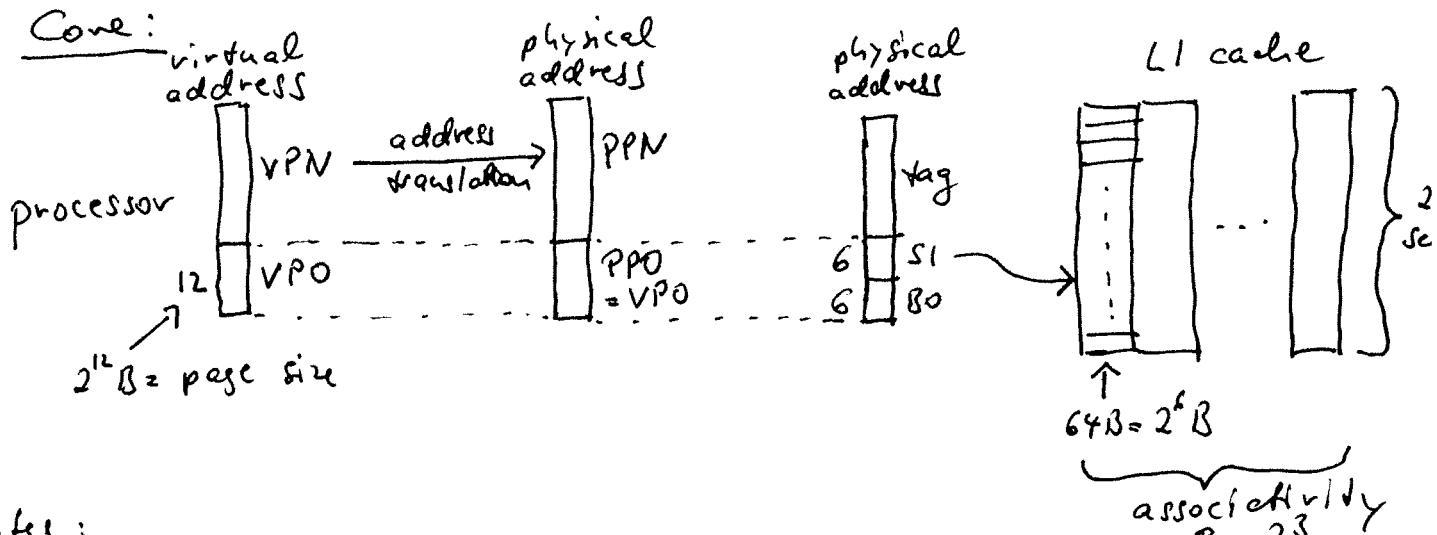


Optimizations Related to the Virtual Memory System

Background: (6re)

- the processor works with virtual addresses
- all caches work with physical addresses
- both address spaces are organized in pages
- typical page size: 4KB
- so the address translation translates virtual page numbers into physical page numbers

Cave:



Notes:

VPN = virtual page number

VPO = " " offset

PPN = physical page number

PPO = " " offset

S1 = set index

B0 = block offset

address translation: $\text{VPN} \rightarrow \text{PPN}$

$\text{VPO} = \text{PPO} = \text{S1} \cup \text{B0} \Rightarrow$ cache lookup can start before $\text{VPN} \rightarrow \text{PPN}$ translation is finished

address translation

- uses a cache called translation lookaside buffer (TLB)

- Case 2: two levels of caches for loads

DTLB0: 16 entries

DTLB1: 256 entries

Case 1: DTLB0 hit: no penalty

DTLB1 hit: 2 cycle penalty

miss: possibly very expensive

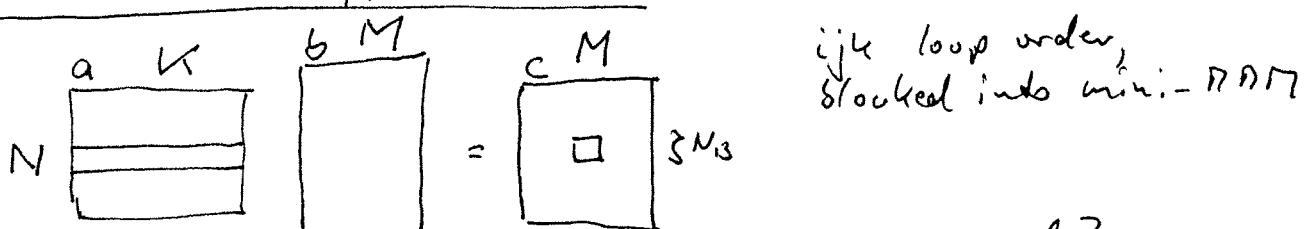
Pentium 4:
- one TLB
- 64 entries

Consequence: Repeatedly accessing a working set that is spread over > 256 pages leads to TLB misses \rightarrow possible severe slowdown

Solution 1: use larger pages
may require different kernel (OS) and C std library

Solution 2 (if possible): copy working set into contiguous memory
 \Rightarrow less pages are used

How does this affect MM2?



which memory regions are repeatedly accessed?

- block rows of a : is contiguous

- all of b : is contiguous

- file of c : can be spread over N_B pages
if $M > 512$ (512 doubles = $4 \times 13 =$ page size)

But: typically $N_B < 100 < \text{size (DTLB)}$

so at most 2 cycles penalty per row

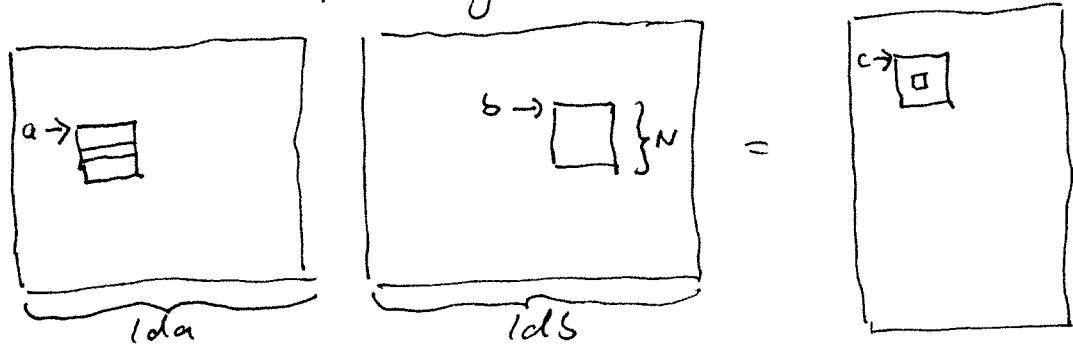
\Rightarrow not worth to copy (on Core)

But: the BLAS 3 function `dgemm` has this interface:

`dgemm(a, b, c, N, K, M, lda, ldb, ldc)`

$\underbrace{a, b, c}_{\text{matrix pointers}}$ $\underbrace{N, K, M}_{\text{matrix dimensions}}$ $\underbrace{\overbrace{\text{leading}}^{\text{"leading"}}$ $\underbrace{\text{dimensions}}_{\text{dimensions}}$

The leading dimensions enable `dgemm` to be called on submatrices of larger matrices:



which memory regions are repeatedly accessed? ldc

- block rows of a : spread over $\leq N_B$ pages

- all of b : spread over $\leq N$ pages

- file of c : spread over $\leq N_B$ pages

here copy may pay for large enough N and $lda > N$

Code :

```
// all of B reused : possibly copy  
for i = 0 : N_B : N - 1  
    // block row of A reused : possibly copy  
    for j = 0 : N_A : M - 1  
        // tile of C reused : possibly copy  
        for k = 0 : N_B : K - 1  
            . . . . .
```