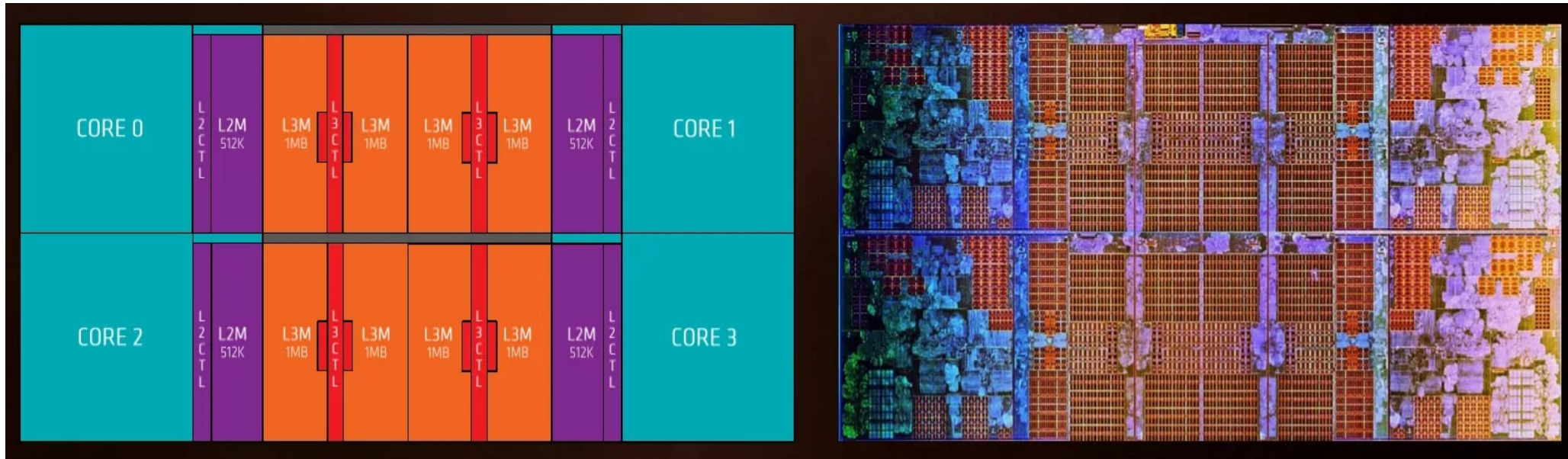


Pamięć

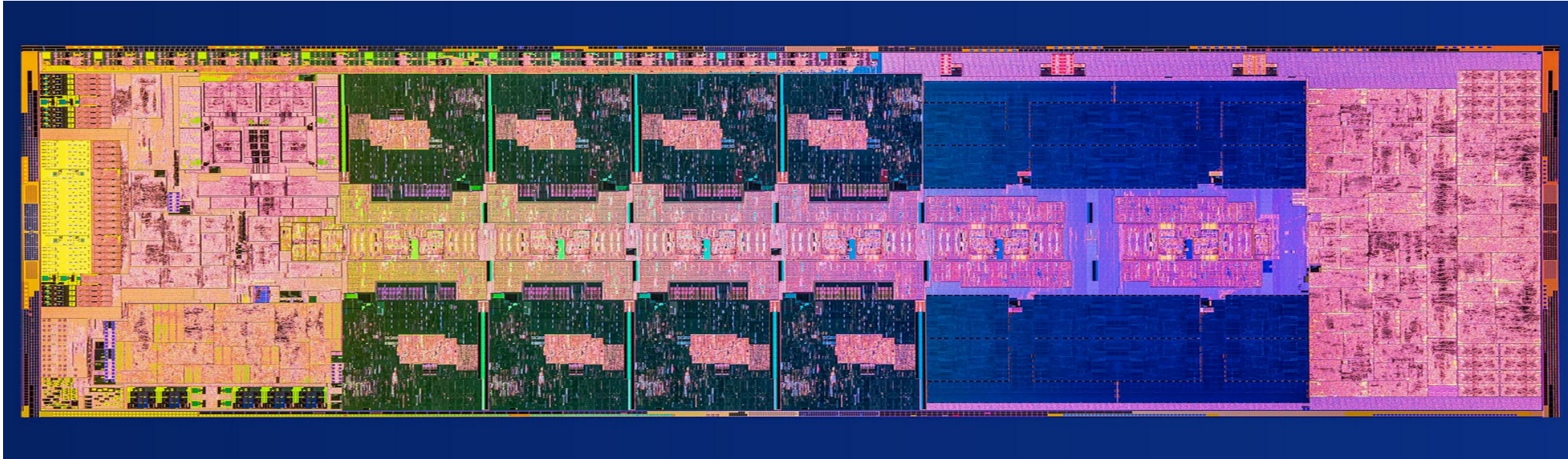
elementarne prawdy
podstawowe komponenty

AMD Threadreaper 1950X / ZEN

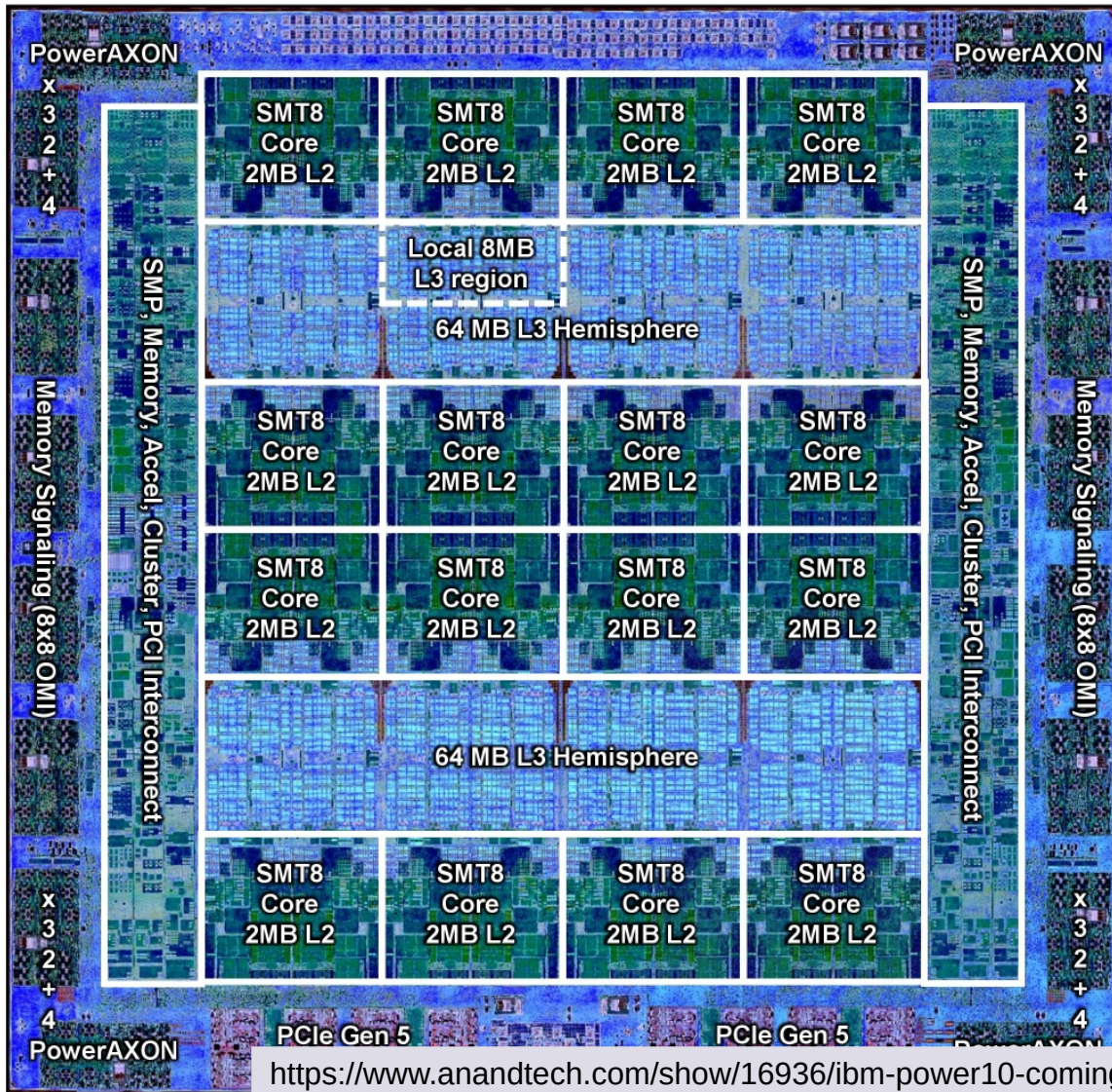


- 4 rdzenie
- L2 – 512kB/rdzeń
- L3 – 8MB – dzielone – dostępne dla innych CCX (core complex)

