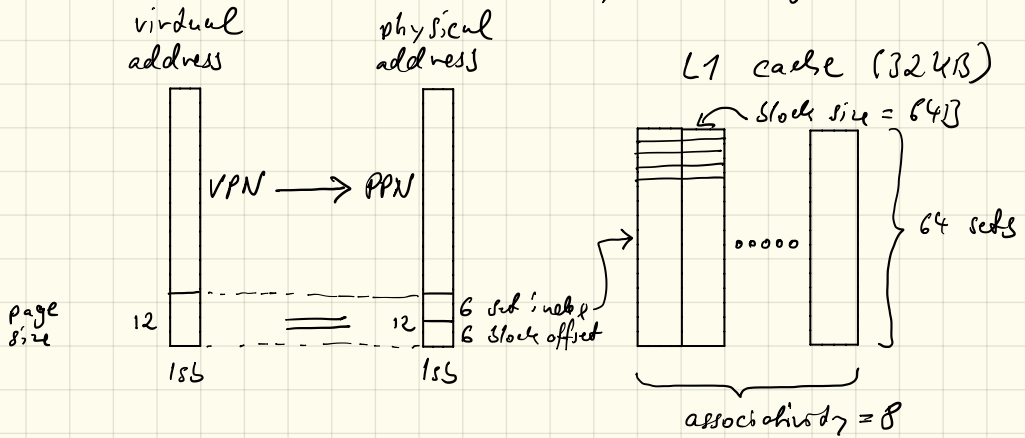


Virtual Memory System (Core Family)

- the processor works with virtual addresses
- all caches work with physical addresses
- both address spaces are organized in pages

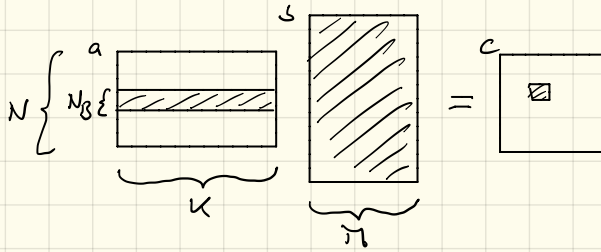
page size: 4 KiB (can be changed to 2 MiB and even 1 GiB on the latest processor; change in OS settings)

- address translation: virtual page number (VPN) \rightarrow physical page number (PPN)



\Rightarrow L1 cache lookup can start concurrently with address translation

Example MMM



working set at highest level is shaded

- We look for parts in working set spread in memory
- block rows of a : contiguous
 - all of s : contiguous
 - block of c : if $m > 512$ (512 doubles = 4KB) then spread over $\geq N_B$ pages

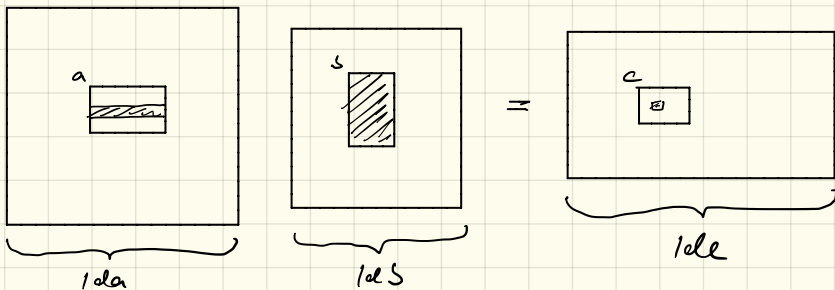
typically N_B is in the 10s, so no big problem

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

$$\text{dgemm}(a, s, c, N, k, m, lda, lds, ldc)$$

matrices
sizes
leading dimensions

leading dimensions: enable use on matrices inside matrices



assume $l_{da}, l_{db}, l_{dc} > 512$

- block row of a : spread over $\geq N_B$ pages
- all of b : spread over $\geq K$ pages
- block of c : spread over $\geq N_B$ pages

So copying to contiguous memory may pay off!

// all of b reused: possibly copy

for $i = 0 : N_B : N-1$

// block row of a reused: possibly copy

for $j = 0 : N_B : M-1$

// block of c reused: possibly copy

for $k = 0 : N_B : K-1$

•••••