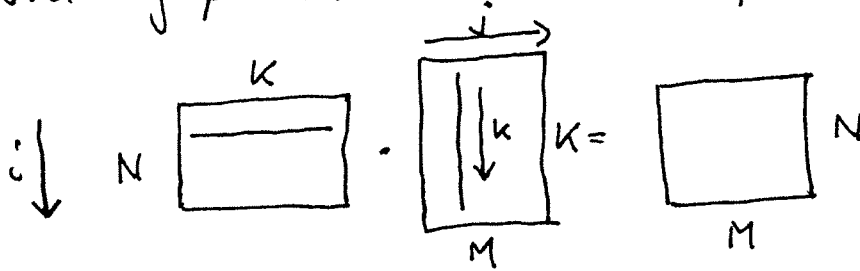


# MM generation with ATLAS

Starting point: standard triple loop



```
for i = 0:1:N-1
  for j = 0:1:M-1
    for k = 0:1:K-1
      cij = cij + aik bkj
```

$$C = C + AB$$

Important cases (from most to least):

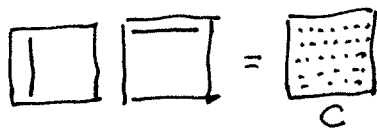
- ~~two~~ two out of  $(N, K, M)$  are small
- one out of  $(N, K, M)$  is small
- none out of  $(N, K, M)$  is small

Reason: based on how it is used inside LAPACK

- 1.) loop order:  $i, j, k$  loops can be permuted into any order
- $ijk$ : B is reused, good if  $M < N$
  - $jik$ : A " " "  $N < M$

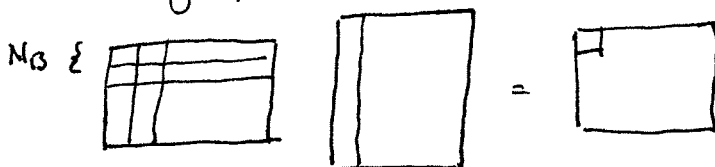
ATLAS generates versions for both

- other choices are bad, e.g.,  $kij$ :



poor temporal locality  
w.r.t. C

2.) blocking for cache



we assume for simplicity  
 $N_3 | (N, M, K)$

```
for i = 0: N3: N-1
  for j = 0: N3: M-1
    for k = 0: N3: K-1
      for i' = 0:1: i+ N3-1
        for j' = j:1: j+ N3-1
          for k' = k:1: k+ N3-1
            ci'j' = ci'j' + ai'k' bk'j'
```

mini-MM

This blocking is done using loop tiling + loop exchange  
 loop tiling: for  $i = 0:1:N-1$   
 (example)  $a_i = \dots$

$\Rightarrow$  for  $i = 0:4:N-1$   
 for  $j = i:1:i+3$   
 $a_j = \dots$

$N_B$  becomes a search parameter.

Bound  $N_B^2 \leq C$  (cache size)

this is a loose bound — the search takes care of the rest.

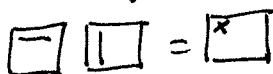
For which cache should one block?

- The original ATLAS blocked for L1-cache
  - L2 was off-die
  - the important cases (see above) of MMM have one or two of  $N, K, M$  small
- For large, square matrices blocking for L2 can make sense
- Blocking for L1 & L2 can make sense
- Versions with different block sizes can make sense

3.) Blocking Mini-MMM for registers

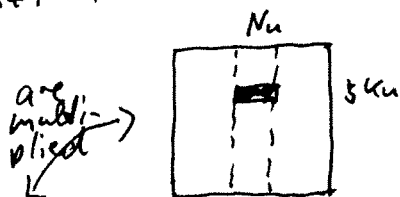
Now:  $k$  as outermost for ILP:

$ijk$

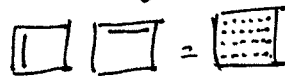


$2n$  instructions:  
 -  $n$  independent multi  
 -  $n$  dependent adds  
 ( $\geq \log_2(n)$  steps)

But:  $2n+1$  "live" variables

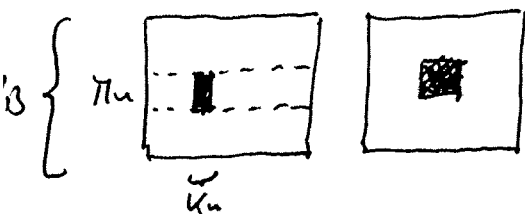


$kij$



$2n^2$  instructions:  
 -  $n^2$  indep. multi  
 -  $n^2$  " " adds  
 $n^2 + 2n$  "live" variables

More "register pressure" since larger working set (always when more ILP)