Intel 14 gen Raptor Lake Refresh



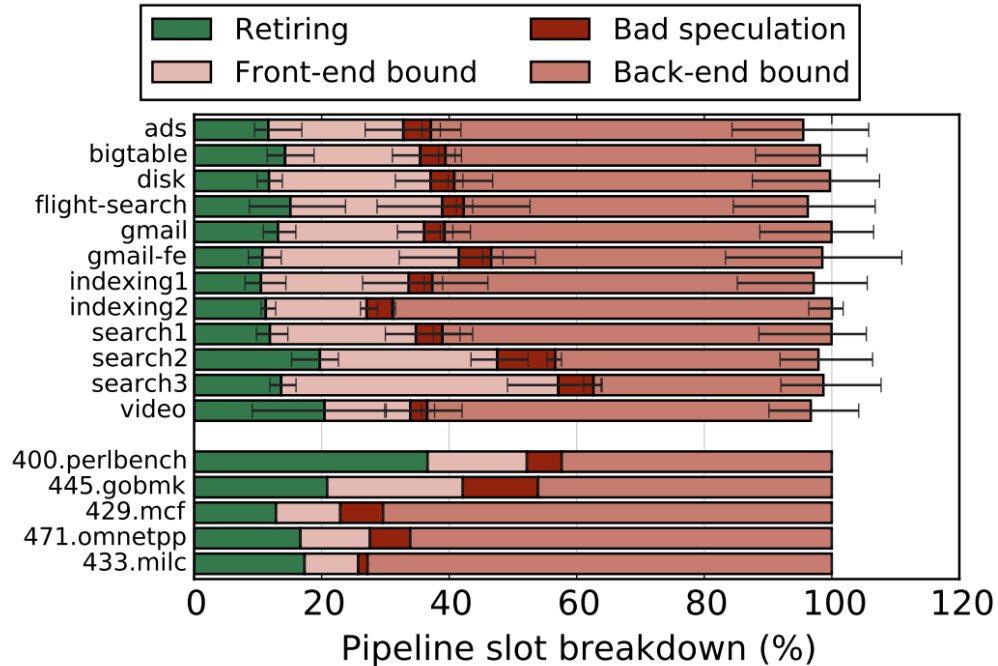
- i7-14701 – 8 rdzeni
- L2 – 2MB/rdzeń
- L3 – 33MB – dzielone



IBM Power10

- 8 rdzeni
- L2: 2MB/rdzeń
- L3: do 120MB
- przepustowość: 1 TB/s

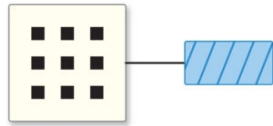
Pamięć jako zatykacz



- Retiring – OK
- Bad speculation
 - predykcja skoków
- Front-end bound
 - pusty potok
 - instruction fetch, decode, ...
- Back-end bound
 - potok niegotowy na kolejną instrukcję
 - data hazard, ILP

Pamięć jako przyspieszacz

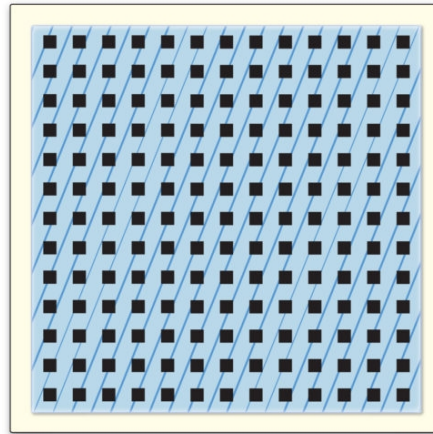
Traditional Memory Architecture



Memory separate from cores

■ Core ■ Memory

Cerebras Memory Architecture



Memory uniformly distributed across cores

■ Core ■ Memory

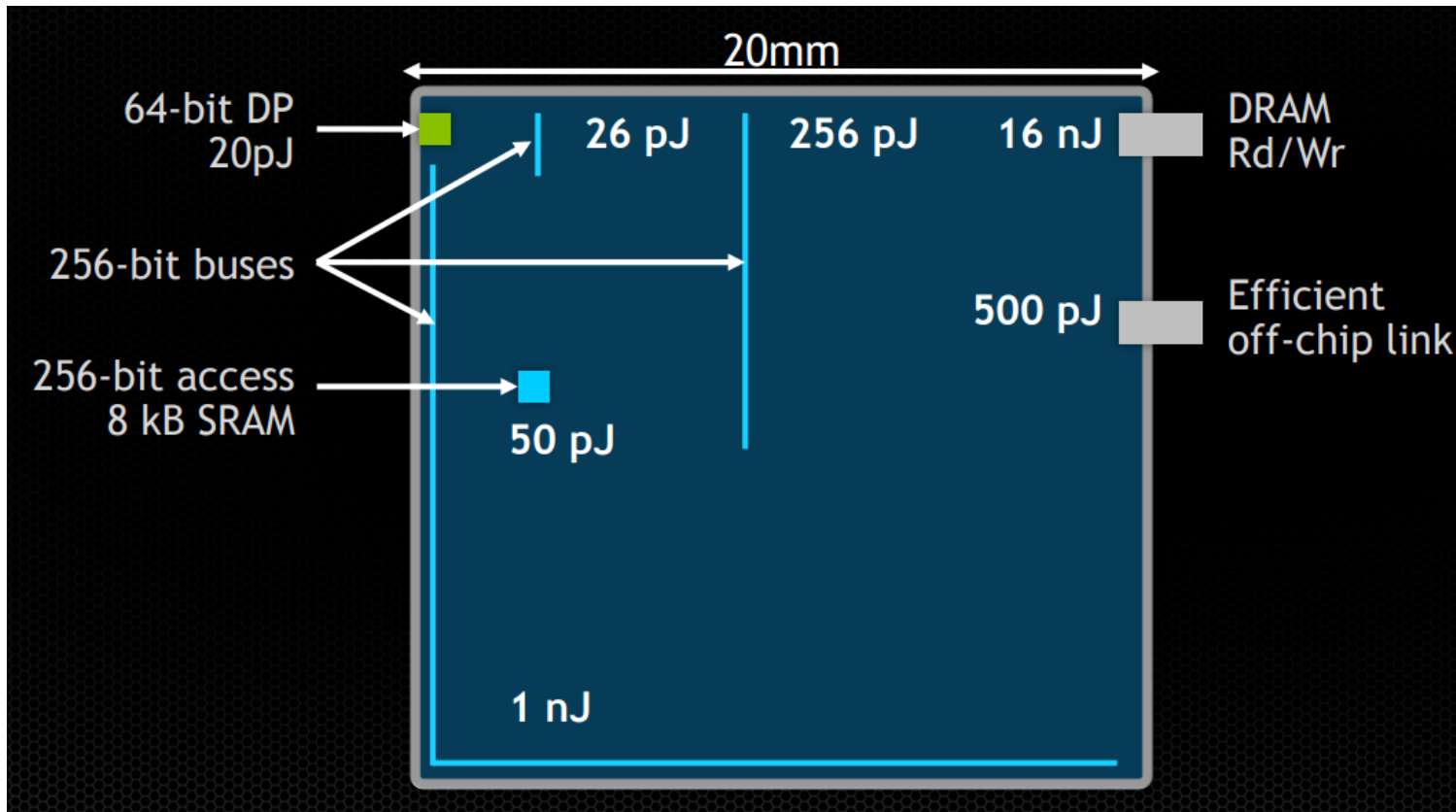
- 18 GB pamięci
- przepustowość: 9,6 PB
- powierzchnia: 46,2 mm²
- 1,2 tryliona tranzystorów
- 400 000 rdzeni dla SI

Pamięć jako pożeracz energii

ENERGY TABLE FOR 45NM CMOS PROCESS [9]. DRAM ACCESS USES THREE ORDERS OF MAGNITUDE MORE ENERGY THAN SIMPLE ARITHMETIC AND 128X MORE THAN SRAM.

Operation	Energy [pJ]	Relative Cost
32 bit int ADD	0.1	1
32 bit float ADD	0.9	9
32 bit int MULT	3.1	31
32 bit float MULT	3.7	37
32 bit 32KB SRAM	5	50
32 bit DRAM	640	6400

Pamięć jako pożeracz energii



Pamięć jako problem

- Sterowanie zapisem/odczytem w potoku
- Utrzymywanie kolejności zapisu w OoOX
- Dla poleceń SIMD: zwielokrotnione banki pam.
- Transfery GPU ↔ CPU
- ...
- Ale jest nieodzownym elementem przetwarzania

Pamięć wymagana

- Sekwencjonowanie DNA, rozwój epidemii
- Symulacje klimatyczne, spalania, przepływu
- Modele językowe, sieci neuronowe
- Astronomia, astrofizyka
- Przetwarzanie wideo, bazy danych
- Analiza (big-)danych, przetwarzanie drzew/grafów
- ...

