

**Enterprise Computing  
Einführung in das Betriebssystem z/OS**

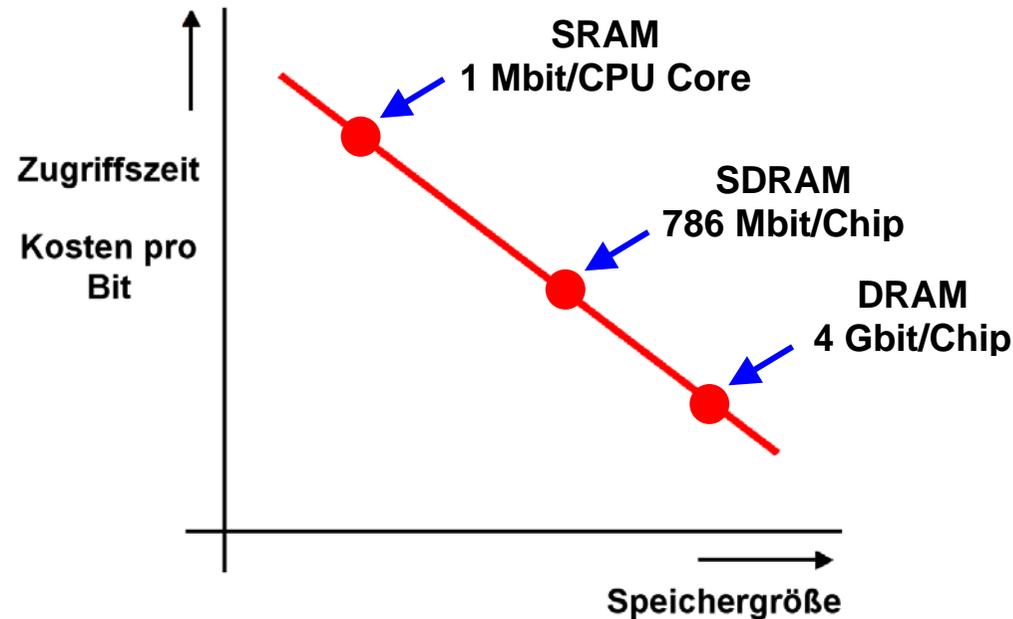
**Prof. Dr. Martin Bogdan  
Prof. Dr.-Ing. Wilhelm G. Spruth**

**WS2012/13**

**Verarbeitungsgrundlagen Teil 4**

**Cache**

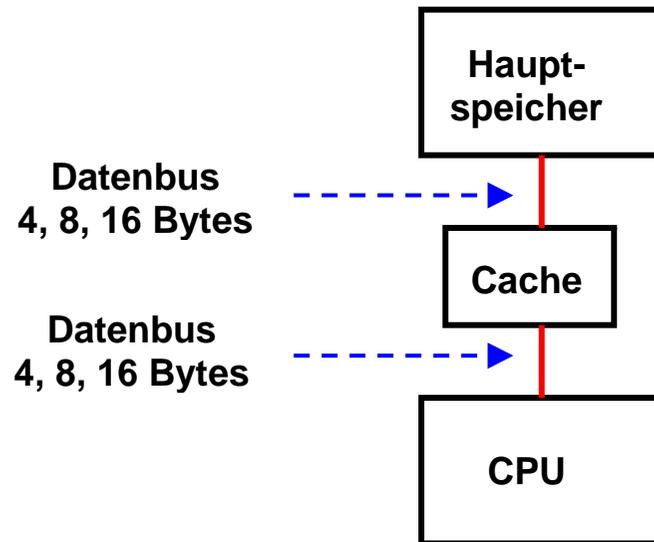
# Halbleiter Speicher Chips



Es existieren eine ganze Reihe unterschiedlicher Technologien, mit denen man einen Speicher auf einem Halbleiterchip realisieren kann. Die Technologien unterscheiden sich in der Zugriffszeit auf den Speicher sowie der Speicherdichte, der Anzahl der Bits, die man pro mm<sup>2</sup> unterbringen kann. Die Speicherdichte beeinflusst die Kosten pro Bit. Allgemein gilt der hier dargestellte Zusammenhang: Kleine Speicher (Schnellspeicher) haben eine schnelle Zugriffszeit, brauchen viel Platz pro Bit und haben hohe Kosten pro Bit. Große Speicher haben eine längere Zugriffszeit, brauchen wenig Platz pro Bit und haben niedrige Kosten pro Bit.

