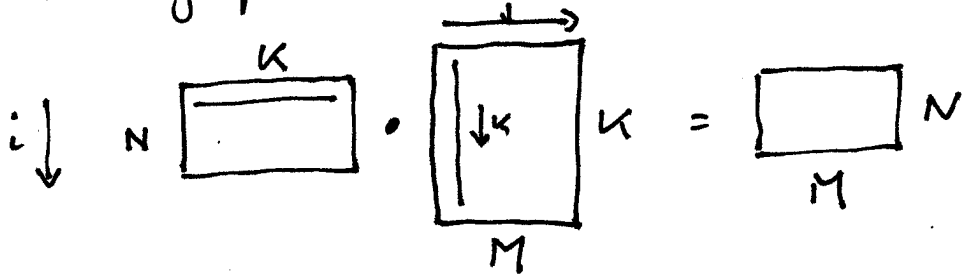


MMM Generation with ATLAS

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 + aikbkj
    
```

1.) loop order

- i, j, k can be permuted into any order
- ijk: B is reused, good if B is smaller than A: $M < N$
- jik: A is reused
- ATLAS generates versions for both
- other choices are bad, e.g., kij

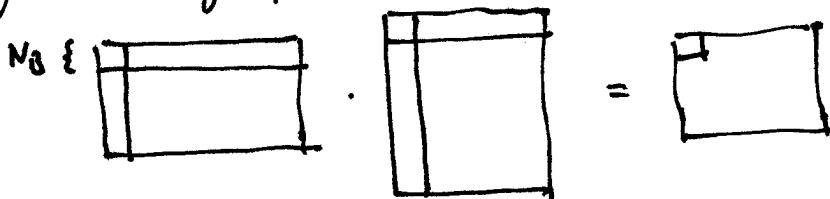


col. row
 \Rightarrow no output locality

```

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

2.) blocking for cache (assume one cache)



for simplicity
 $N_0 | N, M, K$

```

for i = 0:N0:N-1
  for j = 0:N0:M-1
    for k = 0:N0:K-1
      for i' = i+1:i+N0-1
        for j' = j+1:j+N0-1
          for k' = k+1:k+N0-1
            ci'j' = ci'j' + ai'k'bk'j'
          
```

mini-MMM

- formally done using loop tiling and loop exchange

loop tiling: for $i = 0 : 1 : N-1$
 $a_i \dots$

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

- N_B becomes a search parameter in ATLAS
 trivial bound: $N_B^2 \leq C$ (cache size)

3.) blocking for registers

Now: k as outermost loop for instruction parallelism

ijk :



$2n$ instructions:

- n parallel multiplies
- n dependent adds (at least $\log_2(n)$ steps)

and: $2n+1$ live variables

kij :

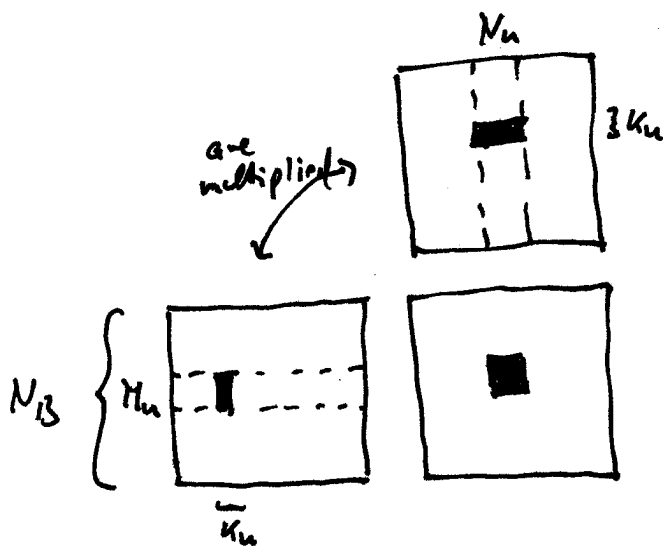


$2n^2$ instructions:

- n^2 parallel multiplies
- n^2 parallel adds

$n^2 + 2n$ live variables

"register pressure" since large working set



Code: $\begin{aligned} &\text{for } i = 0 : N_1 : N-1 \\ &\quad \text{for } j = 0 : N_2 : N-1 \\ &\quad\quad \text{for } k = 0 : N_3 : K-1 \\ &\quad\quad\quad \text{for } i' = i : M_1 : i + M_1 - 1 \\ &\quad\quad\quad\quad \text{for } j' = j : M_2 : j + M_2 - 1 \\ &\quad\quad\quad\quad\quad \text{for } k' = k : M_3 : k + M_3 - 1 \\ &\quad\quad\quad\quad\quad\quad \text{for } k'' = k' : 1 : k' + M_3 - 1 \\ &\quad\quad\quad\quad\quad\quad\quad \text{for } i'' = i' : 1 : i' + M_1 - 1 \\ &\quad\quad\quad\quad\quad\quad\quad\quad \text{for } j'' = j' : 1 : j' + M_2 - 1 \\ &\quad\quad\quad\quad\quad\quad\quad\quad\quad C_{i''j''} = C_{i''j''} + a_{i''k''} b_{k''j''} \end{aligned}$

possibly suffer a tile of C (see s. below)

micro-tiling {

- (M_1, M_2, M_3) become search parameters
 bound $\underbrace{M_1 M_2 + M_1 + M_2}_{\text{live variables}} \leq N_2$ (number of registers)

4.) basic block optimizations

step 1: unroll \otimes and do scalar replacement

unroll: $\left\{ \begin{aligned} C_{ij} &= C_{ij} + a_{ik} b_{kj} \\ &\vdots \end{aligned} \right.$
 many of these a, b, c are arrays

scalar replacement: replace array elements with scalar variables

$\begin{aligned} t_0 &= C_{ij} \\ t_1 &= a_{ik} \\ t_2 &= b_{kj} \\ &\vdots \\ t_0 &= t_0 + t_1 t_2 \\ &\vdots \\ C_{ij} &= \dots \\ &\vdots \end{aligned}$

} loads
 } computation
 } stores

- enables register allocation by compiler and other optimizations (for other numerical problems)

step 2: reorder load, adds, mults, stores to take dependencies and latencies into account

example: $C_{ij} = C_{ij} + a_{in} b_{kj}$

$\Leftrightarrow \begin{cases} t = a_{in} b_{kj} \\ C_{ij} = C_{ij} + t \end{cases}$ } dependent, problem if $f_{mult} \text{ latency} > f_{add} \text{ latency}$

solution: "skewing" with L_s

mul₁
add₁
mul₂
add₂
⋮

standard: $L_s = 1$

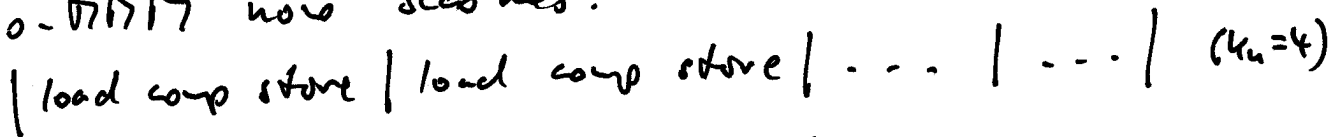
mul₁
mul₂
add₁
mul₃
add₂
mul₄
add₃
⋮
 $L_s = 2$

(L_s) is search parameter

other rounding: interleave loads and adds/mults

step 3: unroll k loop, k_u is limited by size of I-cache

→ micro-OTM now becomes:



- software pipelining: move loads from one iteration of k loop to previous iteration

5.) Buffering:

idea: copy working set that is reused in a way that may cause thrashing into contiguous memory

Example:

