

Pamięć Cache

m@v€K pud3£k0

Urządzenia Techniki Komputerowej

Spis treści

- Definicja pamięci Cache
- Zasada działania pamięci cache
- Integralność danych w pamięci
- Hierarchia pamięci cache

- Test pamięci w programach

Pamięć Cache

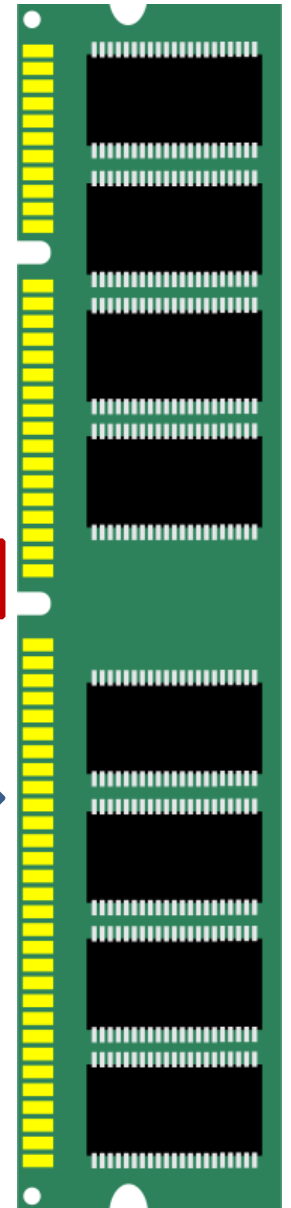
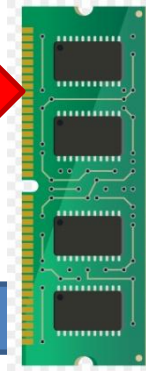
- Szybka pamięć podręczna (ang. Cache Memory), służąca do przechowywania często używanych danych.
- Stanowi bufor pomiędzy wolną dynamiczną pamięcią operacyjną (DRAM) a procesorem.
- Typ pamięci Static RAM (SRAM)
- Wielkość jej wynosi 1 MB – 16 MB

Struktura pamięci RAM i Cache

Pamięć Cache



Procesor



Pamięć RAM

Pamięć Cache w CPU-Z

The screenshot shows the CPU-Z application window. The 'CPU' tab is selected. The processor is identified as an Intel Celeron G530, Sandy Bridge architecture, Socket 1155 LGA package, 32 nm technology, and 2.40 GHz clock speed. The cache information is detailed in the 'Cache' section at the bottom right of the main window.

Cache	Size	Way
L1 Data	2 x 32 KBytes	8-way
L1 Inst.	2 x 32 KBytes	8-way
Level 2	2 x 256 KBytes	8-way
Level 3	2 MBytes	8-way

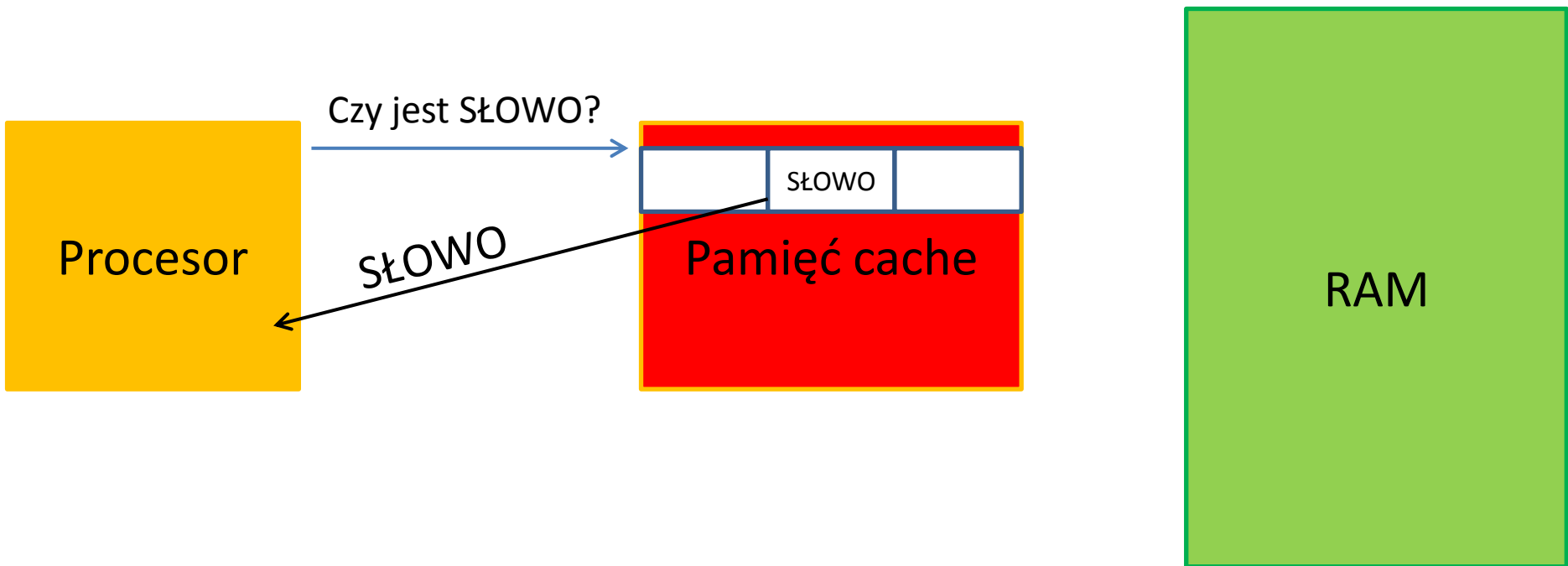
Cache	Size	Way
L1 Data	2 x 32 KBytes	8-way
L1 Inst.	2 x 32 KBytes	8-way
Level 2	2 x 256 KBytes	8-way
Level 3	2 MBytes	8-way

Szybkość pamięci i procesora

- Wzrost wydajności spowodowała, że procesor stał się zbyt szybki dla pamięci RAM.
- Odwoływanie się tylko do niej spowodowałoby, że większość czasu byłby niewykorzystany.
- Dodano dodatkową pamięć Cache, która jest szybsza niż RAM i stanowi magazyn często używanych danych.

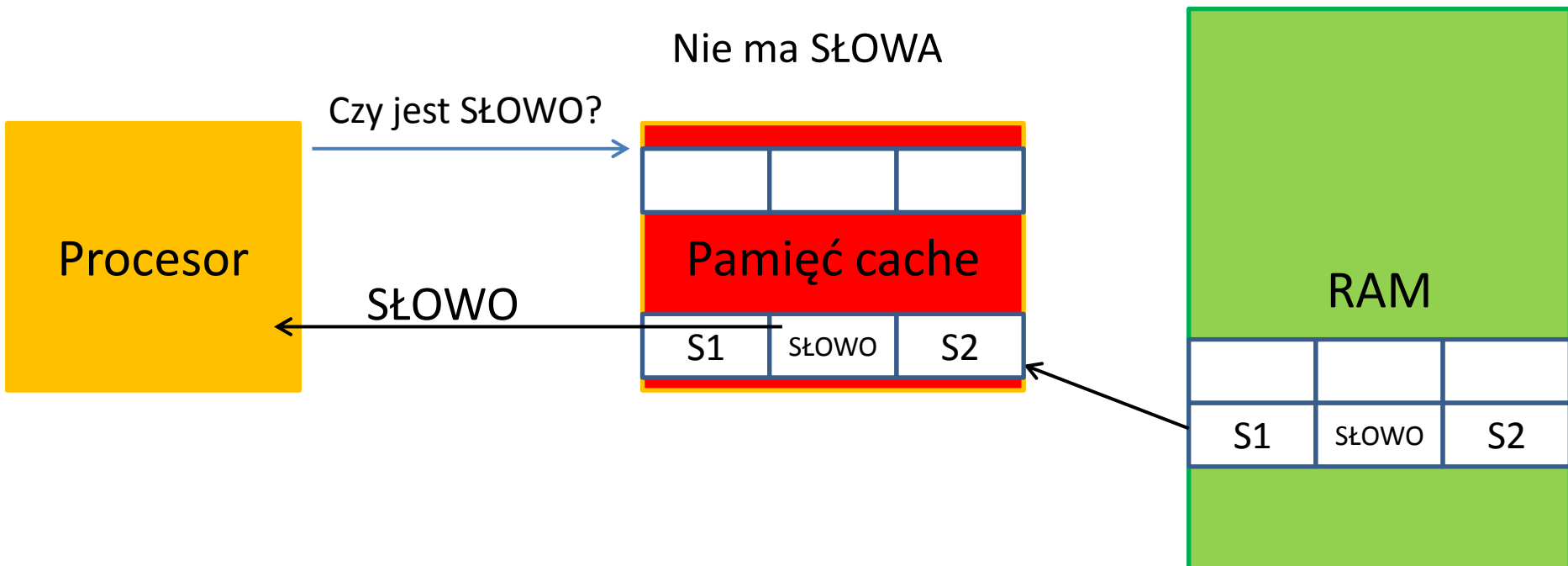
ZASADA DZIAŁANIA PAMIĘCI CACHE

Zasada działania pamięci podręcznej



- Procesor potrzebuje danego słowa.
- Szuka go najpierw w pamięci Cache.
- Gdy jest znajdzie, to ściąga je do obróbki.

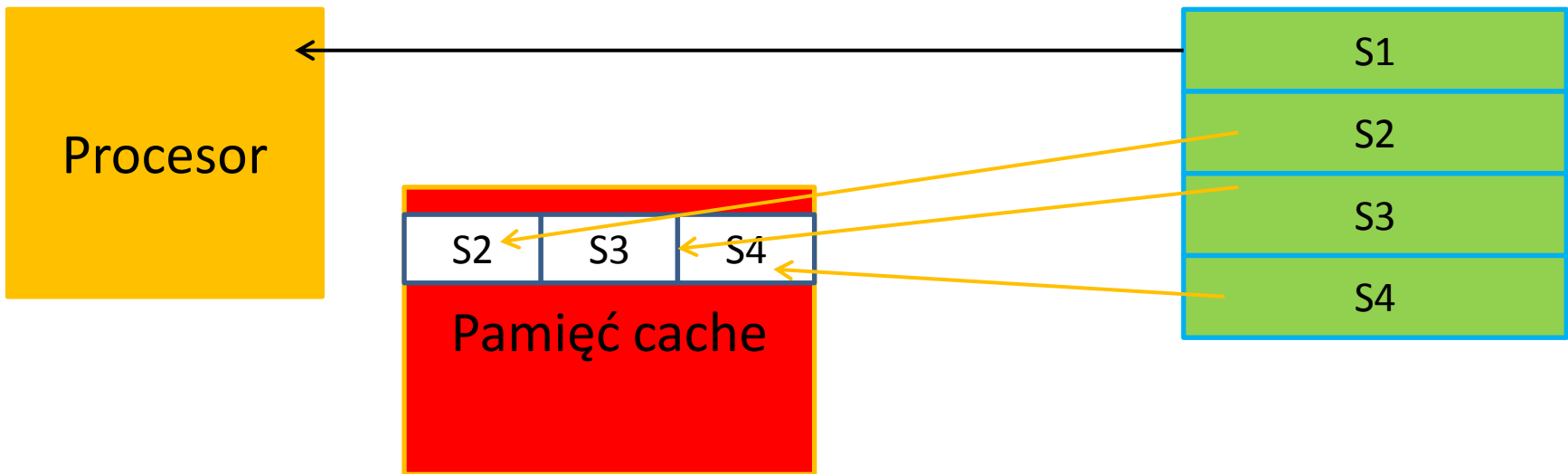
Zasada działania pamięci podręcznej



- Gdy brak danego słowa w pamięci Cache, sięga się do pamięci RAM.
- Blok pamięci RAM zawierający ustaloną liczbę słów (w tym szukane) jest wczytywany do pamięci podręcznej.
- Szukane słowo jest dostarczane do procesora.