Hauptspeicher Chips verwenden die DRAM (Dynamic Random Access Memory) Technologie. Derzeitig (2011) sind 4 Gbit Chips verfügbar. Das speichernde Element ist dabei ein Kondensator, der entweder geladen oder entladen ist. Über einen Schalttransistor wird er zugänglich und entweder ausgelesen oder mit neuem Inhalt beschrieben. Der Speicherinhalt ist flüchtig (volatil), das heißt, die gespeicherte Information geht bei fehlender Stromversorgung oder zu später Wiederauffrischung verloren.

Plattenspeicher sind dagegen statisch. Die gespeicherten Daten bleibt auch im abgeschalteten Zustand erhalten.

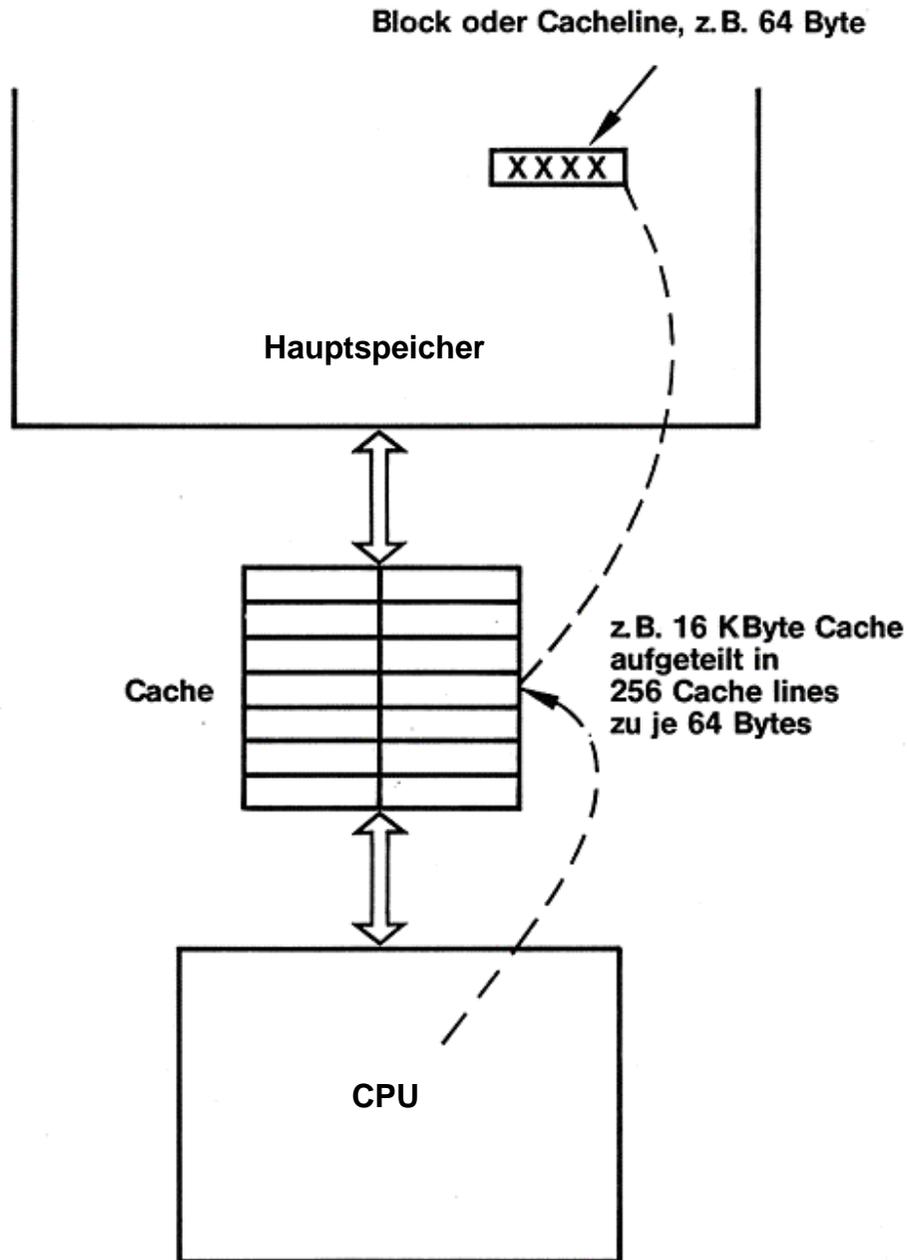


## Cache Speicher

Hauptspeicher Zugriffszeiten sind zu langsam, um mit der heute möglichen Verarbeitungsgeschwindigkeit einer CPU mithalten zu können. Deswegen schaltet man zwischen Hauptspeicher und CPU einen Cache Speicher. Der Cache Speicher verwendet eine SRAM (Static Random Access Memory) Technologie. Das speichernde Element ist dabei ein FlipFlop. Es existieren viele unterschiedliche SRAM Technologien, mit unterschiedlichen Zugriffszeiten, Speicherdichten und Kosten. Während mit DRAMs aufgebaute Hauptspeicher eine Zugriffszeit in der Gegend von 100 ns aufweisen, haben mit SRAMs aufgebaute Cache Speicher Zugriffszeiten zwischen 100 ps und 10 ns .

Wenn man von einem Cache redet meint man meistens einen Hauptspeicher Cache. Da es auch Plattenspeicher-Caches und andere Caches gibt ist die Unterscheidung wichtig.

Jeder moderne Rechner hat einen oder mehrere Caches. Die Existenz des Caches ist für den Benutzer praktisch unsichtbar. Es existieren auch keine Maschinenbefehle, mit denen der Programmierer den Inhalt des Caches manipulieren kann.



