-Kohärenzprotokolle sichern den korrekten ablaufen auch bei nicht konsistenten Datenverteilungen
Cache-Kohärenzprotokolle	<p>Write-Update-Protokoll:</p> <ul style="list-style-type: none"> - Kopien aller Caches sind konsistent (koheränz gesichert, da konsistenz existiert) <p>Write-Invalidate-Protokoll:</p> <ul style="list-style-type: none"> - falls veränderung in einem anderen Cache geschieht, wird die alte kopie Invalid gemacht <p>MESI-Protokoll:</p> <ul style="list-style-type: none"> - beruht auf Mittauschen des Bussystems - Besteht aus 4 Zuständen <ul style="list-style-type: none"> - Modified - Exclusive unmodified - Shared unmodified - Invalid