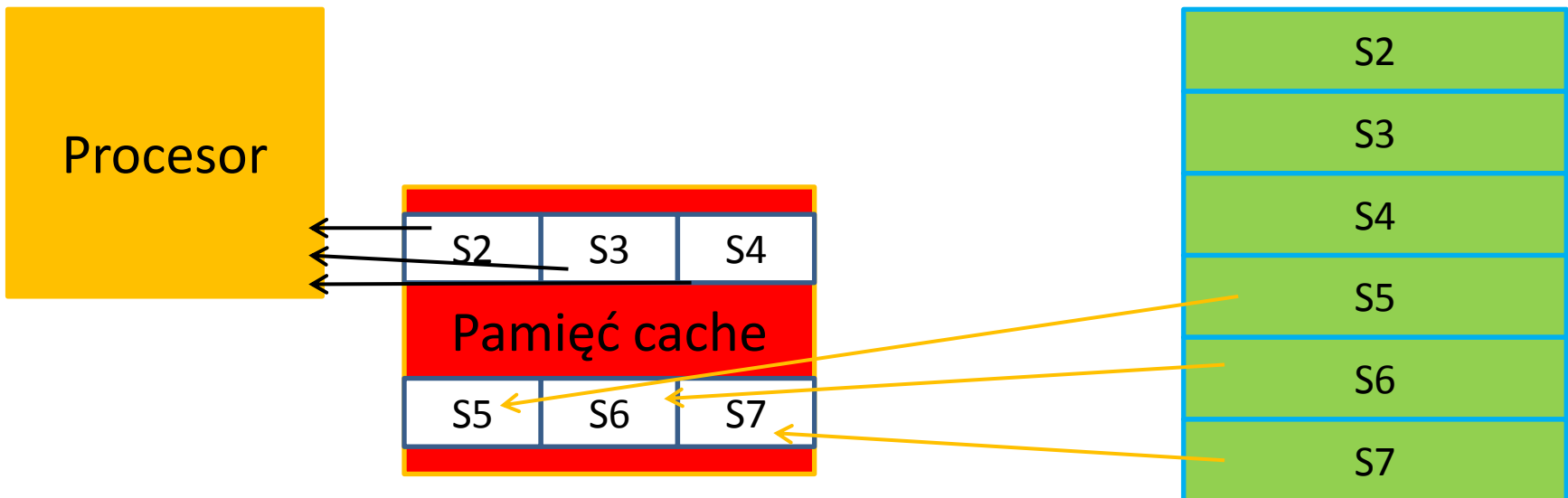
Integralność danych

- Zazwyczaj dane w pamięci RAM są połączone w logiczne grupy.
- Wykonywane są raczej sekwencyjnie.
- Pozwala to przewidzieć jakie dane i polecenia będą potrzebne za chwilę (zazwyczaj następne po aktualnie wykorzystanych) co pozwala na załadowanie ich wcześniej do pamięci cache.



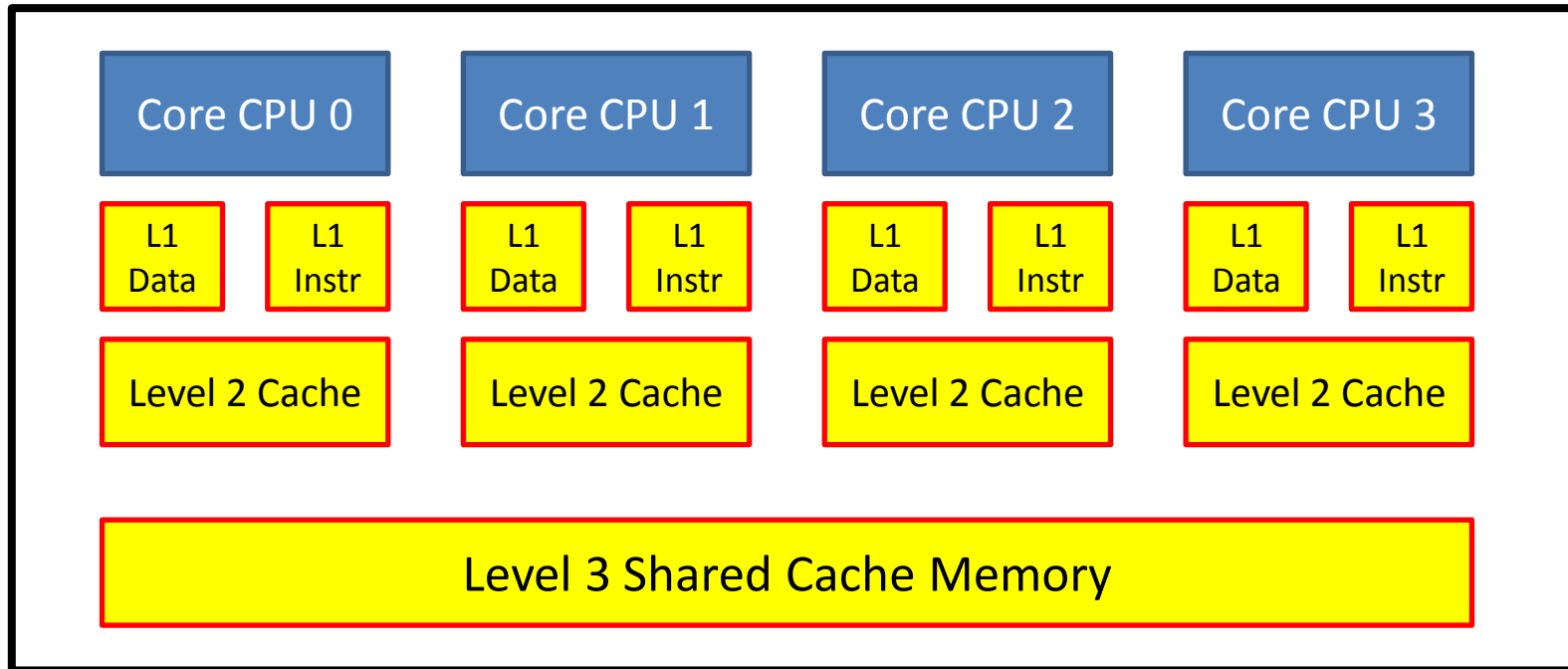
Integralność danych

- Po kolejne dane procesor sięgnie do pamięci podręcznej, co znacznie przyspieszy obliczenia.
- Tymczasem do pamięci cache można kopiować kolejne dane z RAM.



HIERARCHIA PAMIĘCI CACHE

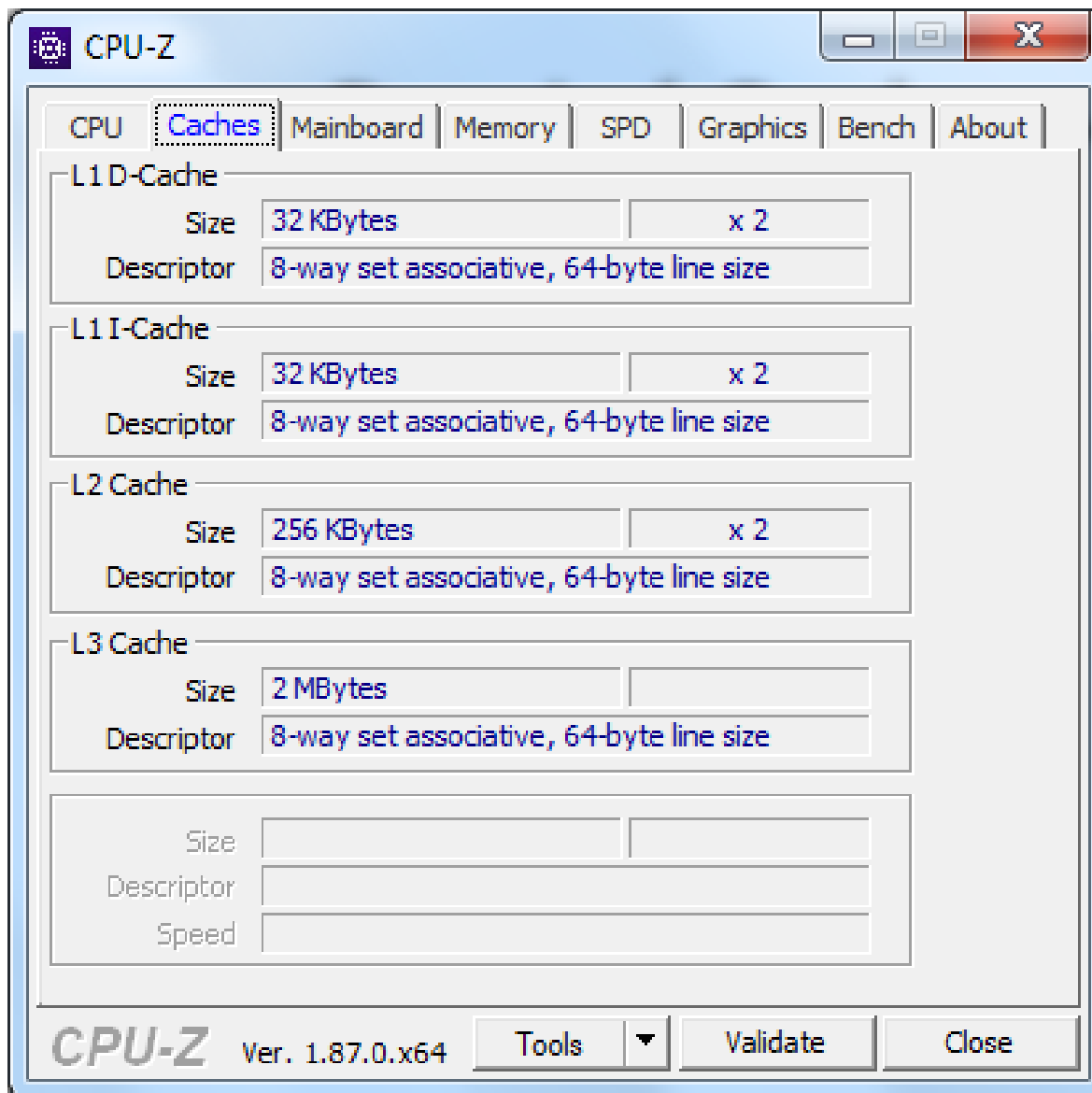
Hierarchia pamięci cache



Pamięć RAM

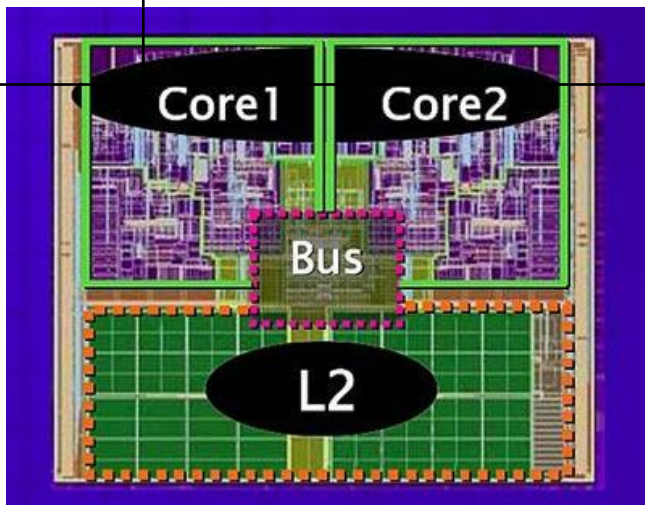
Nośniki pamięci trwałe

Pamięć Cache – CPU-Z

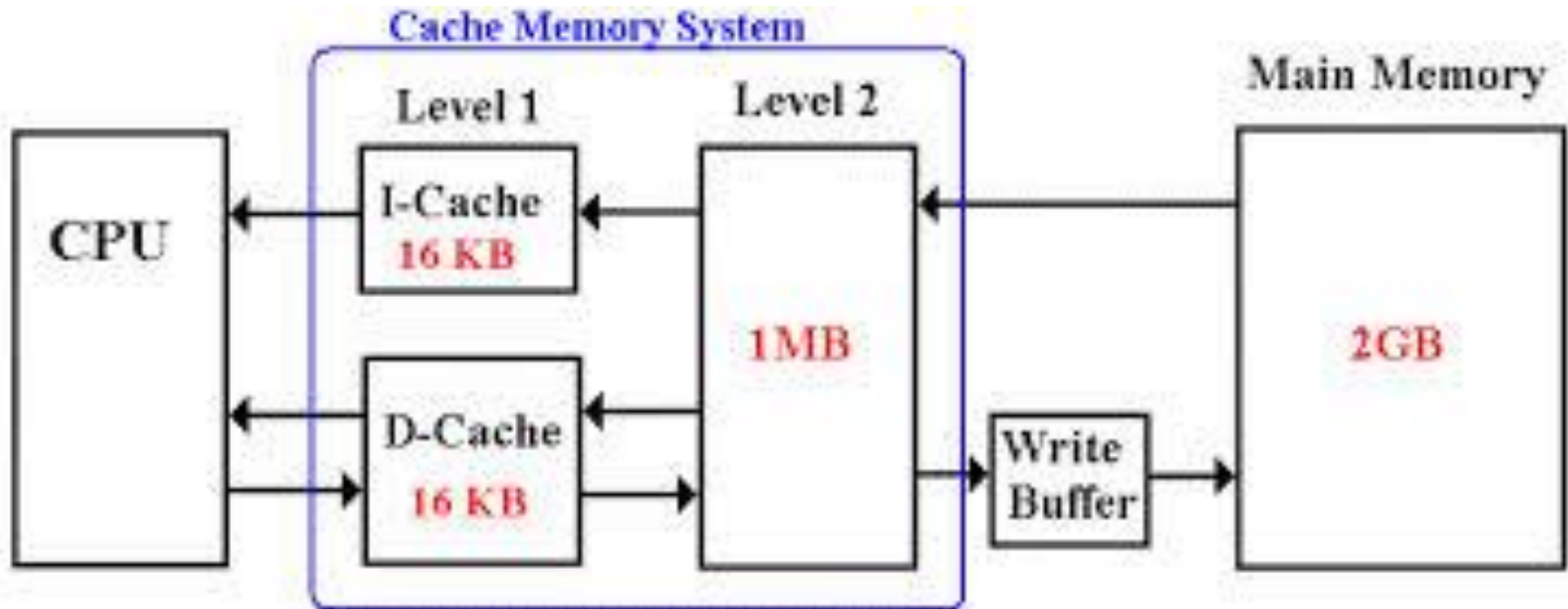


Poziomy pamięci Cache

Poziom pamięci cache	Umiejscowienie	Wielkość
L1	Zintegrowane z jądrem procesora	32kB-128kB
L2	Umieszczona na płycie głównej Może być wewnątrz wielordzeniowego procesora	1MB-8 MB
L3	Umieszczona na płycie głównej	4MB-32 MB W szybkich procesorach serwerowych jest dzielona przez wszystkie rdzenie/procesory