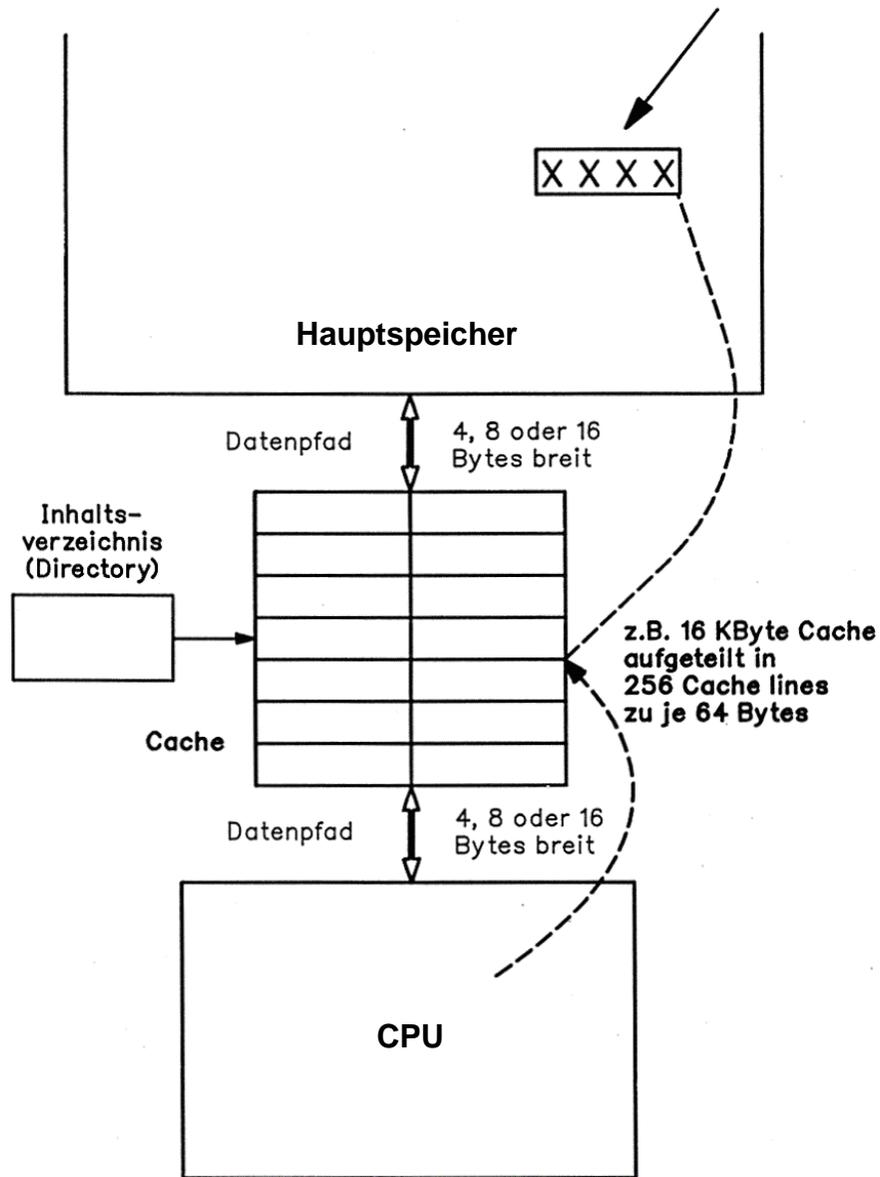
## Struktur des Cache Speichers

Der Cache enthält in jedem Augenblick nur die Kopie einer Untermenge der Daten im Hauptspeicher. Diese Untermenge wird ständig ausgewechselt.

Hierzu werden Hauptspeicher und Cache Speicher in Blöcke mit einer identischen Größe aufgeteilt. Diese Blöcke werden als „Cachelines“ bezeichnet. Die Größe der Cachelines ist implementierungsabhängig. Bei den z10 und z196 Mainframes sind es 256 Bytes, bei anderen Rechnern häufig weniger.