Możliwości przechowywania

- Przerzutniki/zatrzaski → rejestry
 - kilkanaście tranzystorów, szybkie!
- Statyczna pamięć RAM
 - kilka tranzystorów, dość szybkie
- Dynamiczna pamięć RAM
 - tranzystor+kondensator, wolne, upływne, produkcja...
- Dyski, taśmy, płyty, pen-drive
 - niski koszt na MB, bardzo wolne

Mechanika

- Macierz do przechowywania danych
- Każdy wiersz adresowalny
- Każda kolumna adresowalna
- Układ do odczytu/interpretacji danych
- Układ do zapisu danych



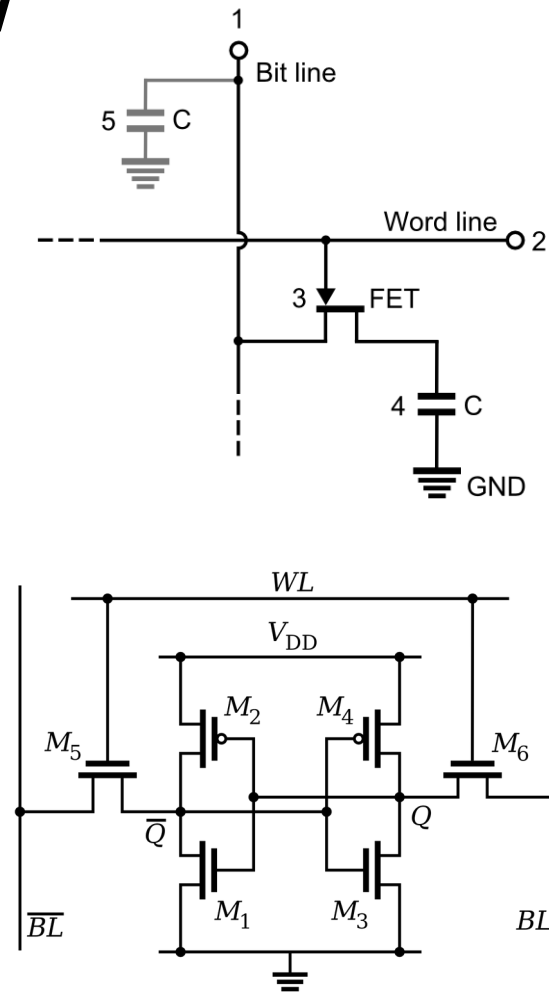
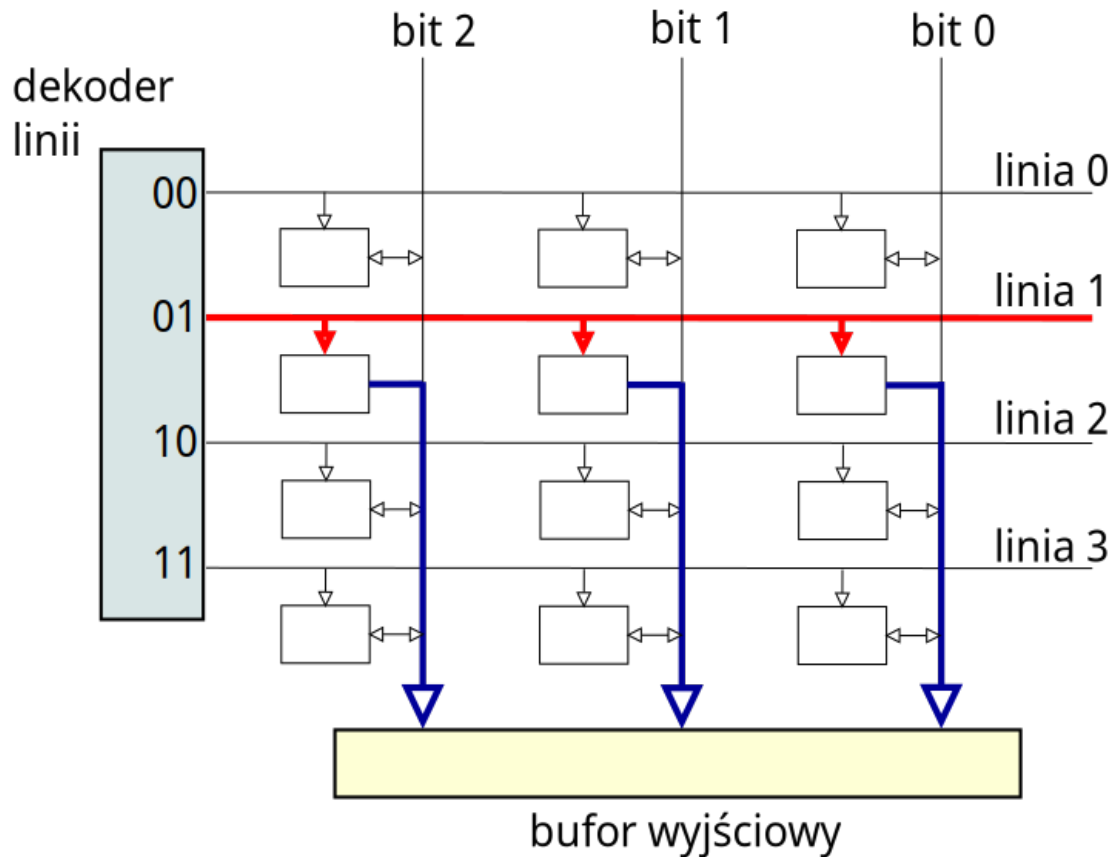
Adres	Dane		
11	1	1	0
10	1	1	0
01	0	1	0
00	1	1	1

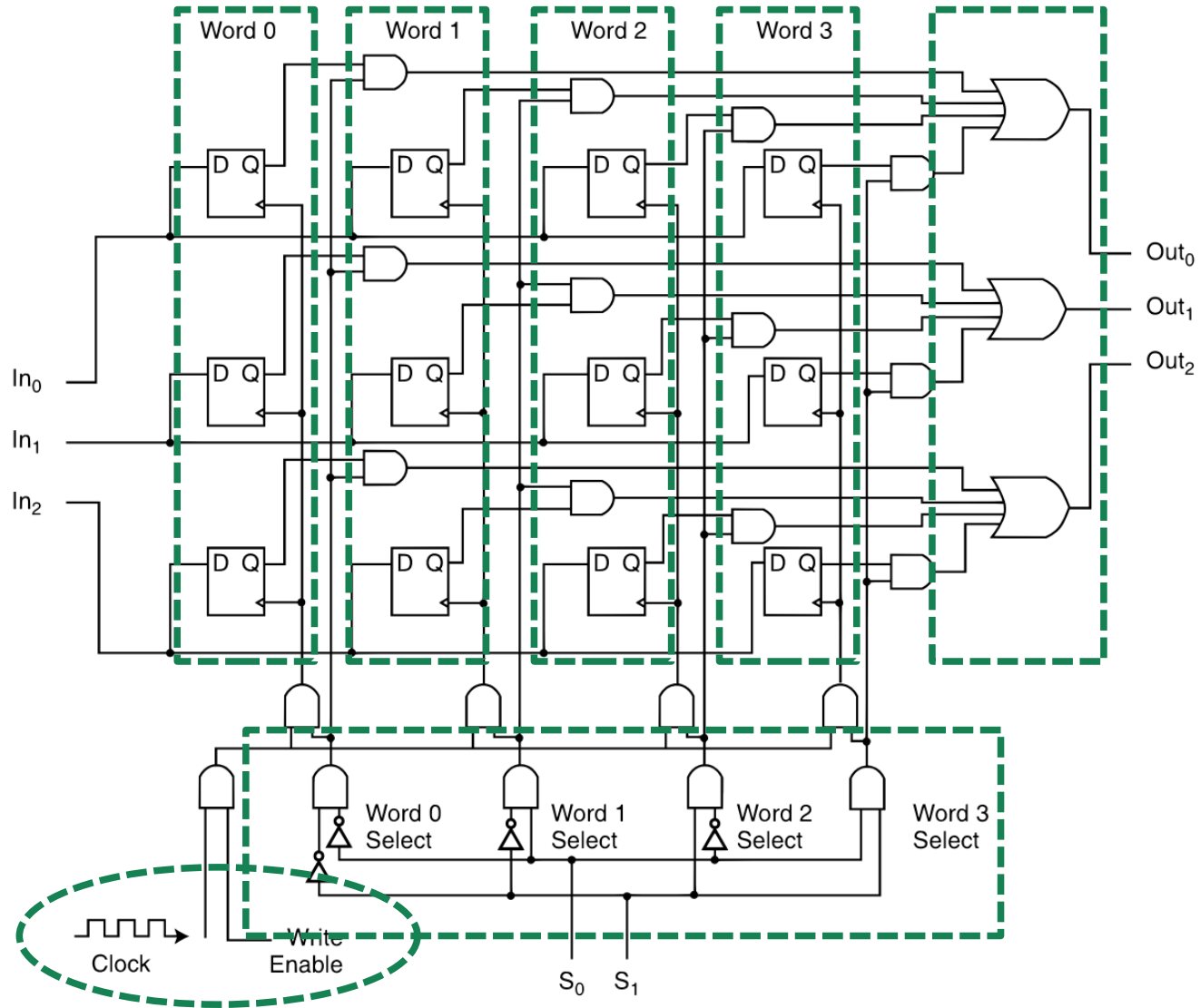
↑
głębokość

←→
szerokość (słowa)