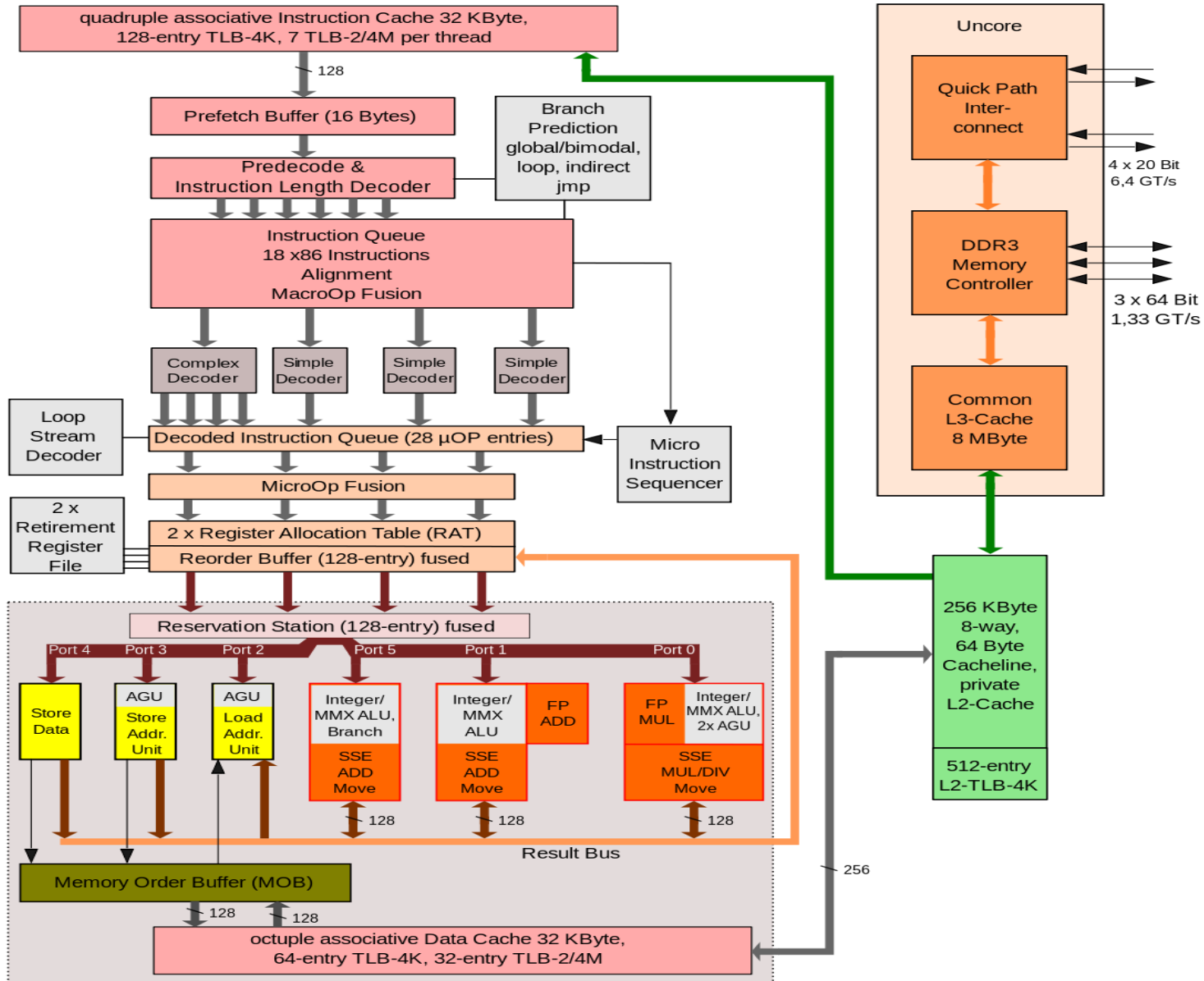


Realizacja cache



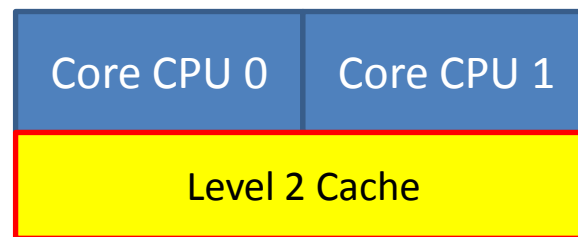
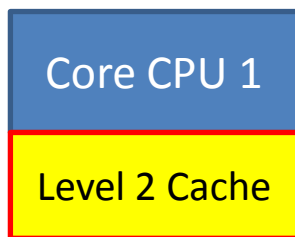
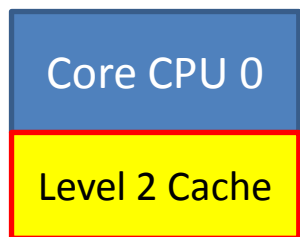
ARCHITEKTURA PAMIĘCI CACHE

Intel Nehalem microarchitecture



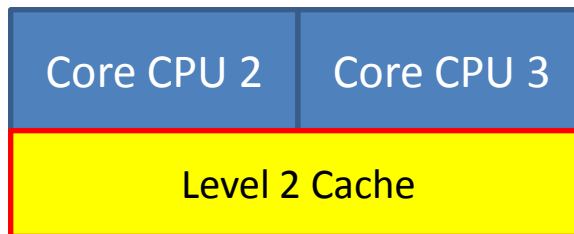
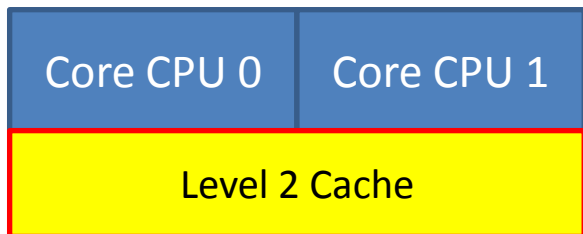
GT/s: gigatransfers per second

Architektura pamięci cache L2



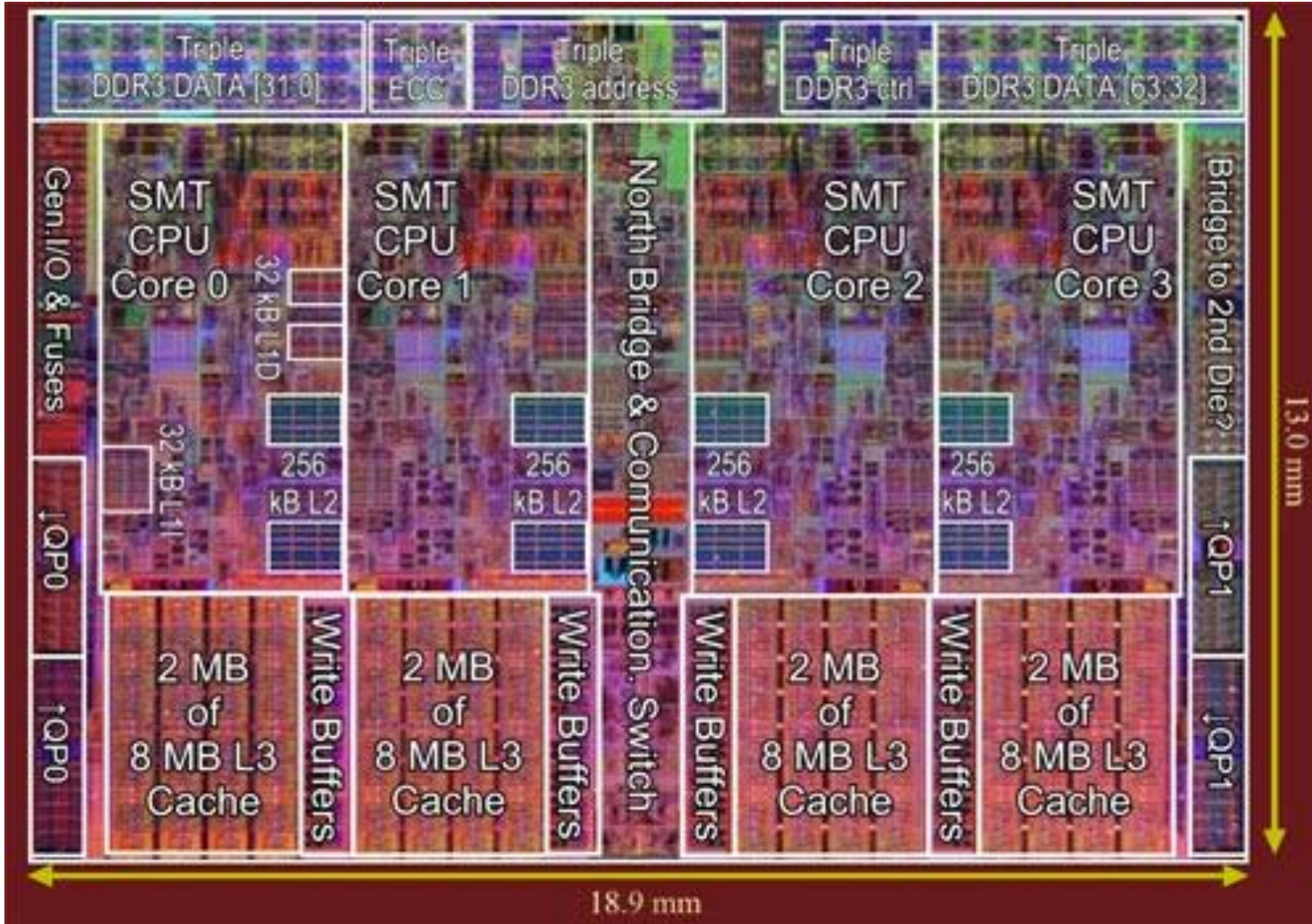
Oddzielne pamięci cache
(Separated Cache)
AMD, Pentium D

