# Optimizations Related to the Virtual Memory System

## Background: (Gre)
- 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

Core:



## Notes:
VPN = virtual page number
VPO = "      "   offset
PPN = physical page number
PPO = "      "   offset
SI  = set index
BO  = block offset

address translation: VPN $\rightarrow$ PPN
VPO = PPO = SI $\cup$ BO $\implies$ cache lookup can start before VPN $\rightarrow$ PPN translation is finished !

## address translation
- uses a cache called translation lookaside buffer (TLB)
- Gre 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 MMM?



ijk loop order,
blocked into mini-MMM

which memory regions are repeatedly accessed?
  - block rows of a: is contiguous
  - all of b: is contiguous
  - tile of c: can be spread over $N_3$ pages
        if $M > 512$ 💬 ($512$ doubles = $4 KB$ = page size

Sut: typically $N_3 < 100 <$ size (DTLB1)
     so at most 2 cycles penalty per row
     $\Rightarrow$ not worth to copy (on Core)

<u>But</u>: the BLAS 3 function dgemm has this interface:
       dgemm (a, b, c, N, K, M, lda, ldb, ldc)
               $\underbrace{\qquad}$   $\underbrace{\qquad}$   $\underbrace{\qquad}$
               matrix      matrix      ~~leading~~ "leading"
               pointers    dimensions   dimensions

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



which memory regions are repeatedly accessed?
  - block row of a: spread over $\leq N_3$ pages
  - all of b: spread over $\leq K$ pages $\leftarrow$ 💬
  - tile of c: spread over $\leq N_3$ pages

Here copy may pay for large
enough K

Code:

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