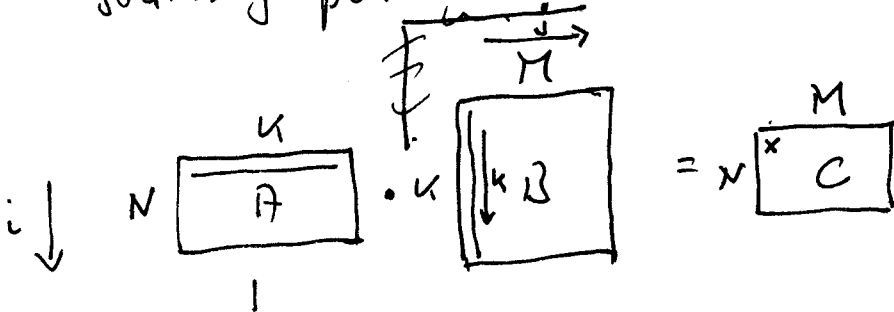


MMM - code generation with ATLAS

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
    
```

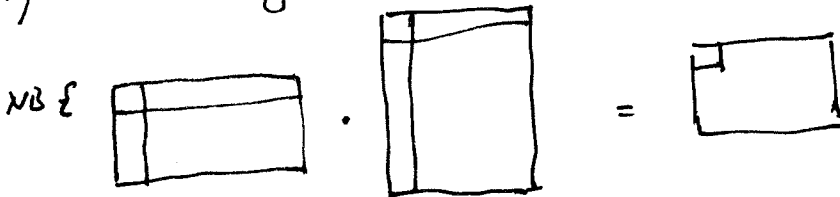
output locality } Cost = 2NMK

1.) loop order ijk or jik?

ijk: B is reused
jik: A is reused

⇒ for locality: if $M < N$: ijk
if $N < M$: jik

2.) blocking into square $N_B \times N_B$ blocks, for cache perfor



```

for i = 0: N_B: N-1
  for j = 0: N_B: M-1
    for k = 0: N_B: K-1

```

assuming $N_B | N, M, K$

mini-MMM

```

for i' = 0:1: i+N_B-1
  for j' = 0:1: j+N_B-1
    for k' = 0:1: k+N_B-1

```

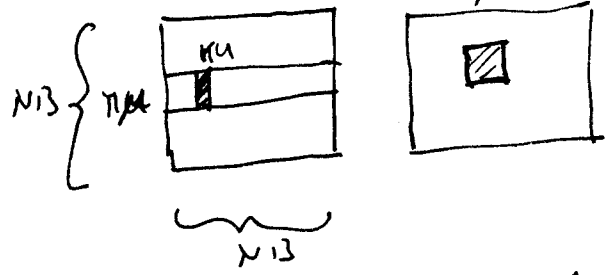
$$C_{i'j'} = C_{i'j'} + A_{i'k'} * B_{k'j'}$$

⇒ for locality

~~N_B is search parameter~~

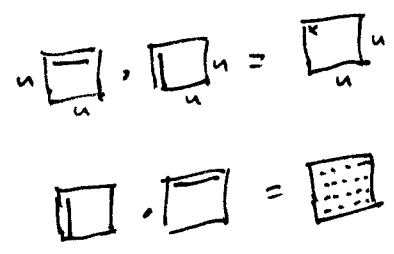
N_B search parameter

3.) Scheduling for registers



now: - k_{ij} order
 \Rightarrow want good locality
 \Rightarrow but better instruction level parallelism

explain instruction level parallelism:



parallelism:
 consecutive instructions:
 - u parallel multiplies
 - u dependent adds (at least $\log_2(u)$ steps)
 $2u-1$ instructions
 $2u^2$ instructions:
 - u^2 parallel multiplies
 - u^2 parallel adds

but: $2u+1$ live vars vs. u^2+2u live vars \leftarrow register pressure
 "larger working set"

need: $u^2+2u \leq N_R$ (# registers)

```

code:
for i = 0: N3: N-1
  for j = 0: N3: N-1
    for k = 0: N3: u-1
      for i' = 0: N3: i+u-1
        for j' = 0: N3: j+u-1
          for k' = 0: N3: k+u-1
            for i'' = i': 1: i' + N3 - 1
              for j'' = j': 1: j' + N3 - 1
                C[i'',j''] = C[i'',k''] + A[i'',k''] * B[k'',j'']
    }
  }
}
    
```

micro- $\pi\pi\pi$ }
 unroll + scalar replacement

necessary: $\pi_u \cdot \kappa_u + \pi_u + \kappa_u \leq N_R$ measured

(π_u, κ_u) search parameters

unrolling 1: "i, j" loop

- reduce loop overhead
- scalar replacement
 - enables register allocation + other optim.
 - do not reuse variable names (artificial dependencies)

4.) instruction interleaving/scheduling

- adds/mults alternately

mul₁
add₁
mul₂
add₂
⋮
L_s = k

mul₁
(mul₂
add₁
mul₃
add_{2k}
mul₄
add₃)

(L_s) search parameter

L_s = 2

$$c_{ij} = c_{ij} + a_{ik} \cdot b_{kj}$$

⌊ $t = a_{ik} \cdot b_{kj}$ $\left. \begin{array}{l} \\ c_{ij} = c_{ij} + t \end{array} \right\}$ dependency, problem with full latency
solution: skewing
increases register pressure

5.) unrolling 2: k' loop (k_u) is search parameter
- limited by size l-cache