Wspólna pamięć cache
(Shared Cache)
Core Duo, Core 2 Duo

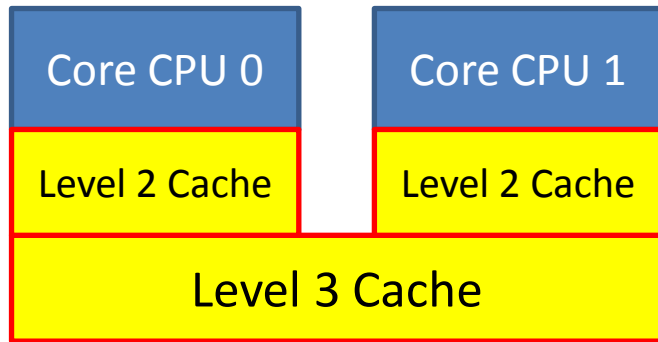


Współbieżna pamięć cache
(Current Cache)
Core 2 Quad, Core 2 Extreme QX

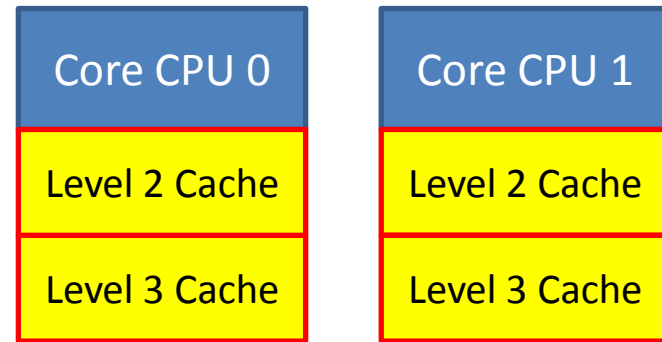
L1, L2 Cache



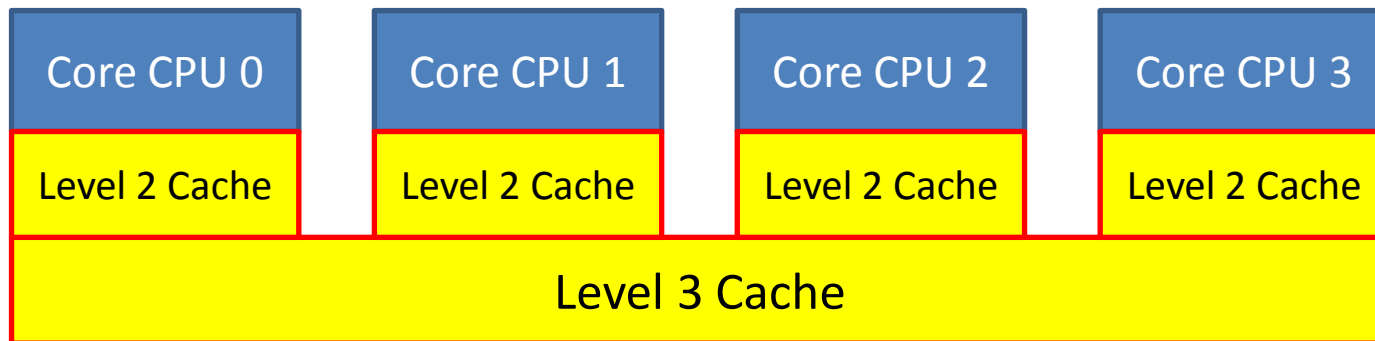
Architektura pamięci cache L3



Dual-Core

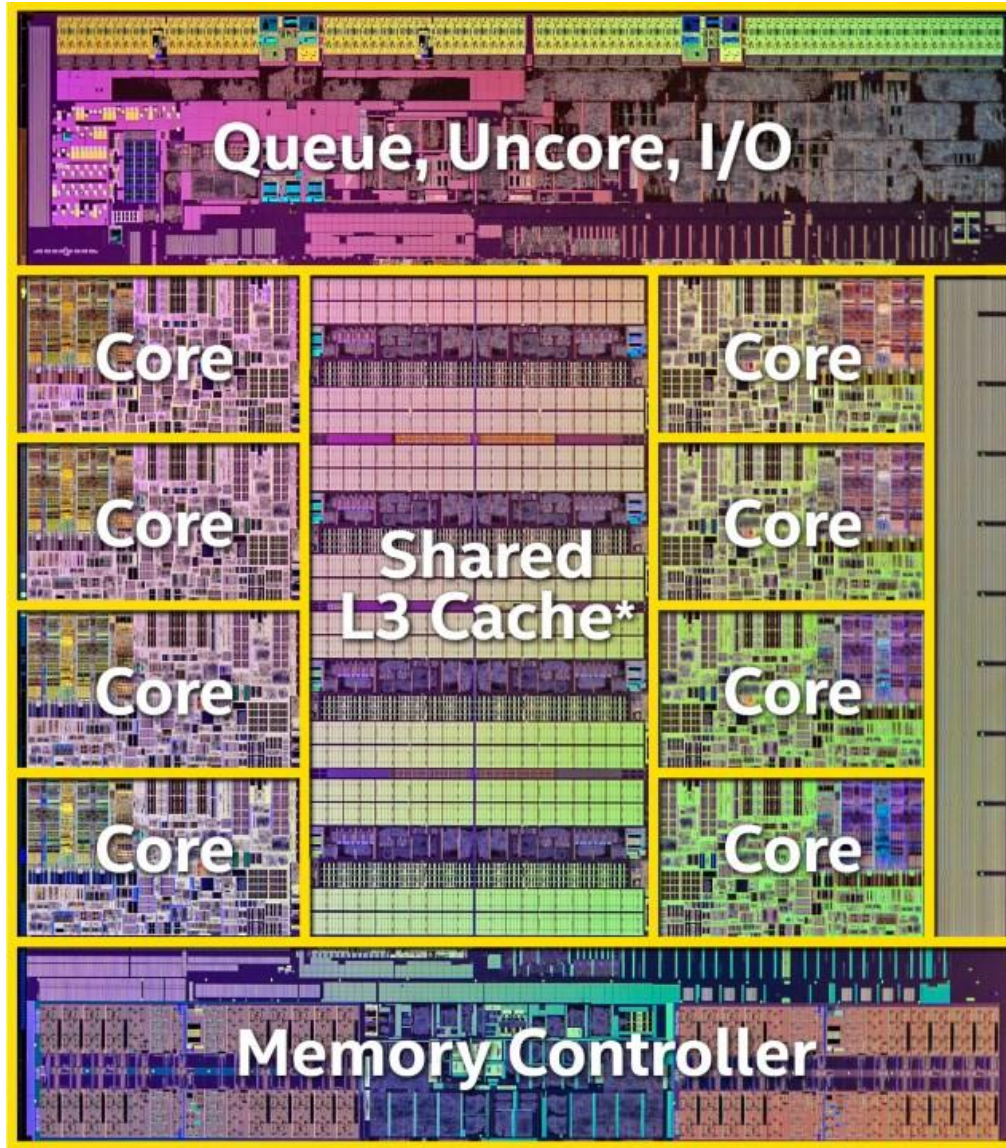


Itanium 2

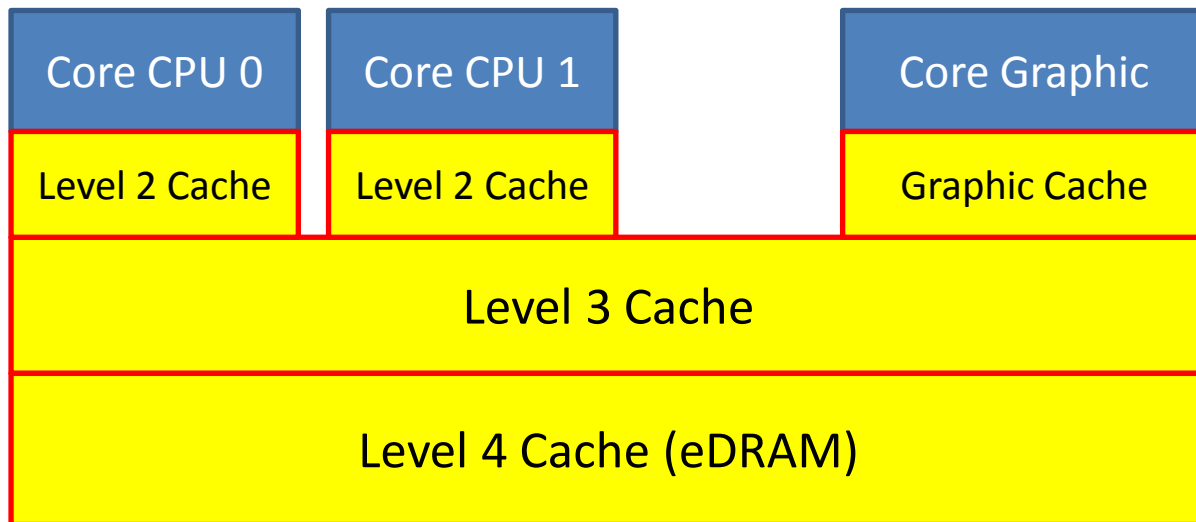


Quad-Core

L3 Cache

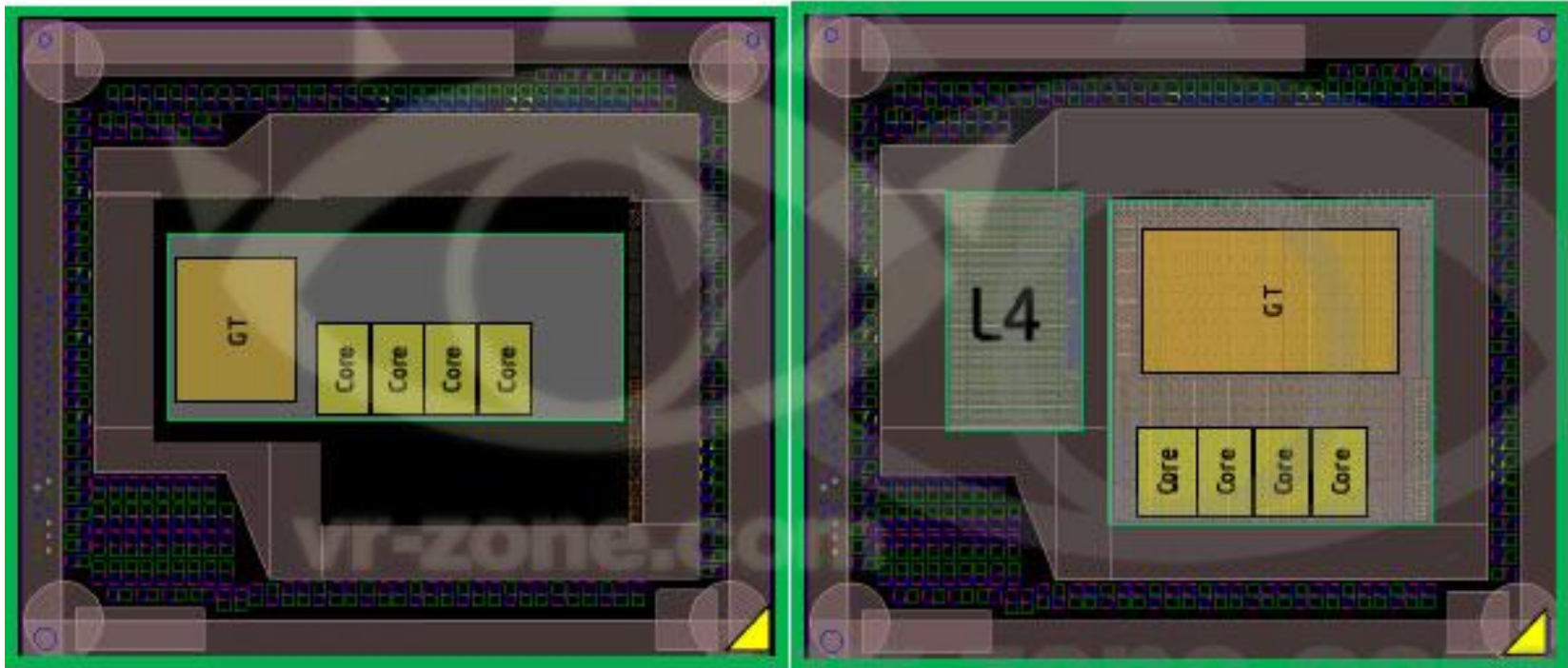


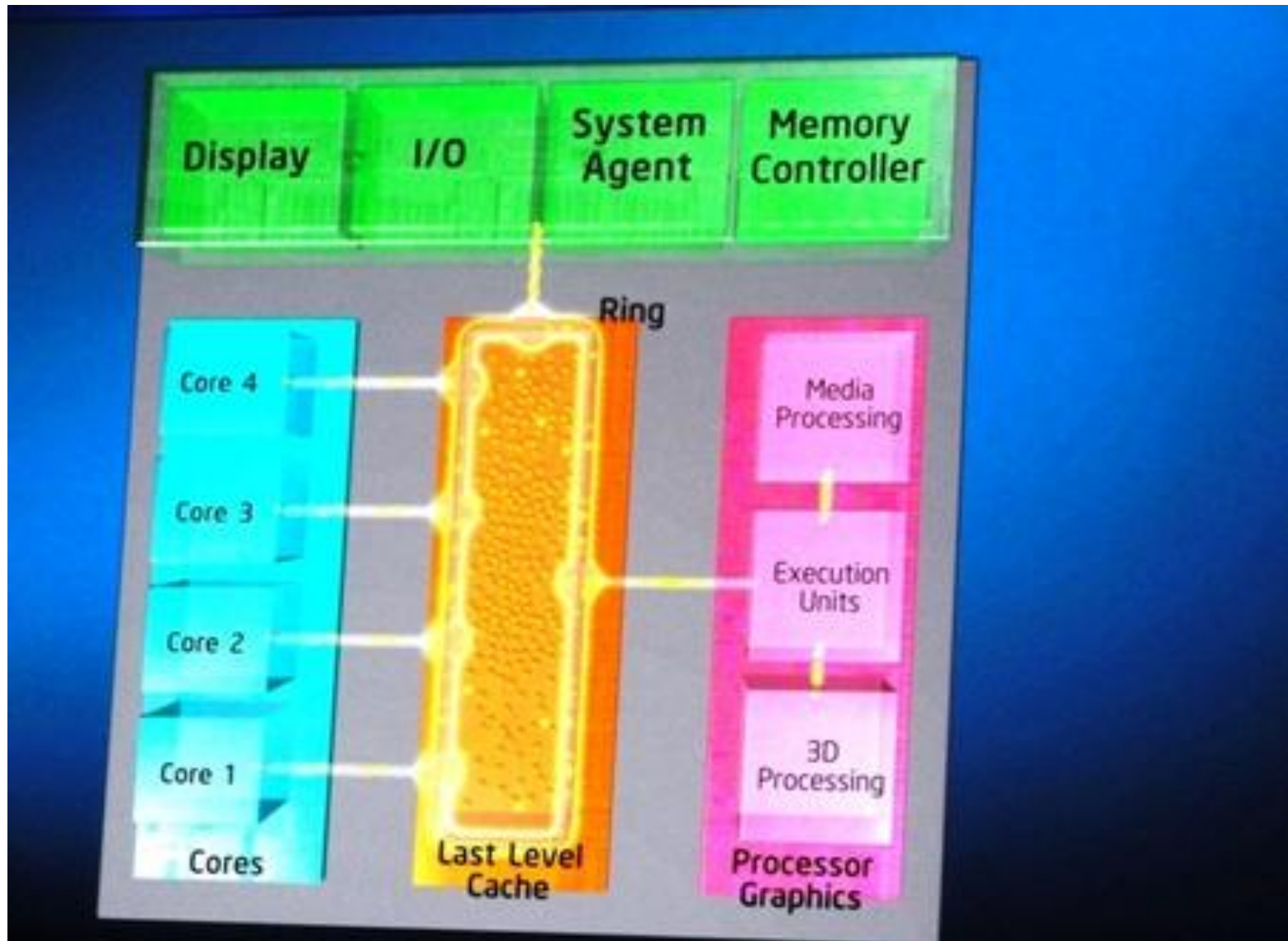
Architektura pamięci cache L4



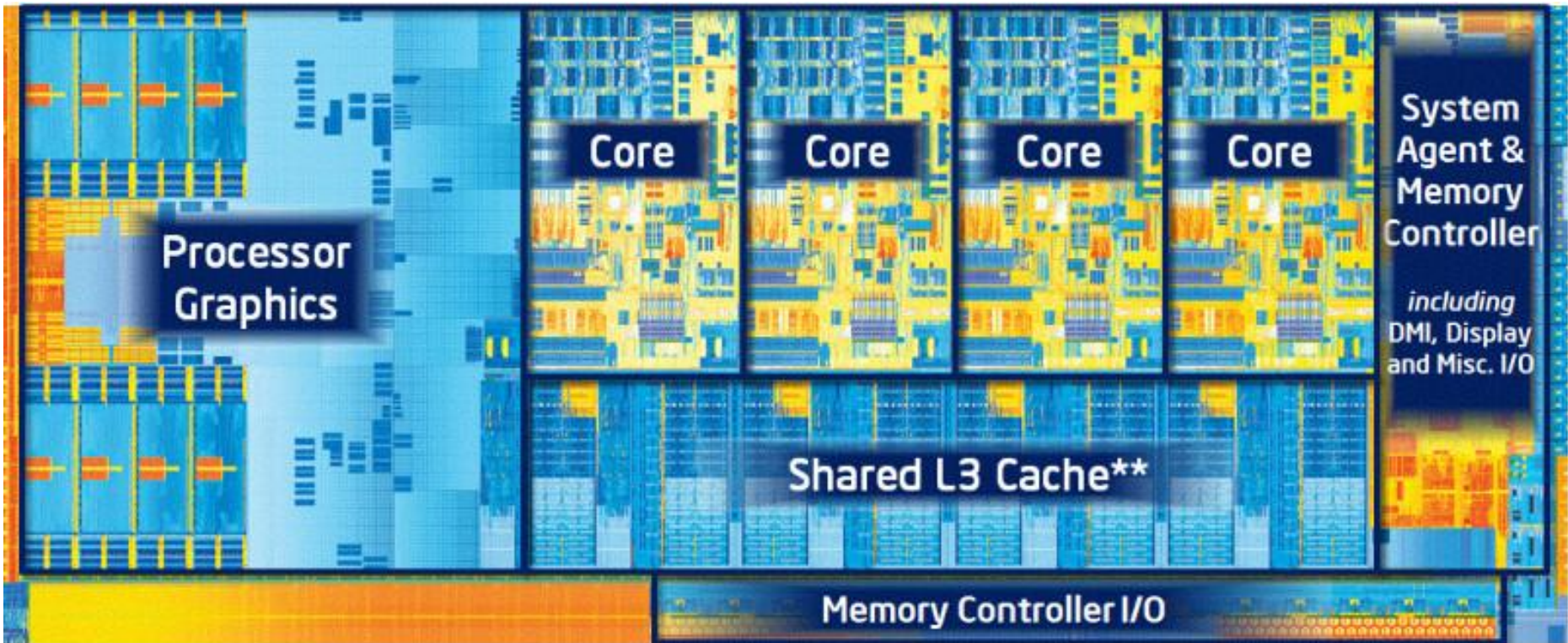
Intel Haswell, SkyLake

Intel Haswell

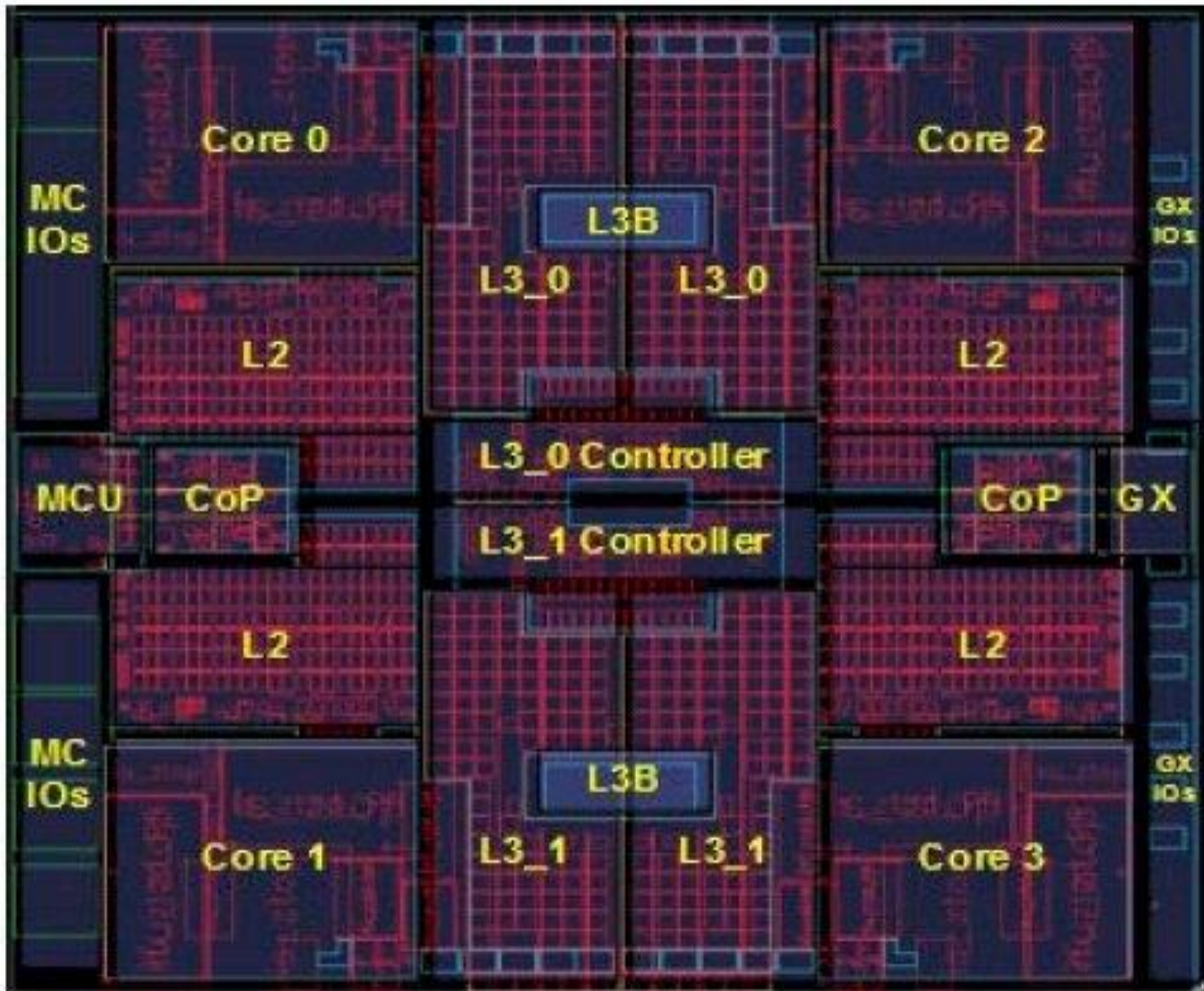




Cache w procesorze



IBM's zEnterprise 196 CPU



Pamięć Cache – CPU-Z

The screenshot shows the CPU-Z application window. The 'CPU' tab is selected. The processor information is as follows:

- Name: Intel Celeron G530
- Code Name: Sandy Bridge
- Max TDP: 65.0 W
- Package: Socket 1155 LGA
- Technology: 32 nm
- Core Voltage: 1.044 V
- Specification: Intel® Celeron® CPU G530 @ 2.40GHz
- Family: 6
- Model: A
- Stepping: [blank]
- Ext. Family: 6
- Ext. Model: 2A
- Revision: [blank]
- Instructions: MMX, SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, EM

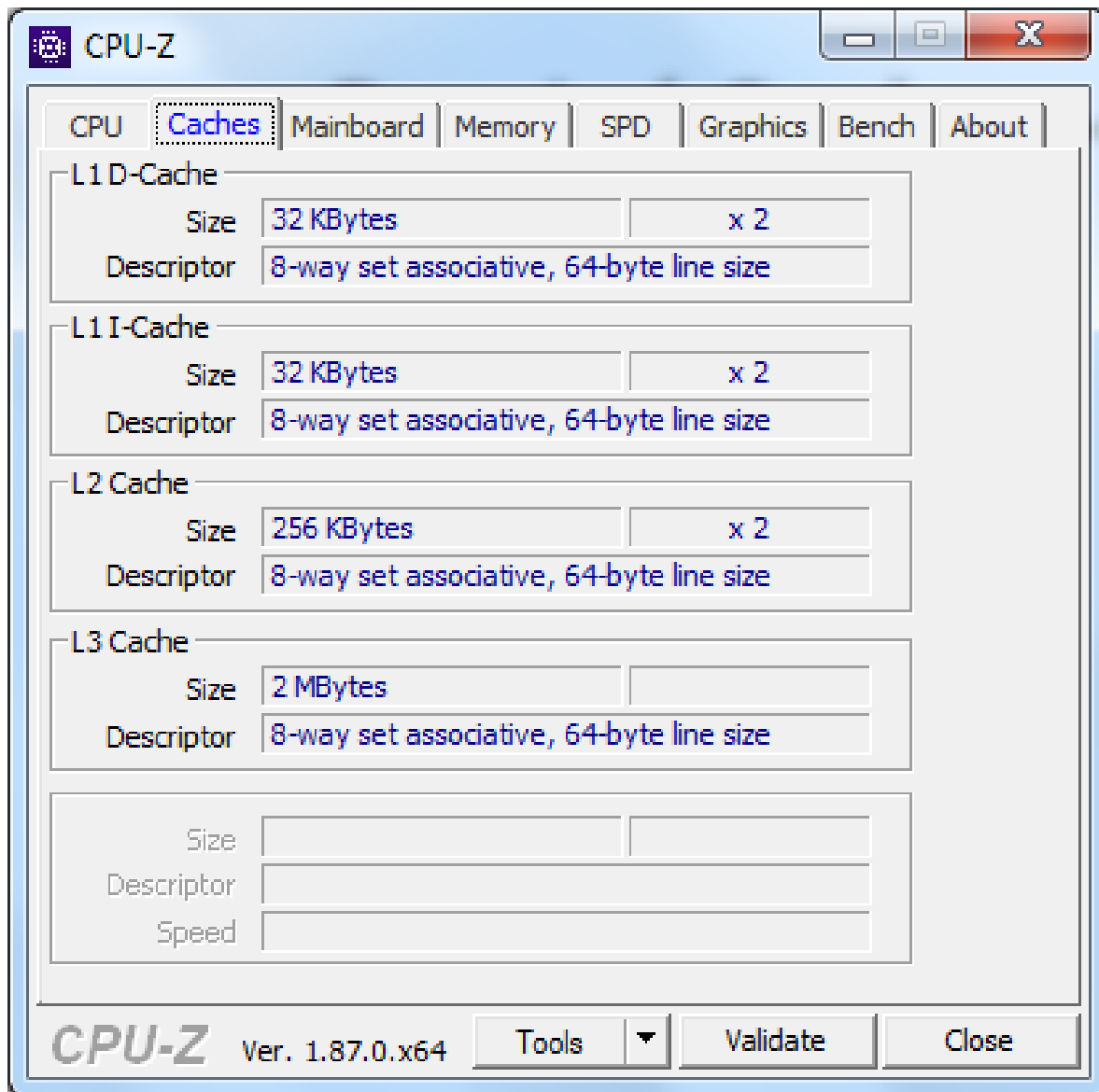
The 'Cache' tab is also visible, showing the following information:

Cache	Size	Way
L1 Data	2 x 32 KBytes	8-way
L1 Inst.	2 x 32 KBytes	8-way
Level 2	2 x 256 KBytes	8-way
Level 3	2 MBytes	8-way

At the bottom of the CPU-Z window, the 'Cores' field shows 2 and the 'Threads' field shows 2. The version is Ver. 1.87.0.x64.

