

Cache: RAM Access

jueves, 7 de abril de 2022 06:55 a. m.

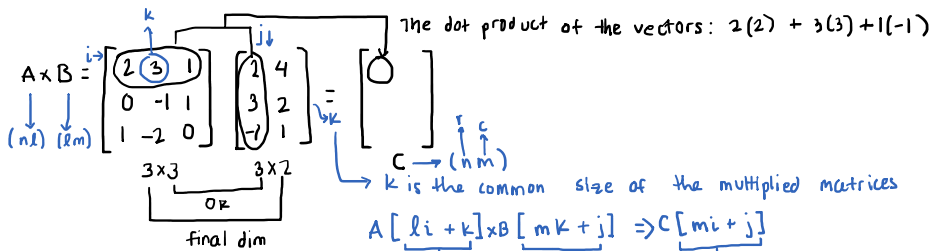
Parallel only: false sharing problem

↳ Due to locality principle, we are obliged to transfer whole blocks, and in order to have all words correct, we are constantly going to RAM: bad performance.

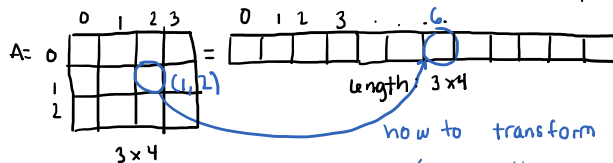
↳ the cause: cache is organized in lines

Dev C++: Run thread-using programs:

Compiler Options > Code Generation > Language standard: GNU C++11

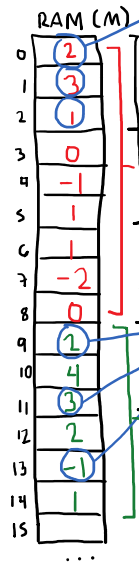


But since RAM is ordered as a line of memory spaces



we have: $(1, 2) = (x, y)$
 we want: $(6) = 4x + y$
 ↳ amount of columns

How is this stored? for this element, the whole block is sent (block 1)

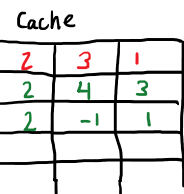


blocks (3) → thus, cache has lines of 3 spaces (words)

0 multiplied elements

block 3 and 4 are copied to cache

we would want that 2, 3, -1 is just one block: only 2 blocks in cache instead of 3
 = less copied RAM blocks to cache, less RAM access
 = better performance



how can we make B matrix elements contiguous?

Transpose B!

$B = \begin{bmatrix} 2 & 3 & -1 \\ 4 & 2 & 1 \end{bmatrix}$ → together in RAM

↳ but $A \times B \neq A B^T$

But we can change the code so that this is $A B^T = A \times B$

$$\begin{bmatrix} 2 & 3 & 1 \\ 0 & -1 & 1 \\ 1 & -2 & 0 \end{bmatrix} \begin{bmatrix} 2 & 3 & -1 \\ 4 & 2 & 1 \end{bmatrix} = \begin{bmatrix} \quad \quad \quad \\ \quad \quad \quad \\ \quad \quad \quad \end{bmatrix}$$

$A \times B^T = C$

Output

```
C:\Users\Administrator\Desktop\04072022\cache-matmul.exe
# elapsed time (Inicializacion): 0.0111964s
# elapsed time (AxB): 8.91092s
# elapsed time (Bt): 0.0275517s
# elapsed time (AxBt): 5.98642s
# elapsed time (Bt_AxBt): 6.01545s

-----
Process exited after 15 seconds with return value 0
Press any key to continue . . .
```

Even though this is a serial program, we see that there is a slower performance due to cache access. The parallel version would have serial parts, and thus carry this problem.