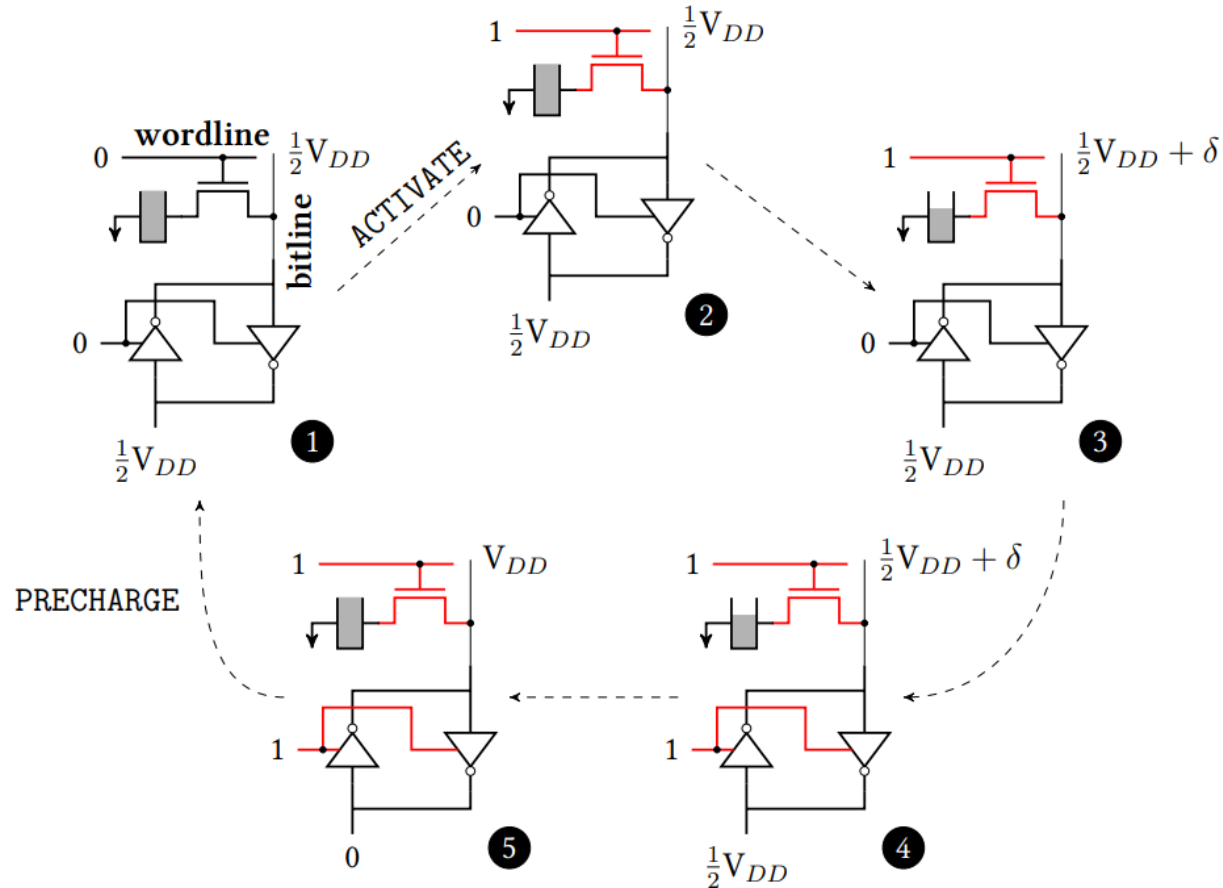
Adres	Dane						
11111111	1	1	0	1	1	0	1
11111110	0	1	0	1	1	1	0
11111101	0	0	0	0	0	1	0
...	0	0	0	0	0	1	0
...
...
...	1	1	1	0	1	1	1
...	0	1	0	0	0	0	0
0000010	0	0	1	1	0	1	1
0000001	1	0	1	1	1	0	0
0000000	1	1	0	1	0	1	1