Cache	Size	Way
L1 Data	2 x 32 KBytes	8-way
L1 Inst.	2 x 32 KBytes	8-way
Level 2	2 x 256 KBytes	8-way
Level 3	2 MBytes	8-way

Pamięć Cache – CPU-Z

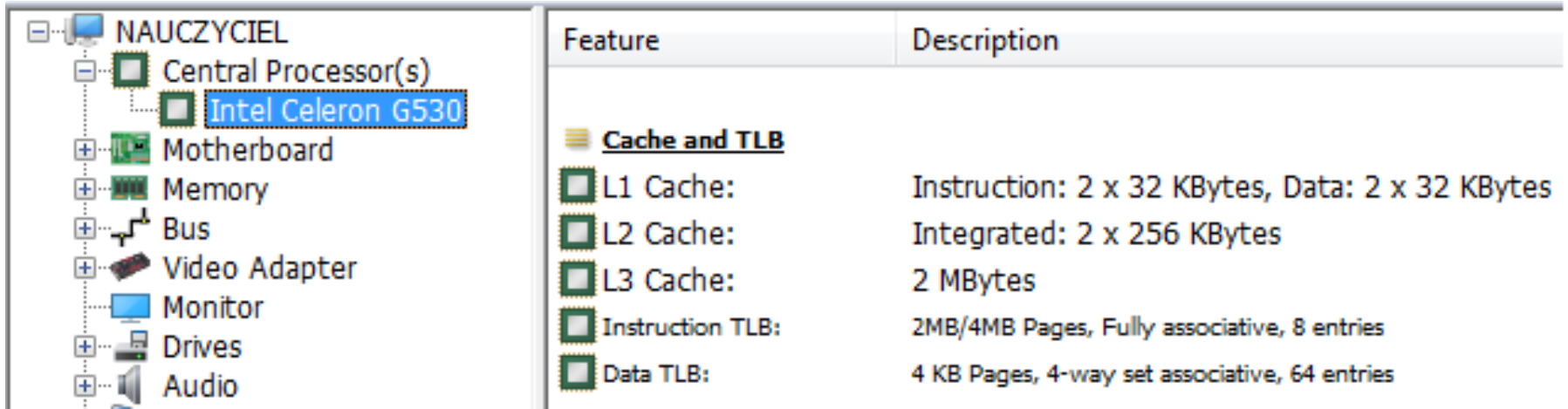


The image shows a screenshot of the CPU-Z application window, specifically the 'Caches' tab. The window title is 'CPU-Z'. The 'Caches' tab is selected, and the following information is displayed:

Cache Type	Size	Descriptor
L1 D-Cache	32 KBytes x 2	8-way set associative, 64-byte line size
L1 I-Cache	32 KBytes x 2	8-way set associative, 64-byte line size
L2 Cache	256 KBytes x 2	8-way set associative, 64-byte line size
L3 Cache	2 MBytes	8-way set associative, 64-byte line size

At the bottom of the window, the version is 'CPU-Z Ver. 1.87.0.x64'. There are buttons for 'Tools', 'Validate', and 'Close'.

Pamięć Cache – HWInfo



The screenshot displays the HWInfo application interface. On the left, a tree view shows system components: NAUCZYCIEL, Central Processor(s) (expanded to Intel Celeron G530), Motherboard, Memory, Bus, Video Adapter, Monitor, Drives, and Audio. The right pane shows a table of hardware features and their descriptions.

Feature	Description
Cache and TLB	
L1 Cache:	Instruction: 2 x 32 KBytes, Data: 2 x 32 KBytes
L2 Cache:	Integrated: 2 x 256 KBytes
L3 Cache:	2 MBytes
Instruction TLB:	2MB/4MB Pages, Fully associative, 8 entries
Data TLB:	4 KB Pages, 4-way set associative, 64 entries

Porównanie osiągnięć pamięci Cache

Cache and Memory Benchmark

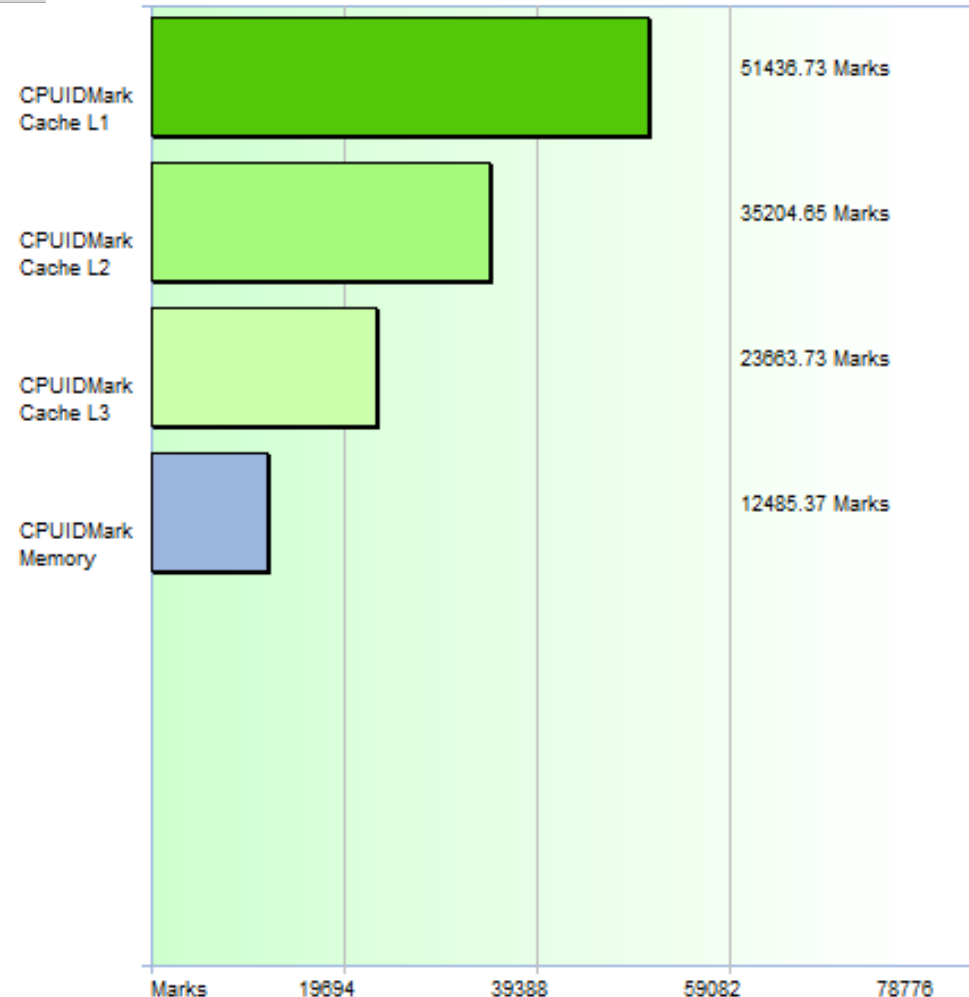
Item	Description
CPUIDMark Cache L1	51436.73 Marks
CPUIDMark Cache L2	35204.65 Marks
CPUIDMark Cache L3	23663.73 Marks
CPUIDMark Memory	12485.37 Marks
Latency Cache L1	4 cycles
Latency Cache L2	12 cycles
Latency Memory	122 cycles

CPUIDMark Memory :
30698

Memory Information :

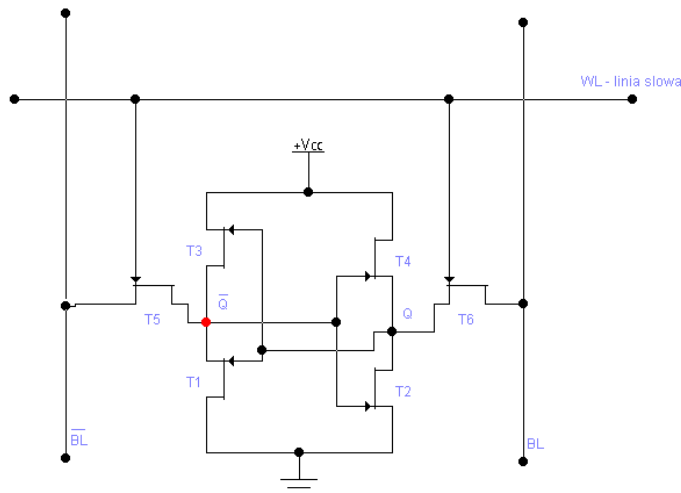
L1 Cache : 2 x 64 KB
L2 Cache : 2 x 256 KB
L3 Cache : 2048 KB
Total Memory : 4016 MB

Cache and Memory Benchmark
Bandwidth

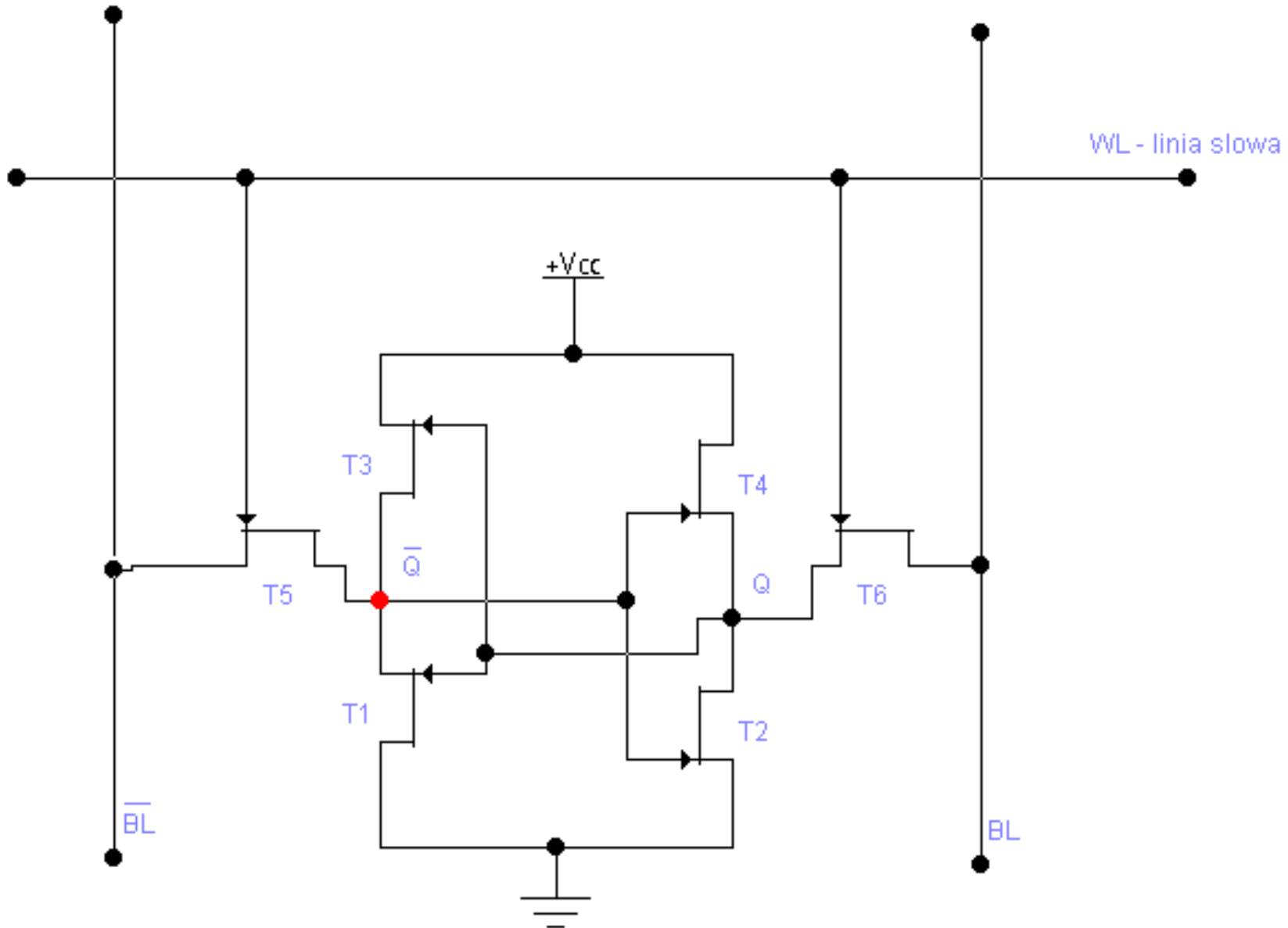


Pamięć statyczna – S-RAM

- Pamięć statyczna – S-RAM (ang. Static Random Access Memory)
 - Pamięć S-RAM przechowuje dane tak długo, jak długo włączone jest zasilanie.
- Każdy bit przechowywany jest w układzie z czterech tranzystorów, które tworzą przerzutnik, oraz z dwóch tranzystorów sterujących.
- Taka struktura umożliwia szybkie odczytanie bitu i nie wymaga odświeżania.



