# Scalable Matrix Architecture (Fundamentals and Vector-Matrix Facility)

*José Moreira*
IBM Corporation

*Abel Bernabeu*
Esperanto

# Comments on the state-of-the-art for matrix computing

- Most high-performance matrix computations rely on the following kernel:

$$C_{m \times n} = A_{m \times K} \times B_{K \times n}$$

$C_{m \times n}$ is a register-resident matrix, with $m \approx n$, and

$A_{m \times K}$ and $B_{K \times n}$ are streamed from memory, with $m \ll K$ and $n \ll K$

- Matrix $A_{m \times K}$ is organized as a vector of $K$ elements, each element a column vector of $m$ elements of $A = \begin{bmatrix} A^0 & A^1 & \cdots & A^{K-1} \end{bmatrix}$
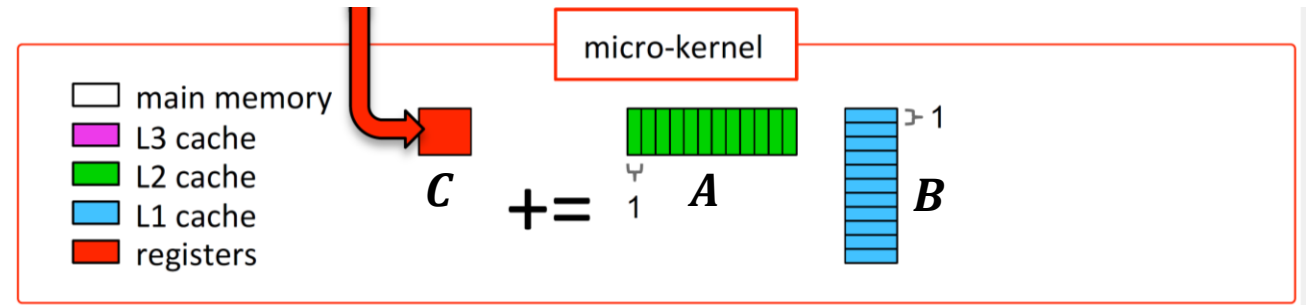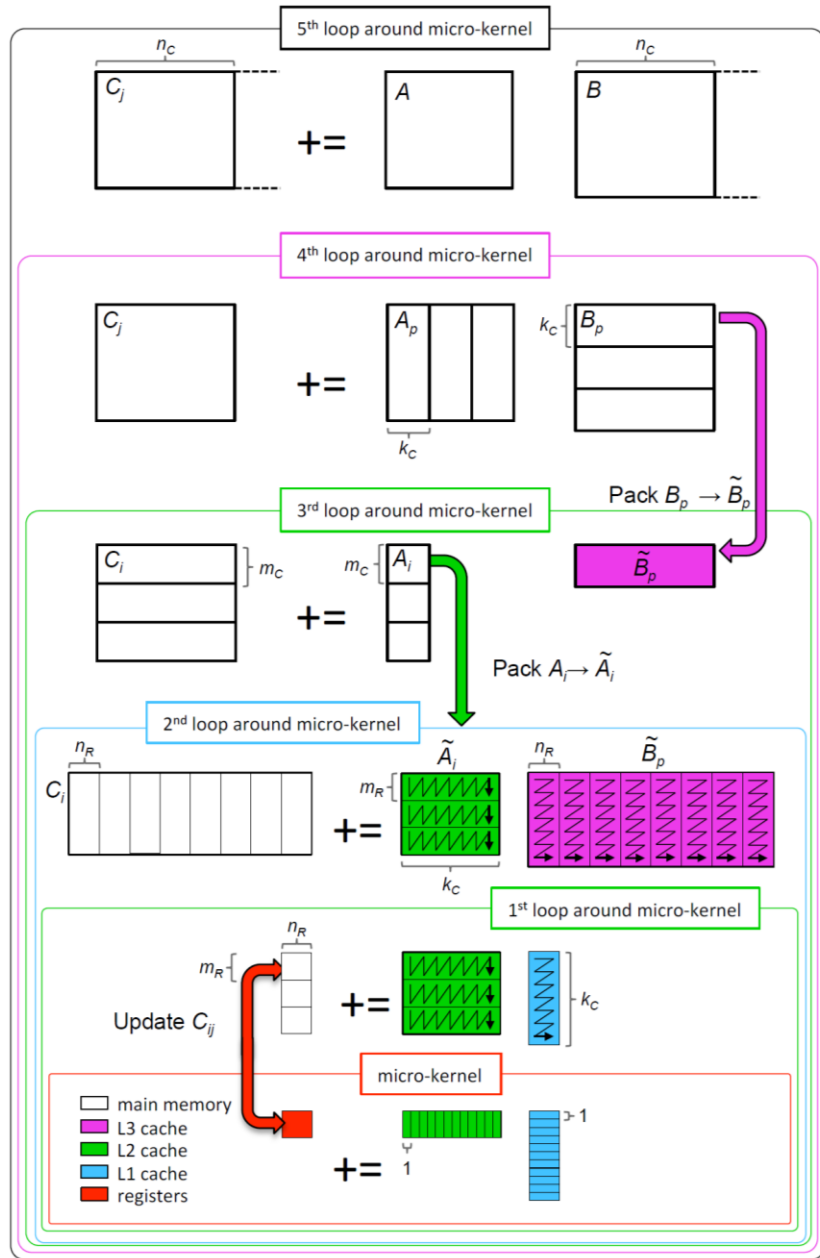
- Matrix $B_{K \times n}$ is organized as a vector of $K$ elements, each element a row of $n$ elements of $B = \begin{bmatrix} B_0 \\ B_1 \\ \vdots \\ B_{K-1} \end{bmatrix}$