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

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

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

for  $i' = i : M_u : i + N_B - 1$

for  $j' = j : M_u : j + N_B - 1$

for  $k' = k : K_u : k + N_B - 1$

for  $k'' = k' : 1 : k' + K_u - 1$

for  $i'' = i' : 1 : i' + M_u - 1$

for  $j'' = j' : 1 : j' + M_u - 1$

$$c[i''j''] = c[i''j''] + a[i''k''] b[k''j'']$$

mini-  
MMM

micro-  
MMM

⊗

$M_u, N_u, K_u$  become search parameters, we assume  $\text{all} / N_B$

Bound:  $M_u + N_u + M_u N_u \leq N_R$  (no. of registers)

line variables in ⊗

$K_u$ : later

#### 4.) Basic block optimizations

step 1: unroll ⊗

$$c_{...} = c_{...} + a_{...} b_{...}$$

$$c_{...} = c_{...} + a_{...} b_{...}$$

⋮

$M_u N_u$  many

$a_{...}, b_{...}$  reused

$c_{...}$  not

step 2: scalar replacement

Typically we would only do it on  $a_{...}, b_{...}$  (since reused) but since we later unroll the  $k''$ -loop,  $c_{...}$  will be reused as well.

$$t_0 = c_{...}$$

$$t_1 = a_{...}$$

$$t_2 = b_{...}$$

⋮

$$t_3 = t_0 + t_1 t_2$$

⋮

$$c_{...} = t_3$$

⋮

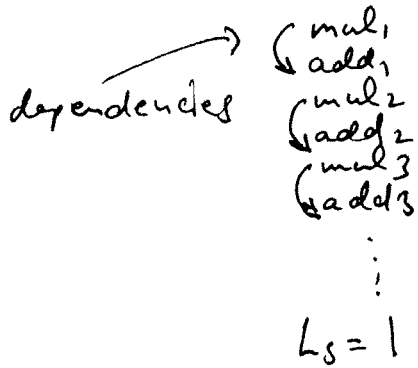
loads  $(M_u + N_u + M_u N_u)$  many  
→ all fit into register

computation  $M_u N_u$  adds  
 $M_u N_u$  multi

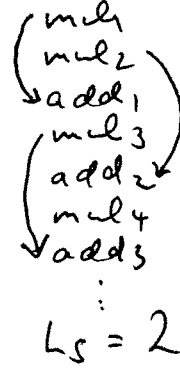
stores  $M_u N_u$  many

step 3: optimize computation part by "skewing"

computation  
as is:



skewed:  
( $L_s = 2$ )



skewed by  $L_s$ :

dependent mul  $\rightarrow$  add  
separated by  $2L_s - 1$   
instructions  
(except for beginning)

$L_s$  is search parameter

Bound:  $M_u + N_u + M_u N_u + L_s \leq N_{I2}$

Also: Not all loads are done in the beginning. Only  
block of  $I_F$  loads and later blocks of  $N_F$  loads  
as needed.

step 4: unroll  $K_u$  loop

$K_u$  is limited by l-cache size

Bound:  $K_u \leq \frac{N_B}{2}$

Example:  $K_u = 3$

```

load
comp
store
- load
comp
store
- load
comp
store
  
```

all those write to  
the same location (c...)

So now we have reuse  
on c...

$\Rightarrow$  get rid of first 2 store

Result:

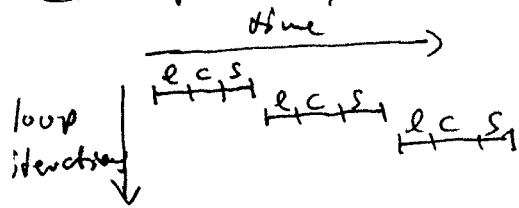
```

load
comp
load
comp
load
comp
store
  
```

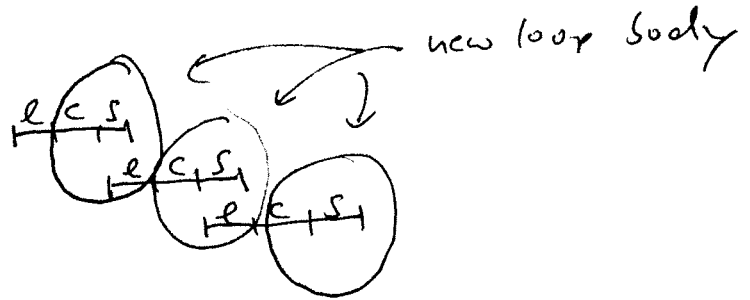
step 5: (software pipelining)

Reorganize k'-loop so loads of one iteration are moved to previous iteration

Conceptually:



$\Rightarrow$



l = load

c = compute

s = store

5.) Buffering to avoid TCIB misses

Comes a little later