Pamięć statyczna – S-RAM

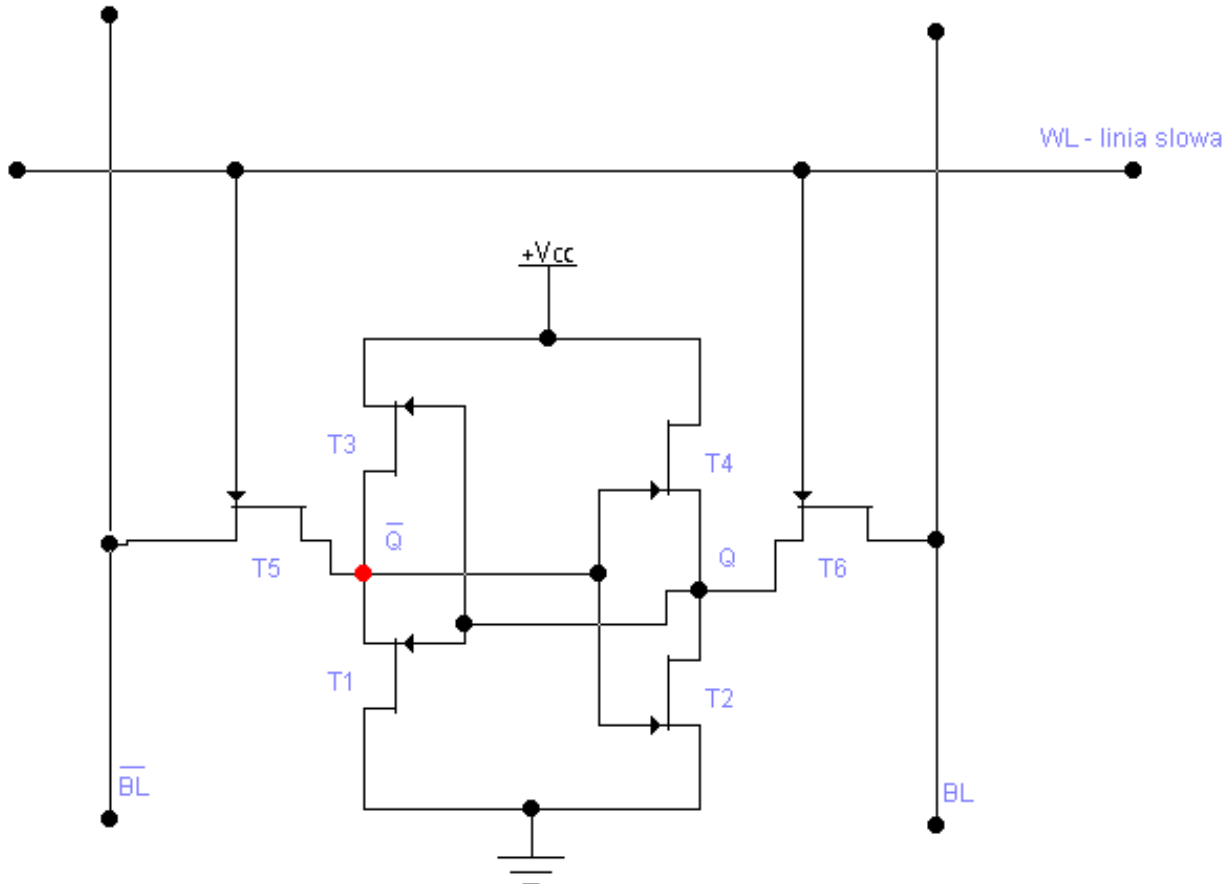


Pamięć statyczna – S-RAM

- Do przechowywania każdego bitu w pamięci statycznej RAM wykorzystywane są dwa krzyżowo sprzężone inwertory, zbudowane z tranzystorów polowych CMOS oznaczonych jako T_1 , T_2 i T_3 , T_4 .
 - Inwertory tworzą prosty przerzutnik bistabilny, który posiada dwa stabilne stany wykorzystywane do zapisu poziomów logicznych 0 i 1.
- Dwa tranzystory T_5 i T_6 służą do sterowania dostępem do komórki podczas zapisu i odczytu danych. Są one podłączone do linii słowa WL.
- Odpowiednio wysterowane sygnałem na tej linii tranzystory T_5 i T_6 łączą wyjście Q przerzutnika z linią bitu BL oraz wyjście komplementarne $\sim Q$ z komplementarną linią $\sim BL$.
 - Dwie linie bitu zwiększają poziom sygnału w stosunku do szumów (sygnałów zakłócających), które pojawiają się w strukturze pamięci półprzewodnikowych.
- Komórka pamięci statycznej RAM może być w jednym z trzech stanów:
 - Oczekiwania
 - Zapisu
 - Odczytu.

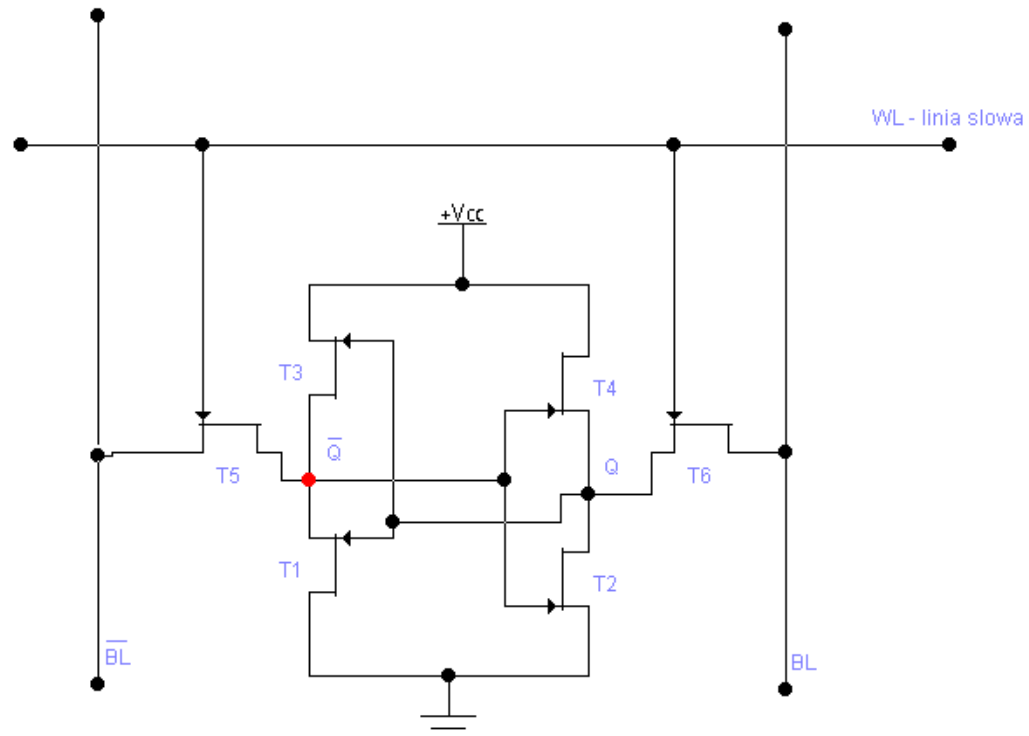
Stan oczekiwania S-RAM

- Gdy linia słowa WL nie jest wystereowana odpowiednim napięciem, tranzystory T_5 i T_6 separują wyjścia przerzutnika od linii bitów.
- Przerzutnik, zbudowany z dwóch sprzężonych wzajemnie inwertorów, pamięta swój stan wewnętrzny ustawiony przy poprzedniej operacji zapisu.



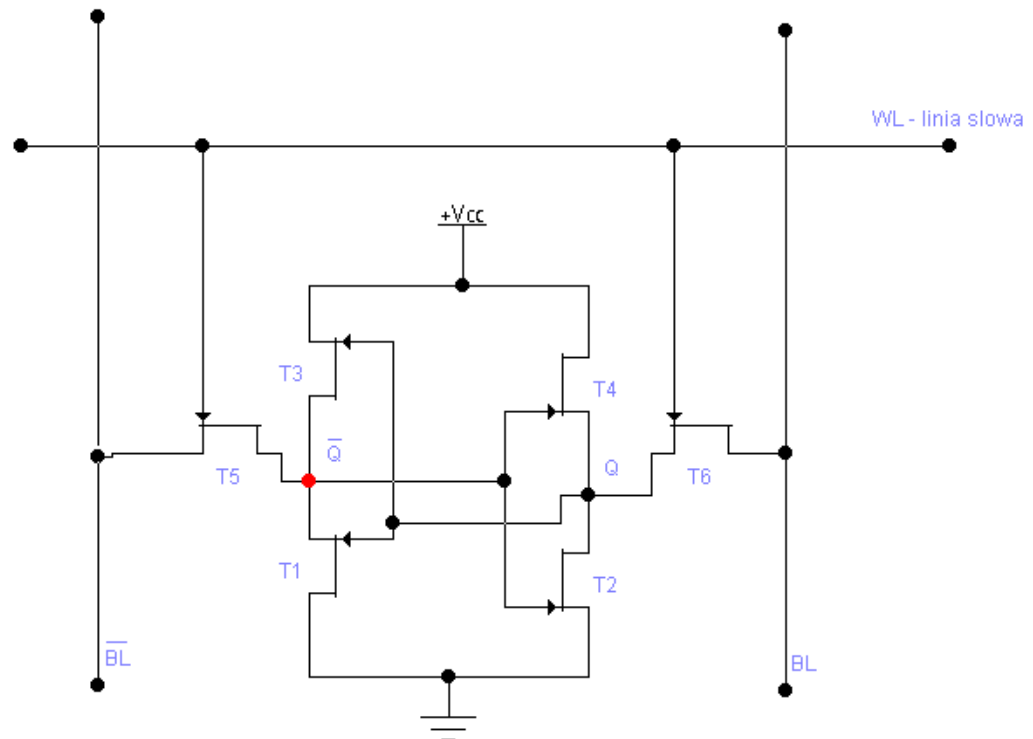
Zapis w S-RAM

- Przy zapisie ustawia się linie BL zgodnie z wartością zapisywanego bitu
 - (dla 1 $BL = 1$, $BL = 0$; dla 0 $BL = 0$, $BL = 1$).
- Następnie linia WL zostaje wysterowana i tranzystory T_5 , T_6 łączą wejścia inwertorów z liniami BL powodując zapis informacji w przerzutniku.
- Jest to możliwe, ponieważ sygnał na liniach BL i \overline{BL} jest wystarczająco mocny, aby wymusić zmianę stanu w tranzystorach $T_1...T_4$ przerzutnika, które są z reguły bardzo małe.



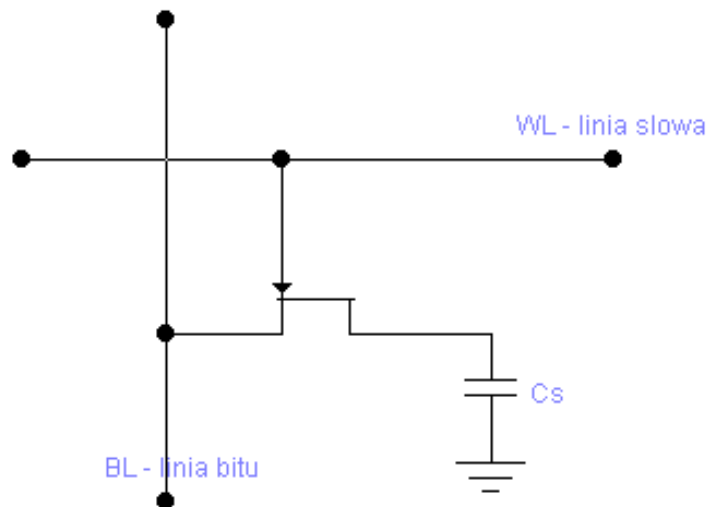
Odczyt z S-RAM

- Załóżmy, iż przerzutnik w komórce pamięci pamięta stan logiczny 1, czyli wyjście $Q = 1$, a $\bar{Q} = 0$ (odblokowane tranzystory M_4 i M_1 , zablokowane M_2 i M_3). Cykl odczytu rozpoczyna się przez naelektryzowanie obu linii bitów BL i \bar{BL} do wartości logicznej 1, a następnie przez wystawienie linii słowa WL, co spowoduje włączenie (odblokowanie) tranzystorów M_5 i M_6 . Stan wyjść Q i \bar{Q} zostaje przeniesiony na linie BL i \bar{BL} . Linia BL pozostaje w stanie 1, gdyż odblokowane tranzystory M_4 i M_6 łączą ją z napięciem V_{dd} . Z kolei napięcie linii \bar{BL} zostanie rozładowane do zera, ponieważ tranzystory M_1 i M_5 łączą ją z masą układu. W efekcie na linii BL pojawi się 1, a na \bar{BL} 0. Jeśli przerzutnik przechowuje wartość 0, to otrzymamy sytuację odwrotną.