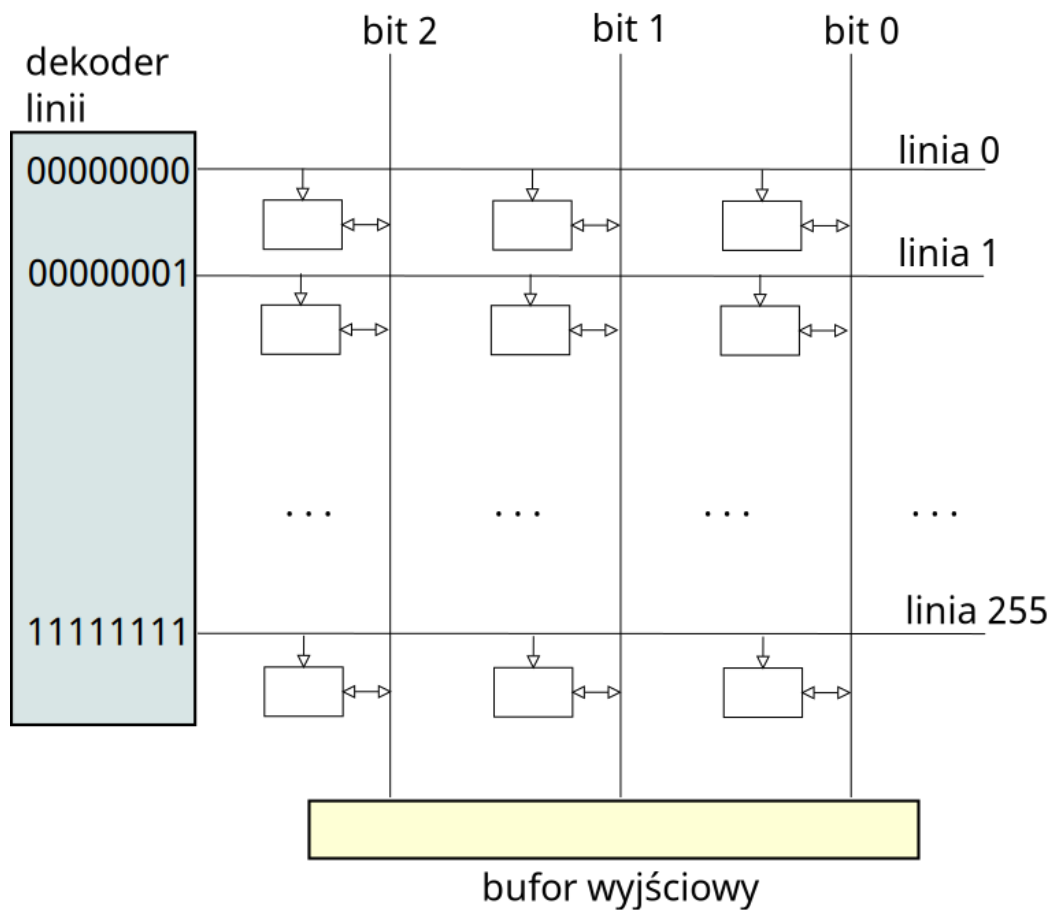
Dostęp do bitów



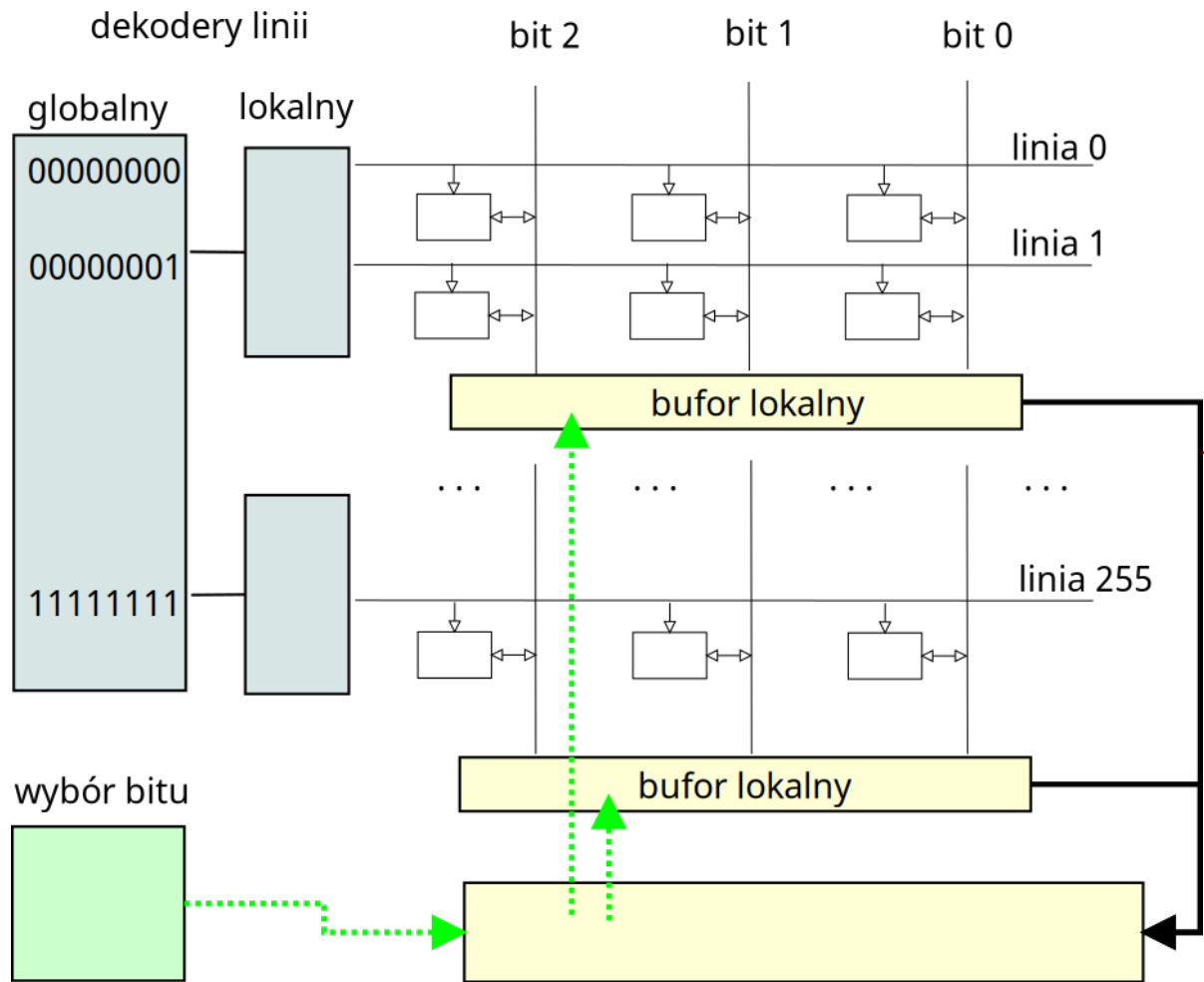


Odczyt bitu – wzmacniacze odczytu





robi się długie

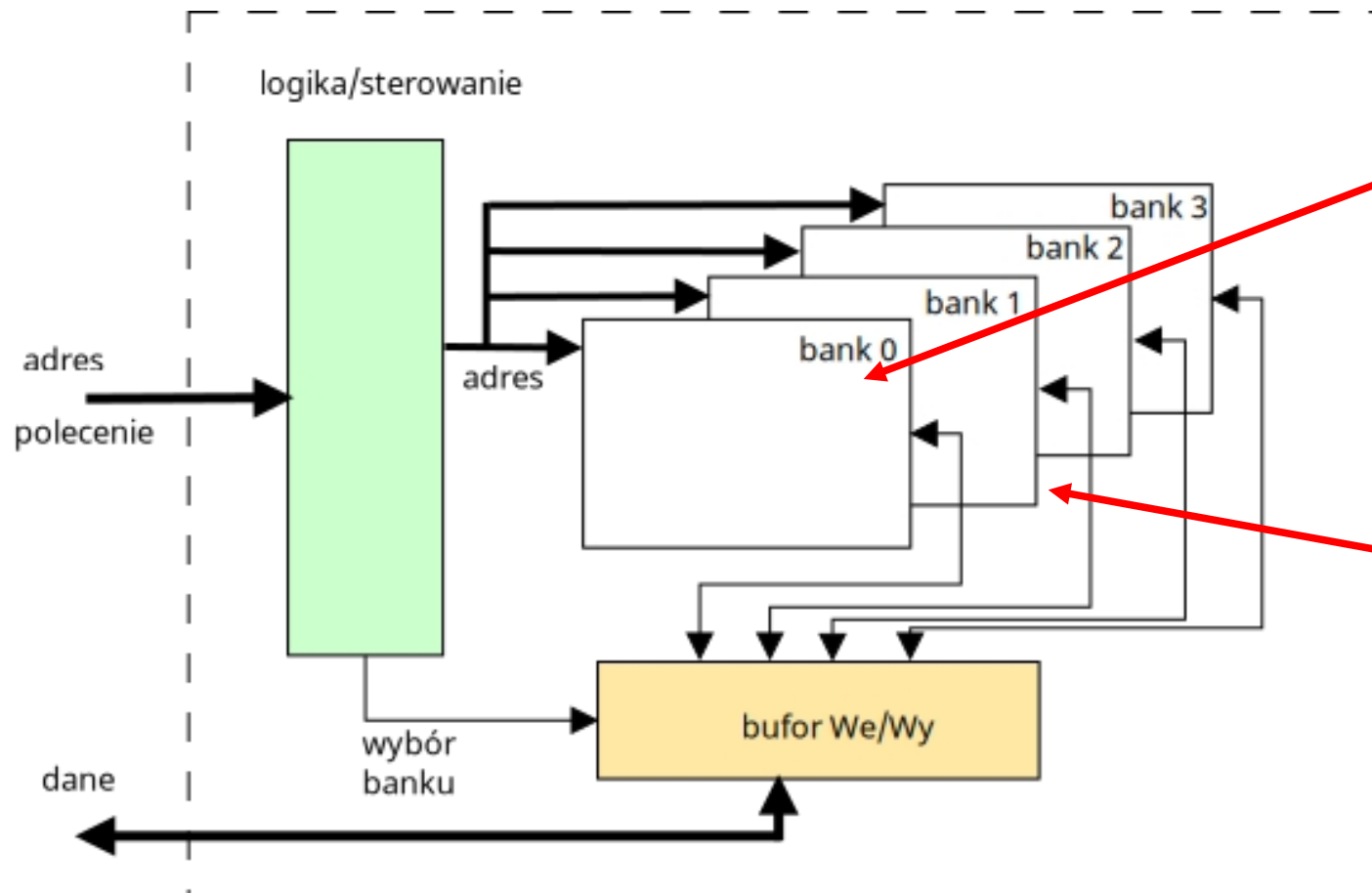


na innej warstwie grubsze

można też rozbudowywać „na szerokość”

no i mamy **bank!**

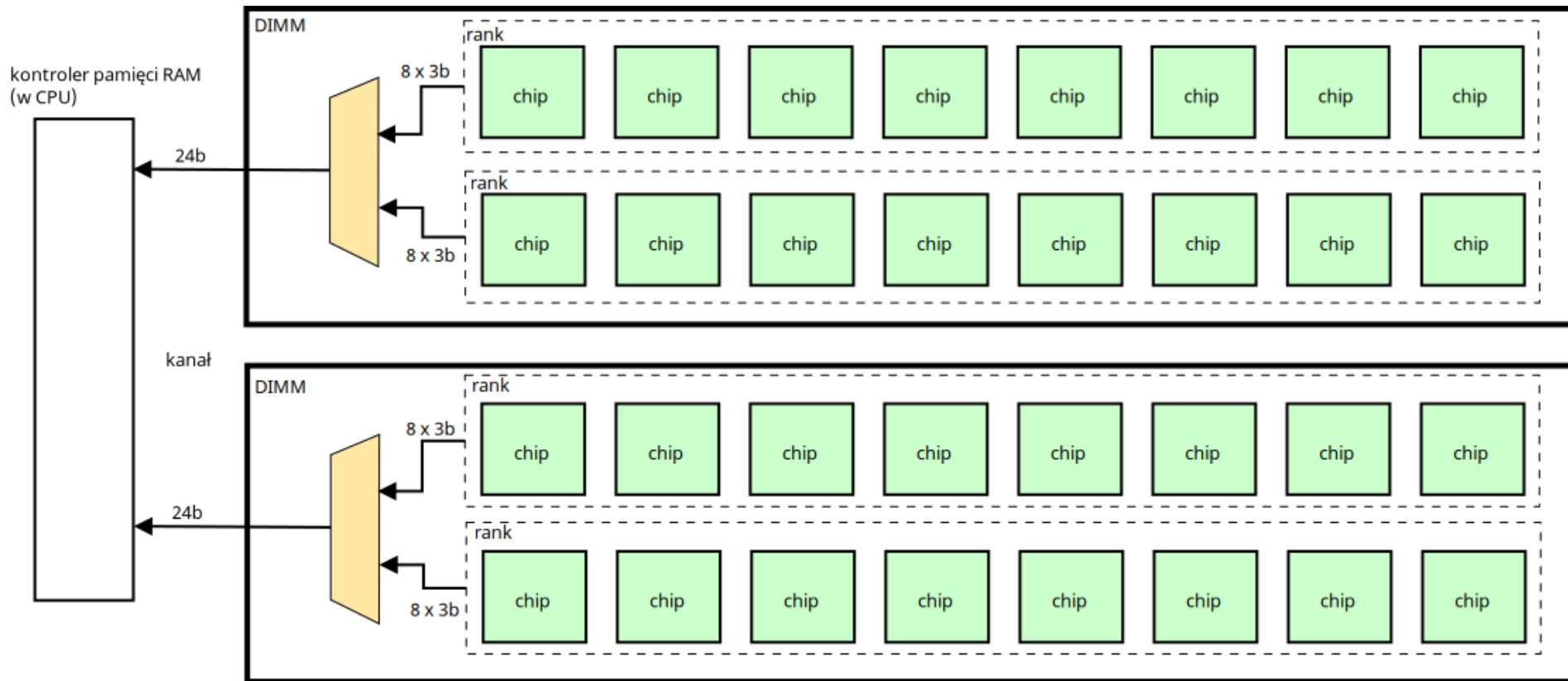
Banki składamy w chipy...

