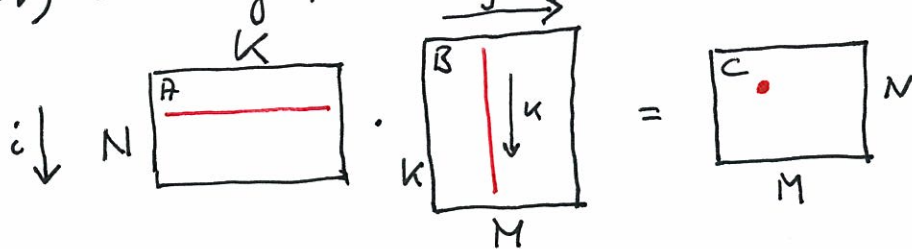


# How to optimize TTTT

The presentation follows the way the program generator ATLAS does it. ATLAS identifies optimization parameters (e.g., the blocking size) and uses search to find the best choices. Alternatively, a model can be used to determine each parameter (see paper on website). We discuss both.

0.) Starting point: standard triple loop



```

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

} computes  $C = C + AB$

Important cases (from most to least): based on usage in LAPACK

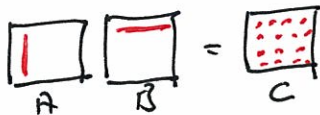
- two out of  $N, K, M$  are small
- one " "
- none " "

1.) Loop order:  $i, j, k$  loops can be permuted into any order

- $i, j, k$ :  $B$  is reused, good if  $M < N$
- $j, i, k$ :  $A$  " " "  $N < M$

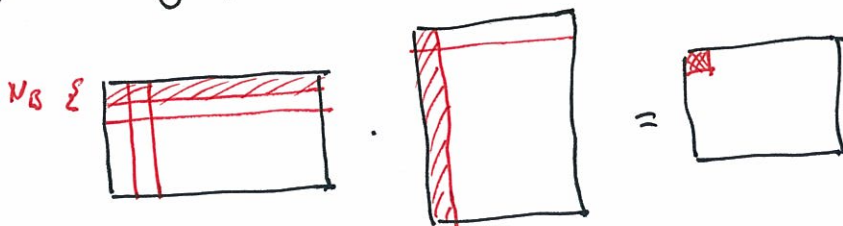
ATLAS includes versions for both

- other choices are bad, e.g.,  $k, i, j$ :



poor temporal locality w.r.t.  $C$

2.) Blocking for cache



assume  $N_B | N, K, M$

results in a six-fold nested loop

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

Formally obtained from the triple loop through loop tiling & loop exchange

mini-NTM multiplies  $N_B \times N_B$  blocks

for  $i' = i : i + N_B - 1$   
 for  $j' = j : j + N_B - 1$   
 for  $k' = k : k + N_B - 1$   
 $C[i', j'] = C[i', j'] + a[i', k'] b[k', j']$

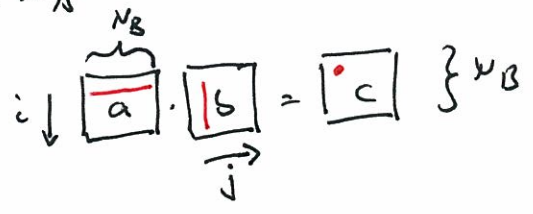
ATLAS: uses search to find best  $N_B$   
 bound:  $N_B^2 \leq C_1$  (cache size)

Model: explained next, model refined in steps

a.) Idea: working set has to fit in cache

easy bound:  $|working\ set| = 3N_B^2$

$\Rightarrow 3N_B^2 \leq C_1$



b.) Closer analysis:

$N_B^2 + N_B + 1 \leq C_1$

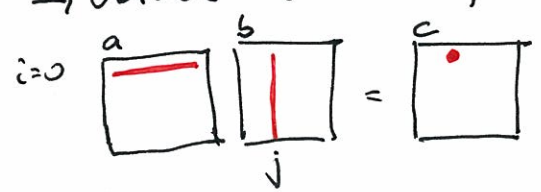
↑ all of b    ↑ row of a    ↑ element of c

c.) Take into account cache block size  $B_1$

$\lceil \frac{N_B^2}{B_1} \rceil + \lceil \frac{N_B}{B_1} \rceil + 1 \leq \frac{C_1}{B_1}$

(this just translates b.) into cache block units)

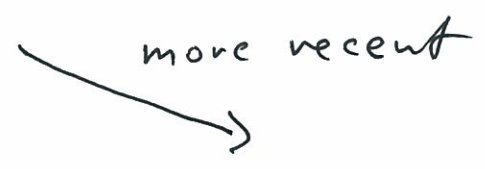
d.) Take into account LRU replacement  
 → build a history of elements being accessed



$i=0$ :  $a_{00} b_{00} \dots a_{0N_B-1} b_{N_B-1,0} c_{00}$  ( $j=0$ )  
 $a_{00} b_{01} \dots \dots \dots b_{N_B-1,1} c_{01}$  ( $j=1$ )  
 ...  
 $a_{00} b_{0N_B-1} \dots \dots \dots b_{N_B-1, N_B-1} c_{0N_B-1}$  ( $j=N_B-1$ )

corresponding history:  
 $b_{00} \dots b_{N_B-1,0} c_{00}$   
 $b_{01} \dots b_{N_B-1,1} c_{01}$   
