

Model-Based ATLAS (Yotov. et al.)

parameters used include

- C_1 L1 cache size
- B_1 cache line size
- C_I l-cache size
- L_x Fl.p. latency

1.) Estimating N_B (assuming L1-cache fully associative)
refine model step by step

$$i \downarrow \begin{matrix} N_B \\ \boxed{} \\ K \end{matrix} \cdot \begin{matrix} \boxed{} \\ M \end{matrix} = \begin{matrix} \boxed{} \end{matrix}$$

Blocked into mini-MMM's:

$$i \downarrow \begin{matrix} \boxed{} \\ \text{outer} \\ \text{most} \\ \text{loop} \end{matrix} \cdot \begin{matrix} \boxed{} \\ j \end{matrix} = \begin{matrix} \boxed{} \\ N_B \end{matrix}$$
$$A \cdot B = C$$

a.) working set $3N_B^2$
 $\Rightarrow 3N_B^2 \leq C_1$

b.) closer analysis

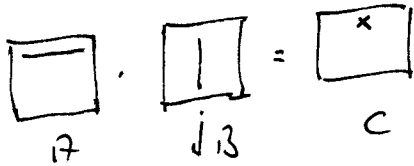
$$N_B^2 + N_B + 1 \leq C_1$$

\uparrow \uparrow \uparrow
 B row of A element of C

c.) take into account cache line size
units are cache line

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

d.) take into account LRU replacement



blocks

$$A_{00}, B_{0j}, A_{01}, B_{1j}, \dots, A_{0, N_B-1}, B_{N_B-1, j}, C_{0j}$$

more recent

for all j : (computing ~~some~~ rows of C)
 $(i=0)$ want to keep for $i=1$

$$A_{00} B_{00}, \dots, A_{0, N_B-1} B_{N_B-1, 0}, C_{00},$$

$$A_{00} B_{01}, \dots, A_{0, N_B-1} B_{N_B-1, 1}, C_{01},$$

$$\dots$$

$$A_{00} B_{0, N_B-1}, \dots, A_{0, N_B-1} B_{N_B-1, N_B-1}, C_{0, N_B-1}$$

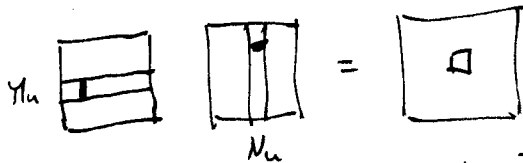
in the next iteration ($i=1$) we want to reuse B

\Rightarrow for $i=1$ also A_{0x} and C_{0x}
1. row A 1. row C

have to fit
 (i.e., entire B
 2 rows of A
 1 row of C
 1 elem. of C)

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

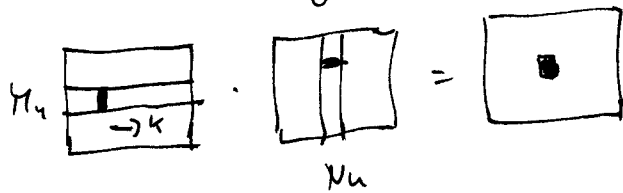
e.) take into account register tiling



$$\Rightarrow \left\lceil \frac{N_B^2}{B_1} \right\rceil + 3 \left\lceil \frac{N_B \cdot M_u}{B_1} \right\rceil + \left\lceil \frac{M_u \times M_u}{B_1} \right\rceil \leq C_1$$

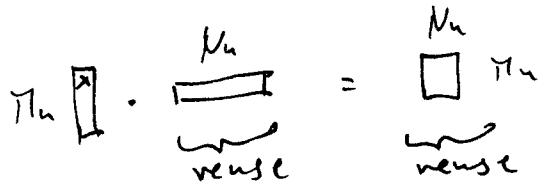
f.) make sure $M_u, N_u \mid N_B$ (no clean-up code)

2.) Estimating π_n and ν_n



(kji-order)
in micro-TMM

a.)
micro
TMM

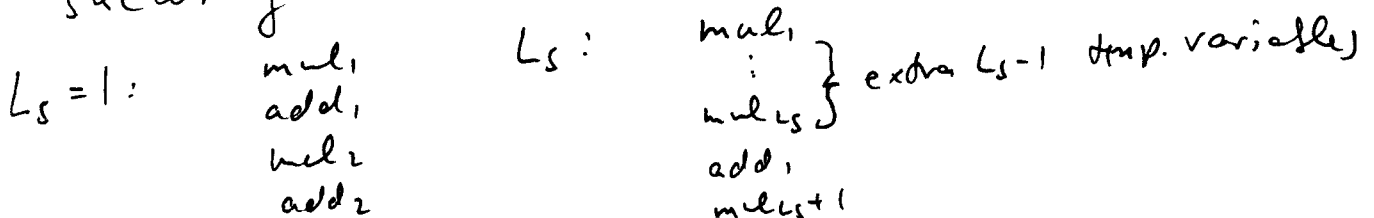


$$\Rightarrow \pi_n \nu_n + \nu_n + 1 \leq N_{12}$$

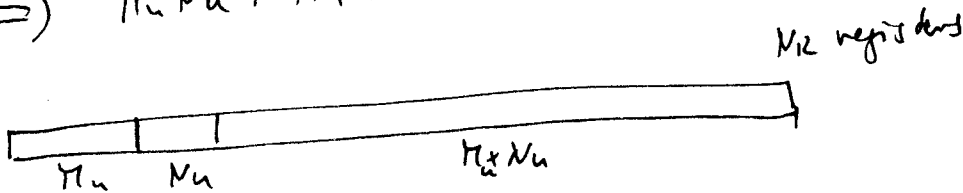
b.) scalar replacement: loads entire ν and π

$$\Rightarrow \pi_n \nu_n + \pi_n + \nu_n \leq N_{12}$$

c.) skewing



$$\Rightarrow \pi_n \nu_n + \pi_n + \nu_n + L_s \leq N_{12}$$



3.) Estimating ν_n

- choose sub that fits into 1-code
- ν_n / N_{12} (avoid clean-up code)

4.) Estimating L_s (assuming one FP unit)

$$\Rightarrow 2L_s - 1 \geq L_x$$

$$\Rightarrow L_s = \left\lceil \frac{L_x + 1}{2} \right\rceil$$