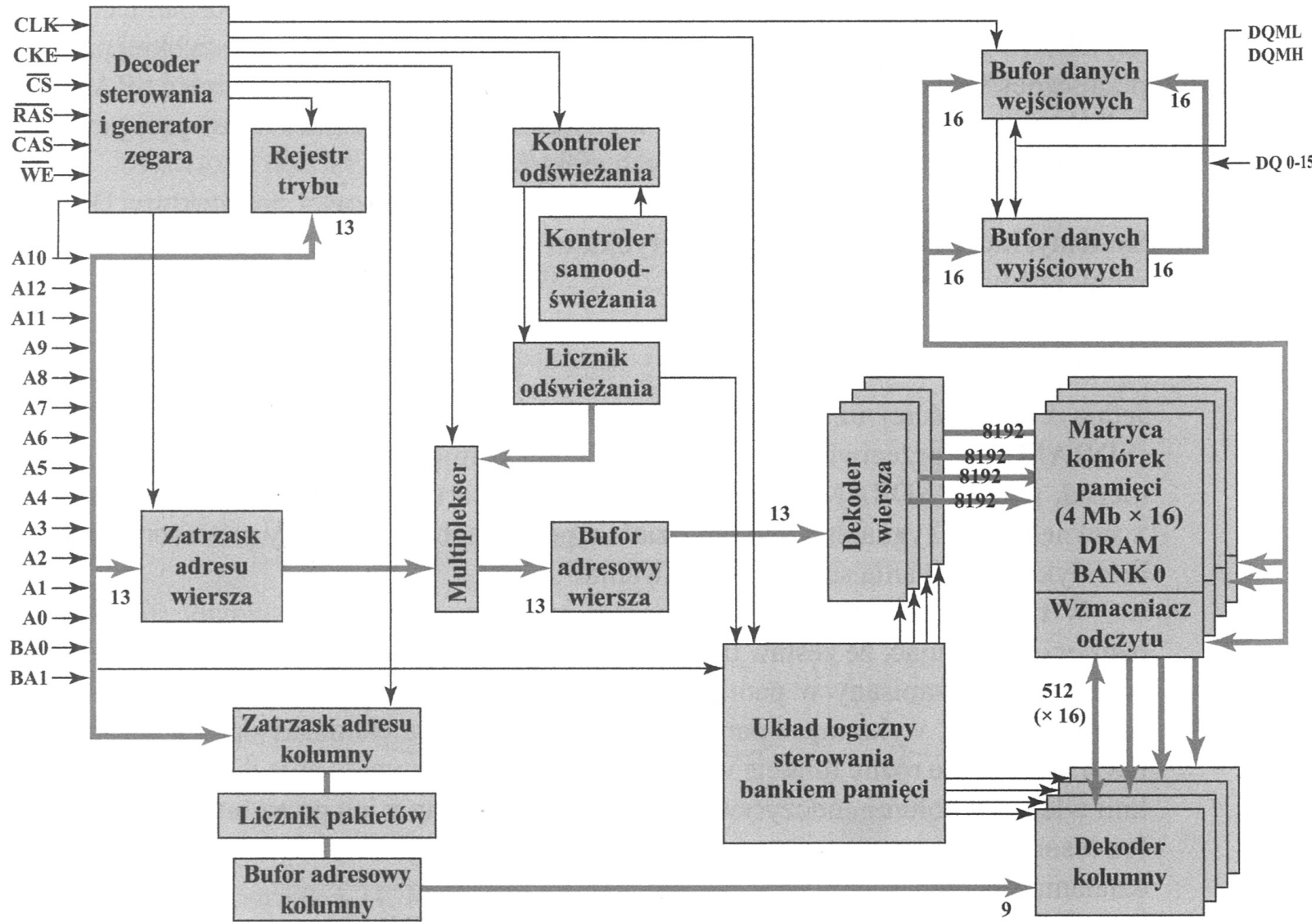


dostęp w linii
jednego banku szybki

multipleksowanie
banków to taki
„potok”

... i tak dalej



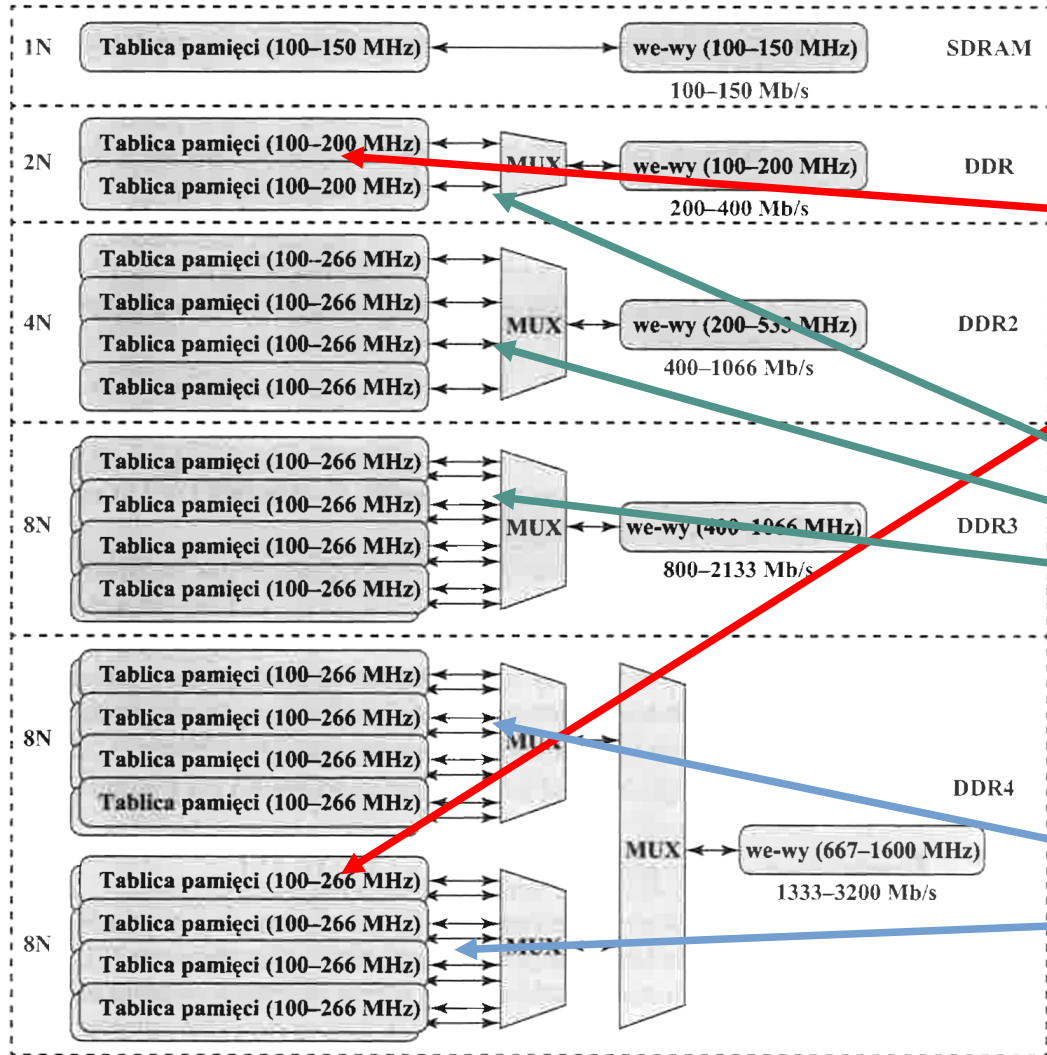


Parametry DRAM

Czas dostępu	- od zaadresowania komórki do pojawienia się danych na Wy - kilka ns: $1/\text{częstotl.}$
Czas latencji (CAS/CL)	- liczba cykli zegara między poleceniem odczytu a pojawieniem się danych - SDRAM – 3, DDR-SDRAM – 2,5 lub 2 (!)
RAS to CAS delay (RCD)	- liczba taktów między adresem wiersza a adresem kolumny
Row Precharge Time (RP)	- liczba taktów między zapisem/odczytem a przejściem do innego wiersza (przygotowanie)
Row Cycle (RC)	- czas dostępu do różnych wierszy w banku