 ...

$a_{00} b_{0N_B-1} \dots a_{0N_B-1} b_{N_B-1, N_B-1} c_{0N_B-1}$



### Observations:

- all of  $b$  has to fit into cache for next iteration  $i=1$
- when  $i=1$ , row 1 of  $a$  will not cleanly replace row 0 of  $a$
- when  $i=1$ , element of  $c$  will not cleanly replace previous elements of  $c$

⇒ This has to fit:

- entire  $b$
- 2 rows of  $a$  (here:  $a_{0*}, a_{1*}$ )
- 1 row of  $c$  (here:  $c_{0*}$ )
- 1 element of  $c$  (here:  $c_{10}$ )

$$\Rightarrow \left\lceil \frac{N_{13}^2}{B_1} \right\rceil + 3 \left\lceil \frac{N_{13}}{B_1} \right\rceil + 1 \leq \frac{C_1}{B_1}$$

e.) Take into account locality for registers (next opt.)

$$\left\lceil \frac{N_{13}^2}{B_1} \right\rceil + 3 \left\lceil \frac{N_{13} N_u}{B_1} \right\rceil + \left\lceil \frac{N_u N_u}{B_1} \right\rceil \leq \frac{C_1}{B_1}$$

Pick largest  $N_{13}$  that satisfies, shrink so  $N_u, N_u | N_{13}$  (avoids clean-up code)

3.) Blocking mini- $MMMs$  for registers into micro- $MMMs$   
revisit the question of loop order:

$ijk$ : 

for fixed  $i, j$ :

- $2n$  instructions
  - $n$  independent multiplies
  - $n$  dependent adds
- ( $\approx \log_2(n)$  steps)

$kij$ : 

for fixed  $k$ :

- $2n^2$  instructions
- $n^2$  independent multiplies
- $n^2$  dependent adds

⇒ good ILP (but ~~flanger~~ working set)

Result: micro- $MMMs$  with  $kij$  loop order for ILP

Code: for  $i = 0 : N_B : N-1$   
 for  $j = 0 : N_B : N-1$   
 for  $k = 0 : N_B : K-1$   
 for  $i' = i : M_u : i + N_B - 1$   
 for  $j' = j : N_u : j + N_B - 1$   
 for  $k' = k : K_u : k + N_B - 1$  ← ①  
 for  $k'' = k' : k' + K_u - 1$  ← ②  
 for  $i'' = i' : i' + M_u - 1$   
 for  $j'' = j' : j' + N_u - 1$   
 $C_{i''j''} = C_{i''j''} + a_{i''k''} b_{k''j''}$

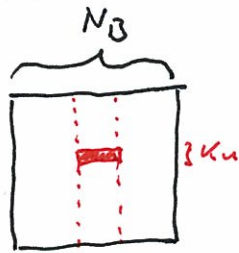
mini-NTM

micro-NTM

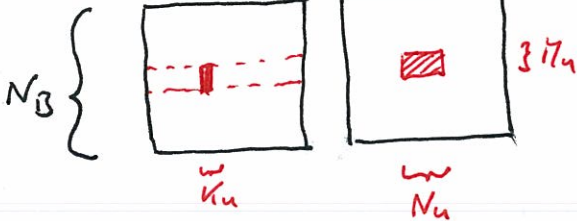
⊗

~~ATLAS~~

are multiplied



mini-NTM  
micro-NTM



ATLAS: Uses search to find best  $M_u, N_u, K_u$   
 Bound:  $M_u + N_u + M_u N_u \leq N_R$  (# registers)  
 size of working set in  $\otimes$   
 = no. of live variables

Model: Use largest  $M_u, N_u$  that satisfy this equation and  $M_u \approx N_u$

4.) Basic block optimizations

step 1: unroll micro-NTM

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

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

⋮

step 2: scalar replacement  
 (all elements  $a_{...}, b_{...}, c_{...}$  are reused)