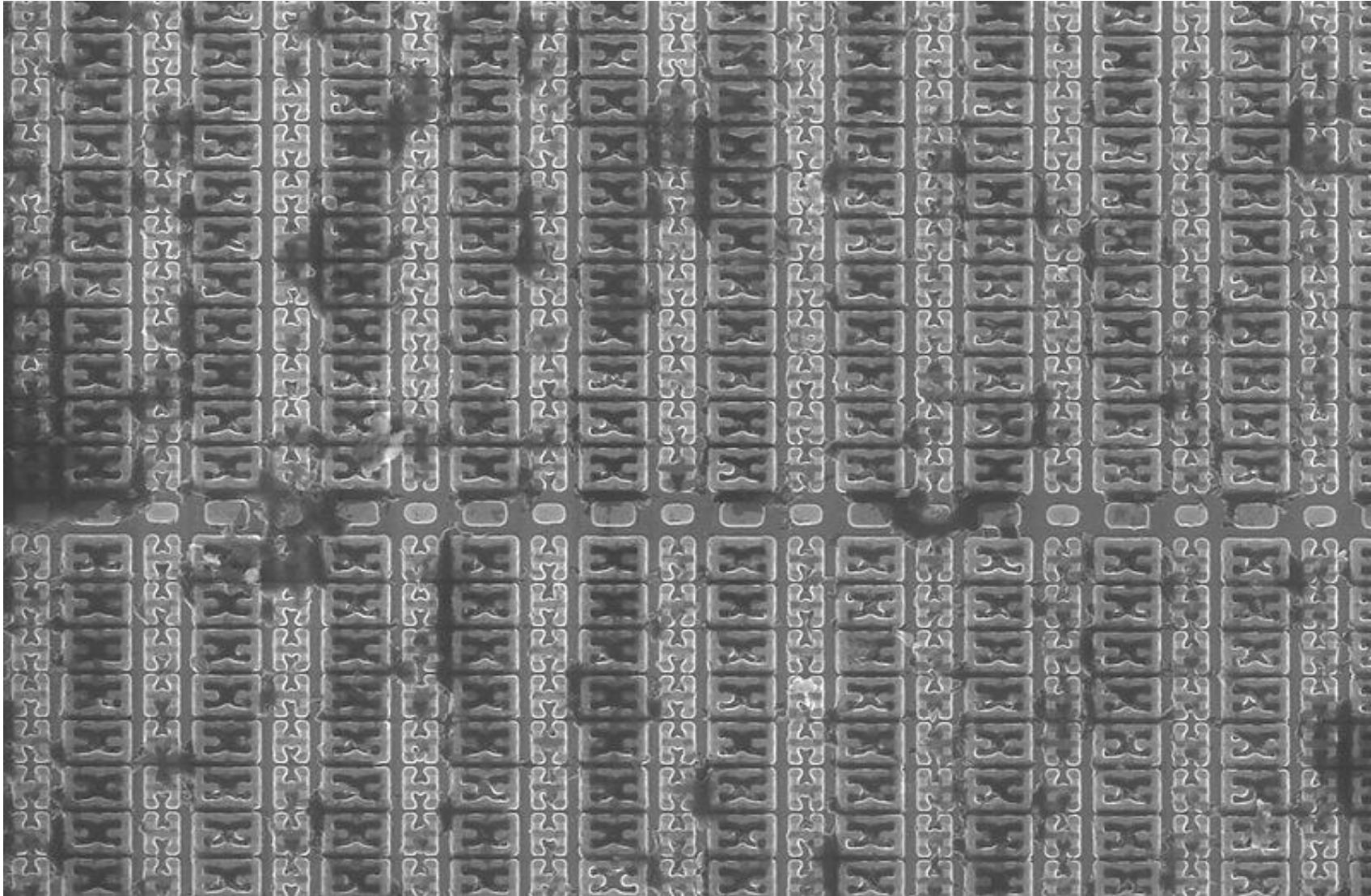


Odświeżanie S-RAM

- http://eduinf.waw.pl/inf/alg/002_struct/0043.php
#Komorka



S-GRAM



Historia pamięci Cache w PC

- Pamięć cache pojawiła się po raz pierwszy w procesorach Intel386 DX. Była zlokalizowana na zewnątrz CPU. Pamięć ta miała 64KB lub 128KB.
 - Niektóre płyty nie były w stanie obsłużyć jej, co pokazywało wyraźną różnicę w wydajności dzięki zastosowaniu pamięci podręcznej.
 - Model zapisu danych to write through gdzie procesor zapisywał wynik do pamięci podręcznej a z niej robiono od razu kopię do pamięci RAM.
- Kolejna generacja 486 miała wewnątrz procesora niewielką ilość (8KB) pamięci cache do dyspozycji CPU. Była to pamięć L1 (poziom1) lub wewnętrzna odróżniała ją to od pamięci zewnętrznej określanej jako L2. Ta druga miała 128KB lub 256KB.
 - Stosowano to architekturę write back. Wyniki operacji są przechowywane w pamięci cache ale uaktualnienie zawartości pamięci operacyjnej odbywa się w określonych momentach. Ten model stosowany jest do dziś.
- Pierwsze procesory Pentium miały 2 oddzielne pamięci cache L1-jedna do instrukcji, druga do danych. Miały one po 8KB.
 - Pamięć L2 w dalszym ciągu znajdowała się na płycie komputera. Rodzaj kości zależał od producenta płyty. Typowe wartości jej wynosiły 256KB lub 512KB.
 - AMD K5, K6 i K6-2 miały taką samą architekturę.
- AMD K6-III wprowadził dodatkowy poziom cache - L3 (level 3).
- Intel w procesorach P6 przesunął pamięć cache L2 do wnętrza CPU.
 - Przyspieszyło to pracę pamięci. Była taktowana wewnętrznym zegarem procesora, szybszym od taktowania zegarem bazowym.
- W 2012 roku Intel wprowadził pamięć cache poziomu 4 – L4. Była to pamięć eDRAM.

Algorytm zastępowania

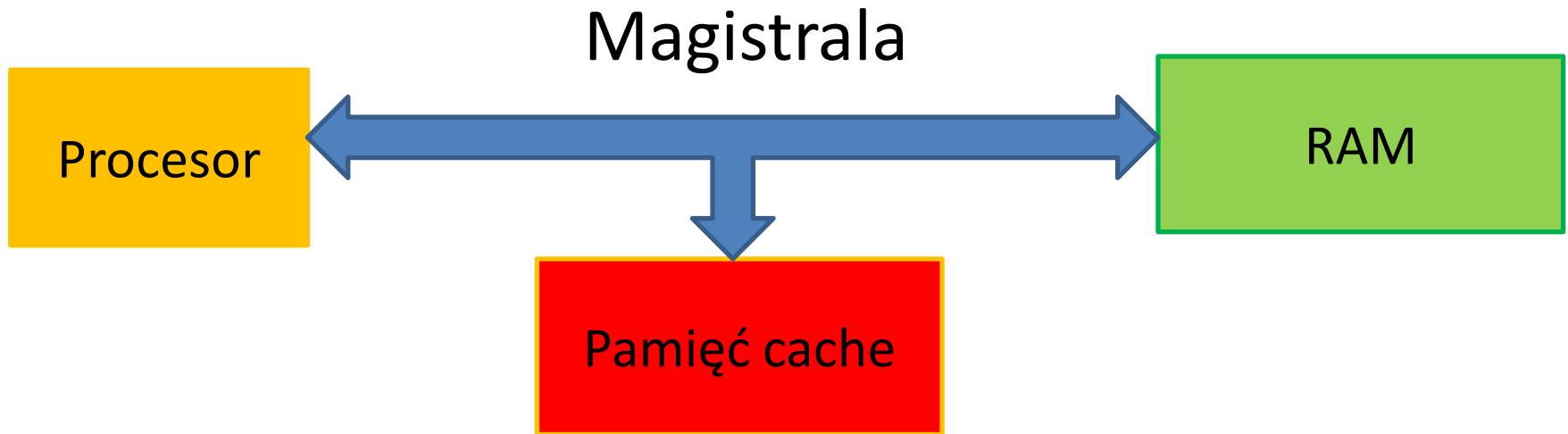
- W przypadku metody skojarzeniowej i sekcyjno-skojarzeniowej wymagany jest algorytm zastępowania.
- **Algorytm „najmniej ostatnio używany”** (ang. least-recently used -LRU).
 - Należy zastąpić ten blok w sekcji, który pozostawał w pamięci podręcznej najdłużej bez odwoływania się do niego.
- **Algorytm „pierwszy wchodzi - pierwszy wychodzi”** (ang. first in –first out- FIFO).
 - Polega na zastępowaniu tego bloku w sekcji, który najdłużej pozostawał w pamięci podręcznej.
- **Algorytm „najrzadziej używany”** (ang. least frequently used-LFU).
 - Zastępowany jest blok w sekcji, którego miał najmniej odniesień.
- Metodą odmienną od ilości odniesień jest **losowy wybór wśród kandydujących wierszy**.

TOPOLOGIE PAMIĘCI CACHE

Topologie pamięci cache

- Look- Aside (dostęp bezpośredni)
 - Procesor odwołuje się do cache wykorzystując magistralę pamięciową.
- Look - Throgh (dostęp „przez”)
 - Układ pamięci podręcznej pośredniczy w dostępie procesora do RAM.
- Backside (dostęp „z tyłu”)
 - Układ pamięci podręcznej jest dołączony do procesora przez oddzielną magistralę nazywaną BSB (Back Side Bus).

Look - Aside (dostęp bezpośredni)



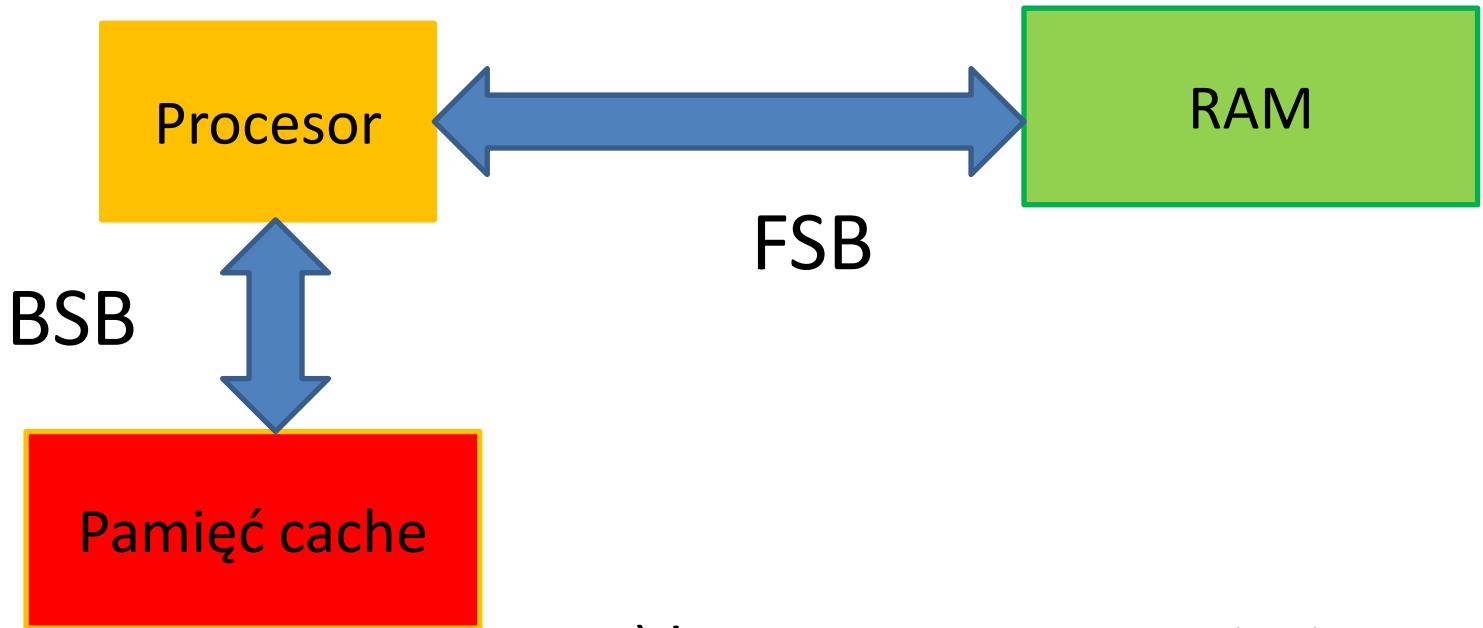
- Procesor odwołuje się do cache wykorzystując magistralę pamięciową.
- Pamięć podręczna jest podłączona równolegle z pamięcią operacyjną RAM.
 - częstotliwość pracy obu pamięci jest taka sama (komunikacja odbywa się po wspólnej magistrali),
 - tylko czas dostępu dzięki szybkości cache może ulec skróceniu.
- Wykorzystanie tej samej magistrali nie jest korzystne. Jest blokowana przy każdym dostępie procesora do cache i nie może być w tym samym czasie udostępniona innym urządzeniom.

Look - Throgh (dostęp „przez”)



- Układ pamięci podręcznej pośredniczy w dostępie procesora do RAM.
- Procesor odwołuje się do układu cache, natomiast ten układ jest dołączony przez magistralę pamięciową do RAM.

Backside (dostęp „z tyłu”)



- Magistrala FSB (Front Side Bus) łączy procesor z pamięcią RAM.
- Pamięć podręczna jest dołączona do procesora przez oddzielną magistralę BSB (Back Side Bus).
- Częstotliwości obu magistral są niezależne od siebie.
- Możliwe jest wykorzystanie magistrali FSB przez inne urządzenia zapisujące do pamięci RAM, w czasie gdy procesor komunikuje się z pamięcią podręczną po BSB.