Ein Prozessor mit einer Cacheline-Size von 256 Byte wird aus dem Hauptspeicher immer nur Pakete dieser Größe in den Cache transportieren. Bei einem Hauptspeicher-Interface mit z.B. 256 Datenleitungen bedeutet das, dass jeder Cache Miss (wenn also angeforderte Daten nicht im Cache stehen) eine Adressierung und 8 Daten-Transferzyklen nach sich zieht.

Block oder Cacheline, z.B. 64 Byte



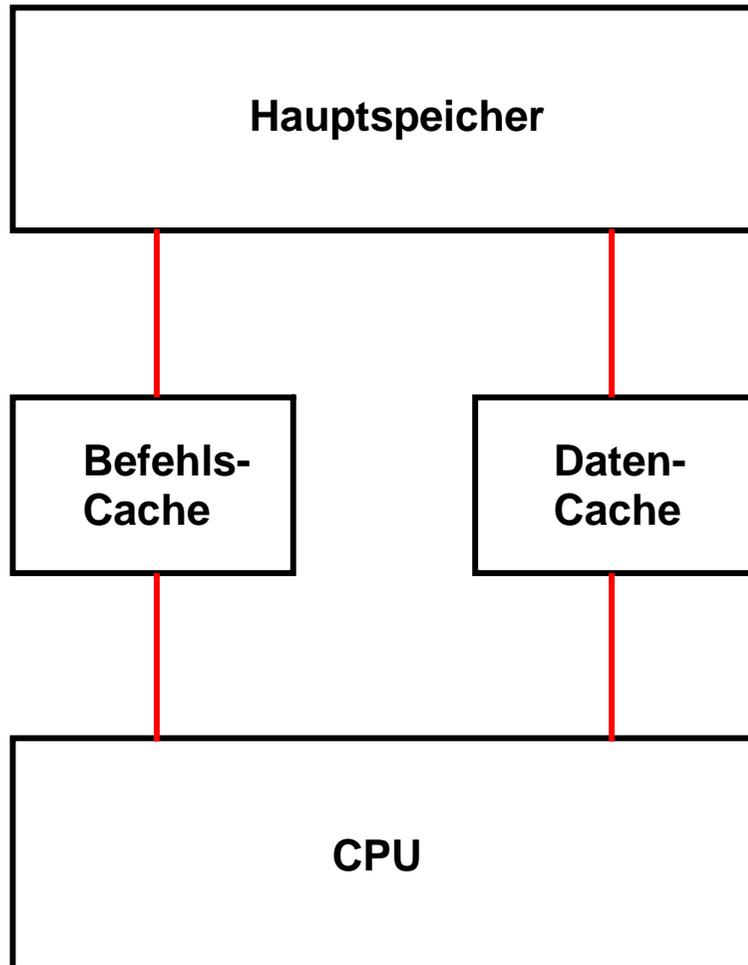
Die Anordnung der Cachelines ist willkürlich und ändert sich ständig. Die CPU konsultiert ein „Cache Directory“ (Inhaltsverzeichnis) um eine bestimmte Cacheline innerhalb des Caches zu finden.