Etapy DRAM

FPM DRAM	adresacja linii, potem poszczególnych komórek
EDO RAM	niezależny odczyt i wybór kolejnego wiersza
SDRAM	synchroniczny z głównym zegarem praca pakietowa (burst): od-do banki, niezależny dostęp „na zakładkę”
DDR	synchronizacja na obu zboczach zegara niższe napięcie (2,5 V)
DDR2	przesyłane 4 bity na takt niższe napięcie (1,8 V) 240 pinów
DDR3	8 bitów na takt technologia 90 nm, napięcie 1,5 V
DDR4	proces 32/36 nm, napięcie 1,1 – 1,2 V 288 pinów
DDR5	32 banki, większa długość pakietu (16 B), jednoczesne odświeżenie i działanie na banku ECC na pokładzie, układy 3D

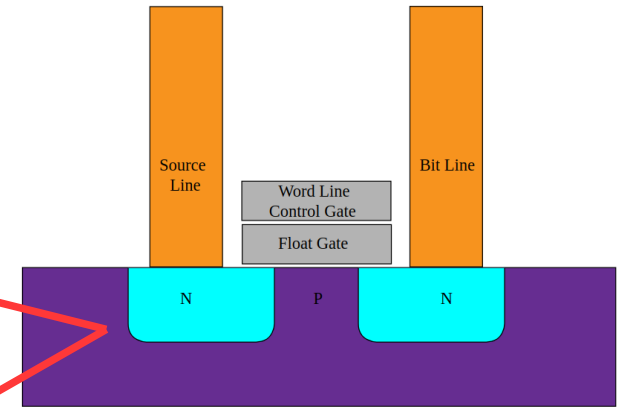
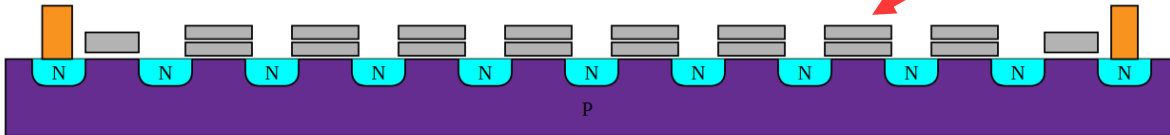
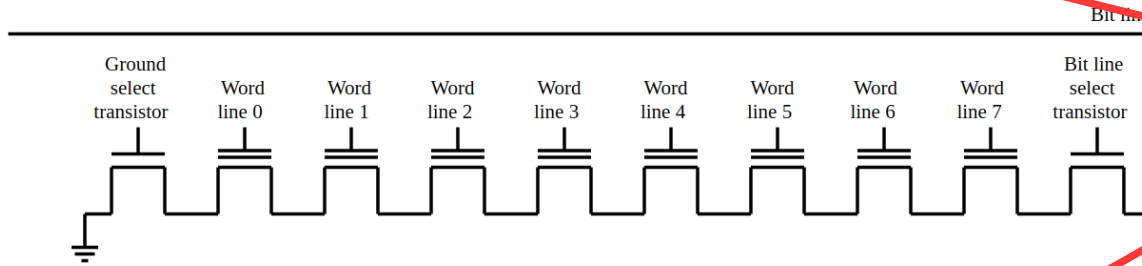
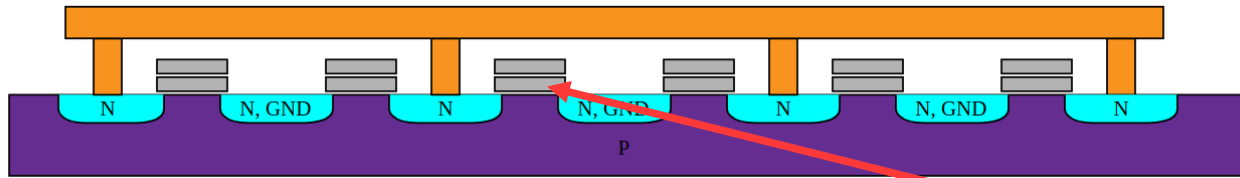
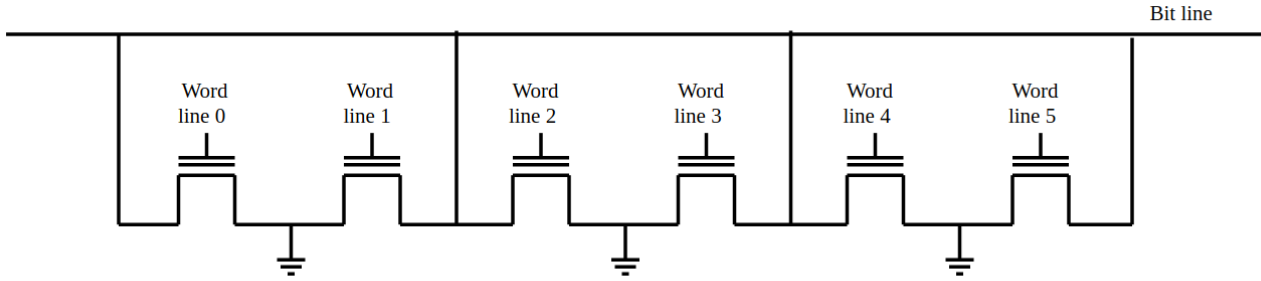


zegary rdzenia DDR
nie przyspieszają...