- We will consider only the case of deterministic values of $C$, as produced by the following algorithm:

$$C \leftarrow 0; \text{for } (k = [0, K)) \; C \leftarrow A^k \times B_k + C$$

# Anatomy of a high-performance matrix multiply

# Performance bounds

- Let $\Delta$ be the latency (in cycles) of an elemental multiply-add operation

- Computing each element of $C$, as described above, requires evaluating a dependence chain of $K$ multiply-adds, and therefore takes time $K\Delta$

$$C \leftarrow A^k \times B_k + C$$

$$\Delta$$

- The computation of $C$ requires $mnK$ multiply-adds

- The maximum computation rate that can be sustained is

$$R = \frac{mnK}{K\Delta} = \frac{mn}{\Delta} \quad \text{madds/cycle}$$

- This sets an upper bound on performance, dictated by the size of the $C$ panel that can be kept in registers (an architectural parameter) and the latency of a multiply-add (a design/technology parameter) – *Note*: integer arithmetic is more forgiving

- For modern server processor and 32-bit floating-point elements, a reasonable value of $\Delta$ is 4 cycles – the range $\Delta \in [2,8]$ should cover most cases

# Some configurations (just to get a feel)

- Let the architected space consist of 32 vector registers of size $L$ words (32-bit)
  - Then, $mn \leq 32L$ and the maximum computation rate is $R = {}^{32L}/_4 = 8L$ madds/cycle (single-precision floating-point)

- The maximum computation rate of a vector-based matrix facility, using vector registers only, scales with $L$

- For $L = 4$, $R = 32$ madds/cycle

- For $L = 16$, $R = 128$ madds/cycle

- *Note*: In the above, we use all architected space for **C**, but in a load/store architecture we need to reserve space for **A** and **B**, so cut the bound in half
  - $mn \leq 16L$ and the maximum computation rate is $R = {}^{16L}/_4 = 4L$ madds/cycle
  - For $L = 4$, $R = 16$ madds/cycle
  - For $L = 16$, $R = 64$ madds/cycle

- *Note*: $4L$ madds/cycle is the performance equivalent of 4 vector pipes operating concurrently

# Conclusions

- The size of the architected register space imposes an upper bound on the performance of deterministic matrix multiply kernels

- If we use the architected scalable vector register space ($32L$ words) to hold the panel of $\boldsymbol{C}$ our performance upper bound, for a reasonable choice of latency parameters, is $8L$ madds/cycle – a more realistic value is $4L$ madds/cycle (single-precision floating-point)

- We advocate to pursue the vector register-based matrix extensions first, starting with three possible options:
  - Option A : 1 matrix/1 vector register     ($R = 4L$ madds/cycle with $4$ matrix pipes)
  - Option B : 1 matrix/$n$ vector registers    ($R = 4L$ madds/cycle with $1$ matrix pipe)
  - Option C : $n$ matrices/1 vector register    ($R = 4L$ madds/cycle with ${}^4/_\lambda$ matrix pipes)