Das Cache Directory enthält je einen Eintrag für jede im Cache gespeicherte Cache Line. Der Eintrag enthält die Adresse der Cacheline im Hauptspeicher und im Cache.

Bei jedem Hauptspeicherzugriff durchsucht die CPU das Cache Directory in der Hoffnung, dass die benötigte Cacheline sich im Cache befindet. Wenn das nicht der Fall ist entsteht ein Cache Miss. Dies bewirkt, dass die benötigte Cache Line in ihrer Gesamtheit vom Hauptspeicher in den Cache geladen wird. Vermutlich muss dabei eine andere (nicht mehr benötigte) Cache Line aus dem Cache ausgelagert werden um Platz zu schaffen.

Einzelheiten hierzu in:

Udo Kebschull, Paul Herrmann, Wilhelm G: Spruth:  
Einführung in z/OS und OS/390.  
Oldenbourg-Verlag, 2002, ISBN 3-486-27214-4



## **Split Cache**

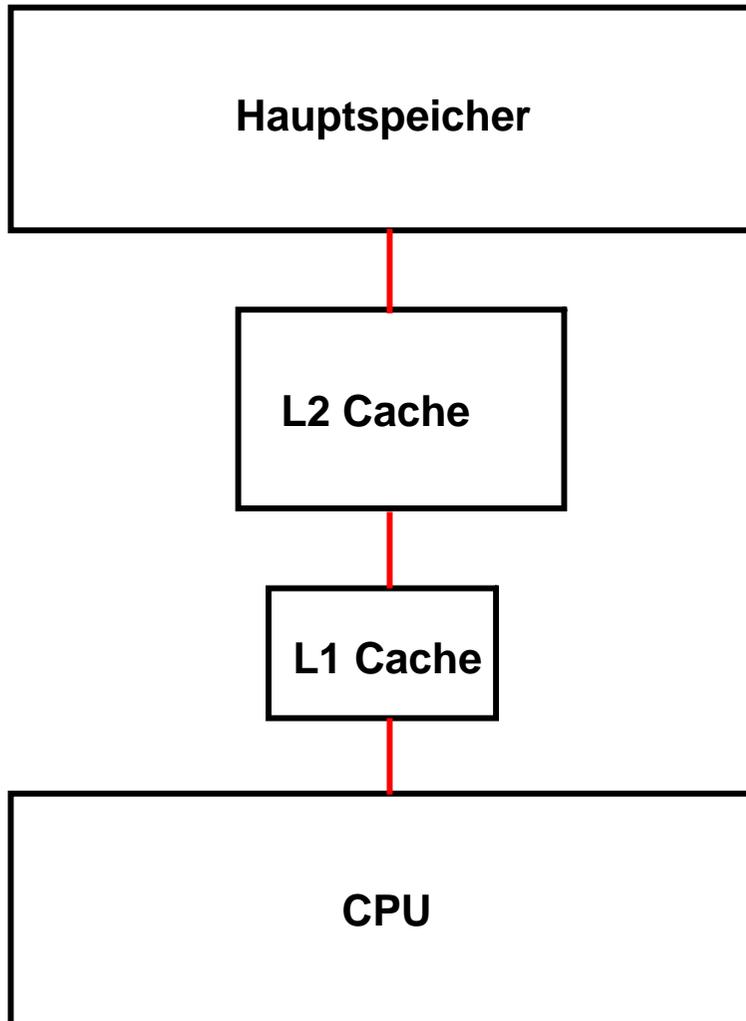
**Dual Cache, Harvard Architecture Cache**

Im Hauptspeicher sind Programme (bestehend aus Maschinenbefehlen) und Daten gespeichert. Diese befinden sich in unterschiedlichen Hauptspeicherbereichen und damit auch in unterschiedlichen Cachelines.