zwiększa się liczba
buforowanych bitów
w bankach

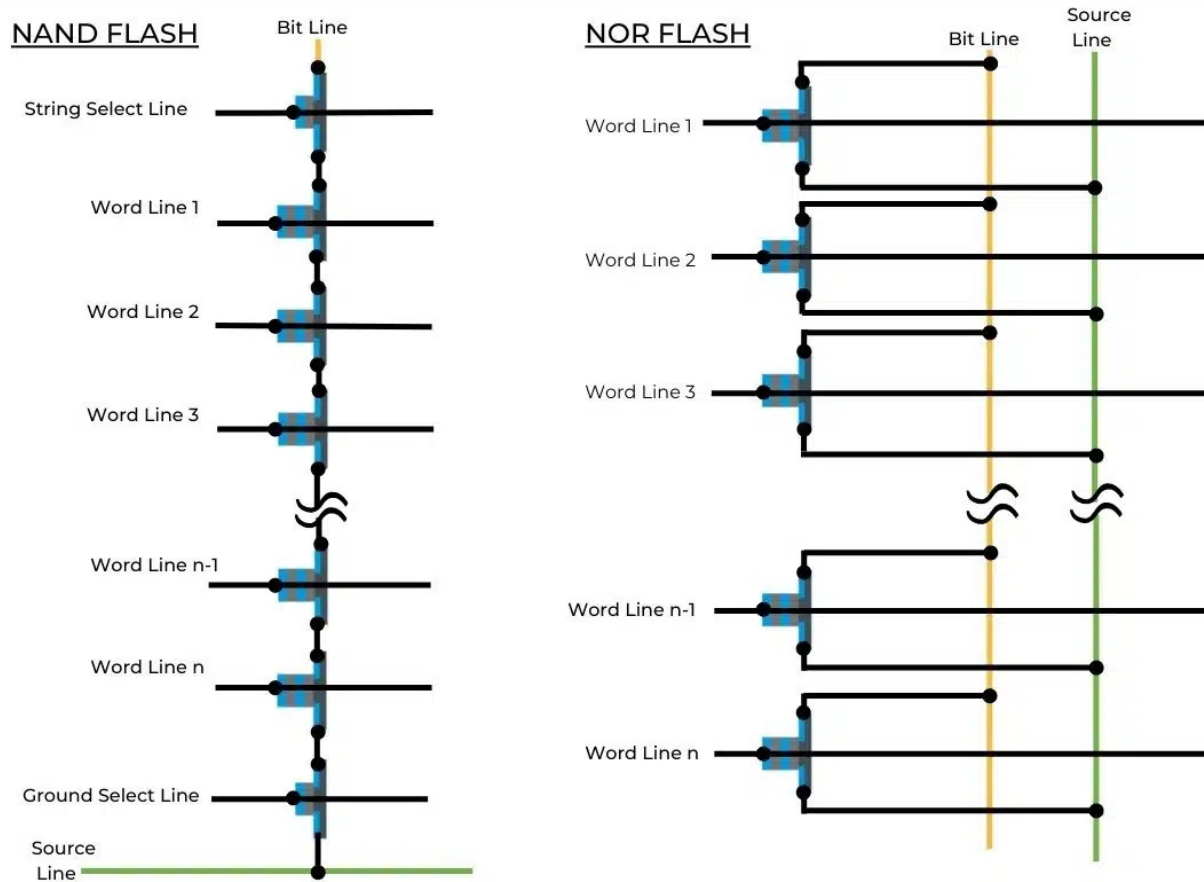
pojawiają się grupy
niezależnych banków

Flash: NOR/NAND

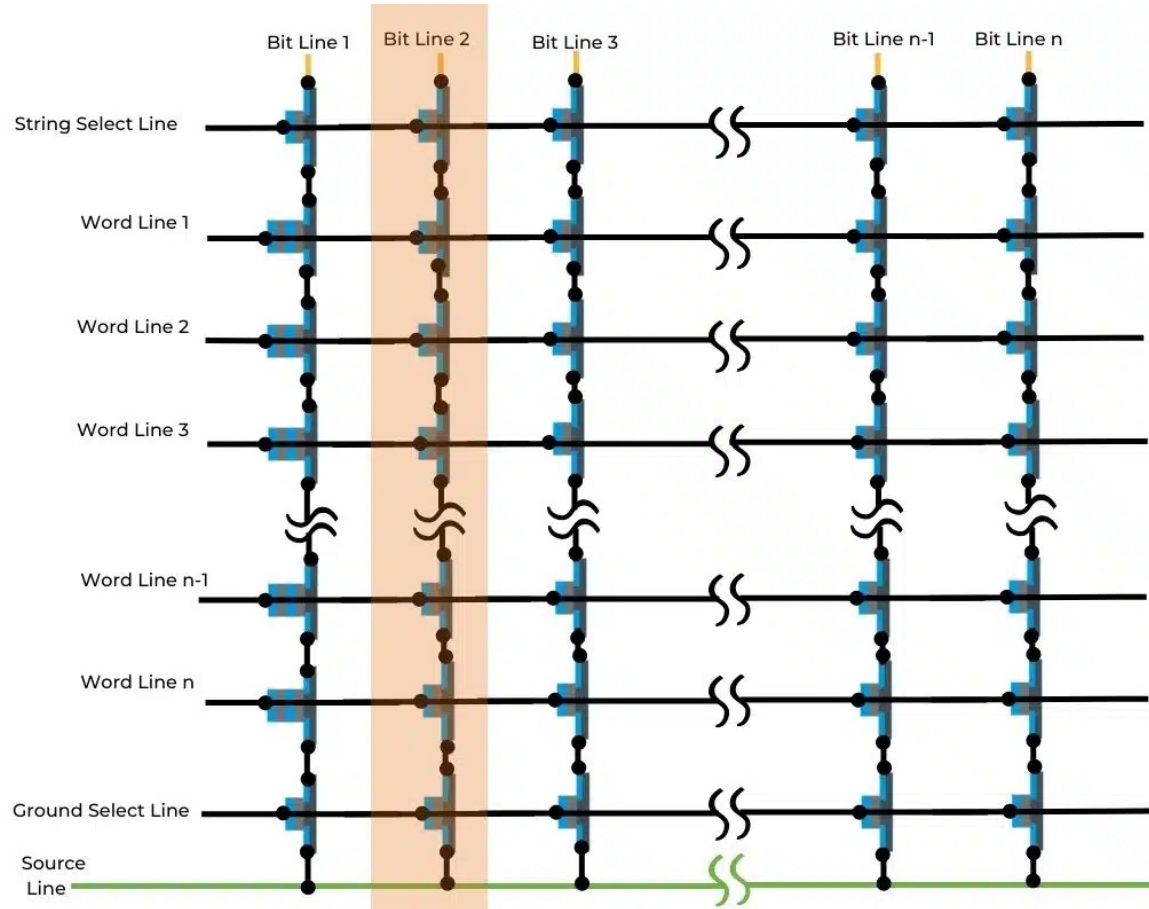


MOSFET
z pływającą bramką

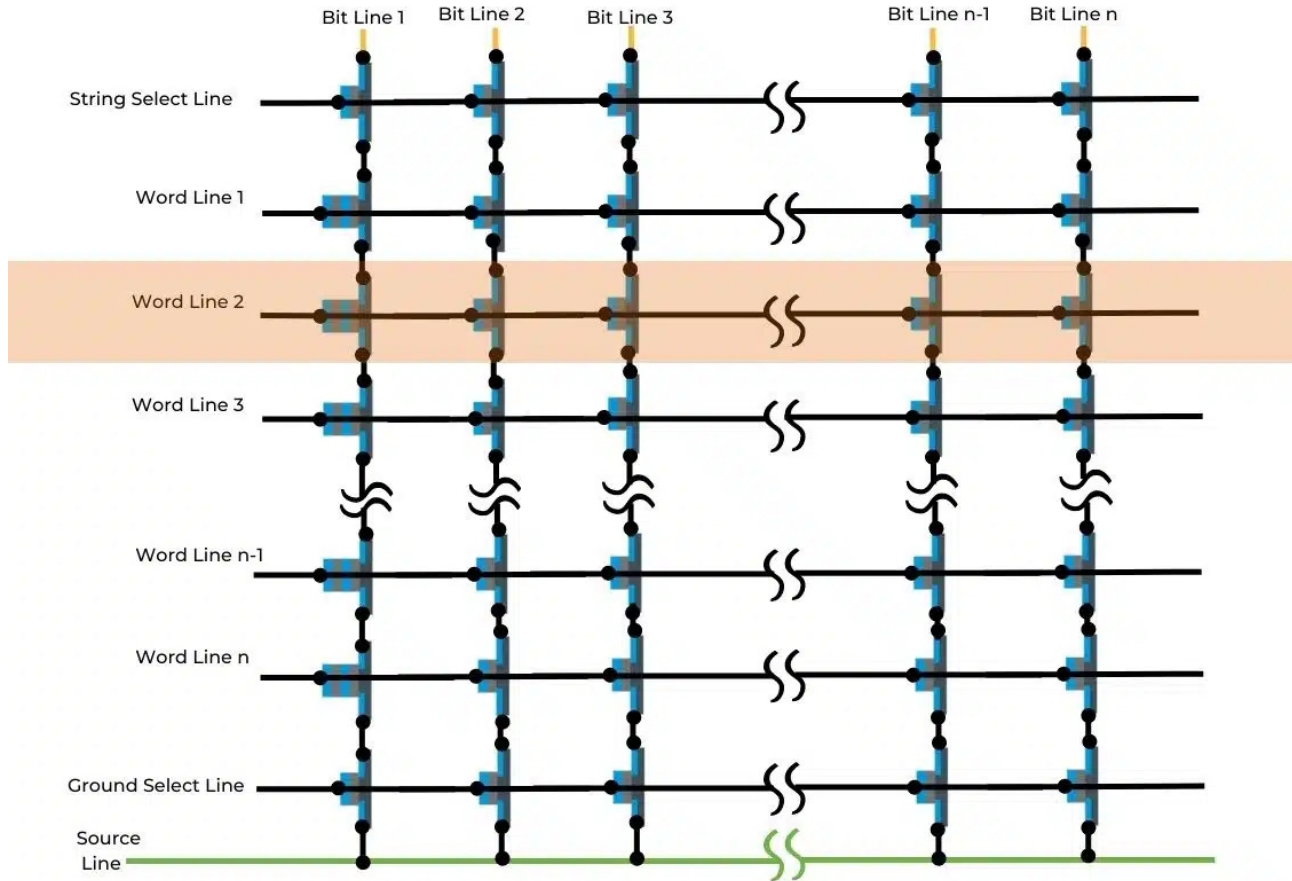
Pamięć flash



Pamięć flash



Pamięć flash



	NOR	NAND
Gęstość/upakowanie	mniej	więcej
Czas odczytu	szybciej	wolniej
Moc przy odczycie	niższa	wyższa
Czas zapisu/kasowania	wolniej	szybciej
Moc przy zapisie/kasowaniu	wyższa	niższa
Moc bezczynności	niższa	wyższa
Żywotność	większa	mniejsza
Dostęp	losowy, bity	blokowy
Koszt/MB	wyższy	niższy
Zastosowanie	pamięć uP, wbudowane	pamięć masowa, SSD