Um den Durchsatz zu verbessern besteht der Cache häufig aus zwei unabhängigen Cache Speichern (Dual Cache): Der Befehls-cache enthält nur Maschinenbefehle und der Datencache enthält nur Daten.

Der z9, der z10 und der z196 Cache ist ein Dual Cache.

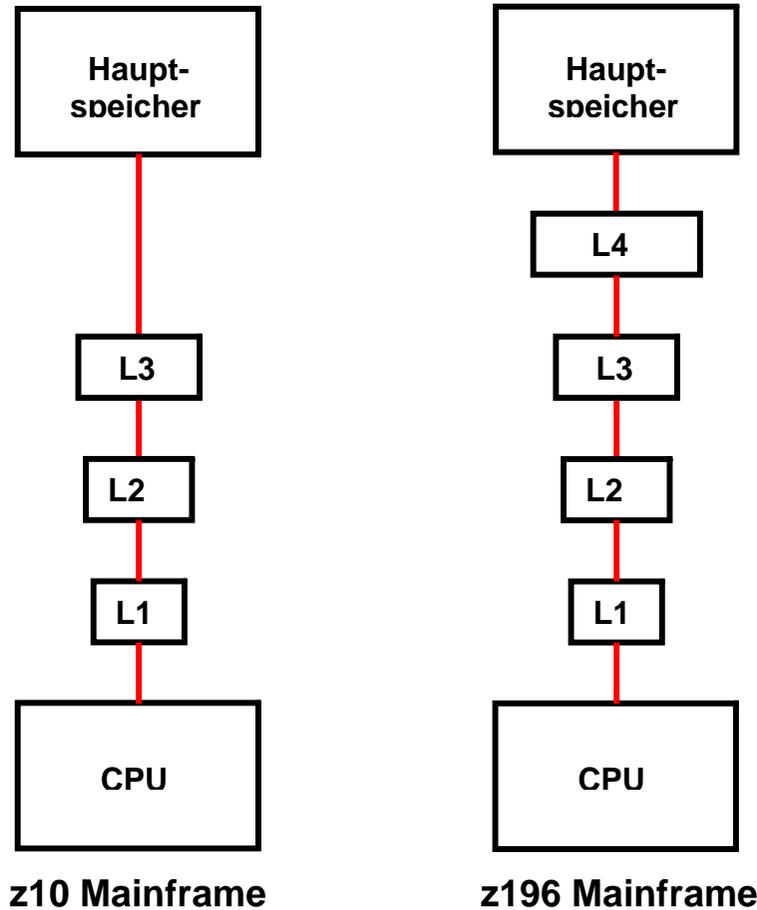
Der Befehls-cache wird häufig als I-Cache (Instruction Cache) und der Datencache als D-Cache bezeichnet.



## Second Level Cache

Es existieren zahlreiche SRAM Technologien um Caches zu implementieren. Heute ist es üblich, mit einer Cache Hierarchie zu arbeiten. Hierbei wird ein „Level 1“ Cache (L1) mit besonders schnellen, aber teuren und platzaufwendigen SRAM Speicherzellen implementiert. Ein „Level 2“ Cache (L2) verwendet langsamere aber dafür kostengünstigere SRAM Zellen.

Hierbei wird häufig, z.B. beim z9 Mainframe, der L1 Cache als Split Cache und der L2 Cache als Uniform Cache implementiert



## Mainframe Cache Hierarchien

Moderne Cache Hierarchien sind noch komplizierter. Gezeigt sind die 3-stufige z10 Cache Hierarchie und die 4-stufige z196 Cache Hierarchie. Sinnvoll wurde dies durch die Erfindung einer neuartigen als eDRAM bezeichneten Cache Speicherzellen-Technologie, die im z196 Mainframe die SRAM Zellen in den L3 und L4 Caches ersetzt.