

Integrated Matrix Extension Option C (Common-Type Variant)

José E. Moreira (IBM)

May 2024

1 Introduction

This is a strawman for Option C (Multiple fixed-size matrix tiles per vector register) of the Integrated Matrix Extension (IME). As such, it is meant to guide and facilitate discussion by providing a concrete draft document. There should be no expectation that any content in the strawman will make into the final specification, particularly in the early stages of the work. As we evolve our work in the IME Task Group, refined concepts from this strawman may be promoted to an actual draft specification.

This document focuses on the *common-type* variant of the extensions. That is, all matrix data types involved in one instruction are the same. This is distinct from a *mixed-type* variant that supports matrices of different data types in one instruction.

Notation: Matrices are represented by bold-face capital letters (**A**, **B**, **C**). The element at row i , column j of matrix **A** can be denoted as either $\mathbf{A}(i, j)$ or a_i^j . $\mathbf{A}(i, :)$ represents the i -th row of matrix **A** and $\mathbf{A}(:, j)$ its j -th column. $\mathbf{A}(i : +h, j : +w)$ denotes the two-dimensional section of matrix **A** consisting of the h rows beginning in row i and the w columns beginning in column j .

2 Matrix tiles

A *matrix tile* is a contiguous rectangular section of a matrix. Matrix tiles are defined by the scalar data type of their elements (same as the type of the matrix elements) and their *shape* $\lambda \times \lambda$, where λ is the number of rows and columns of the tile. All tiles supported in this variant of Option C are square.

Notation: The tile A_i^j corresponds to the matrix section $\mathbf{A}(i\lambda : +\lambda, j\lambda : +\lambda)$.

The shape ($\lambda \times \lambda$) and number (L) of tiles stored in a vector register is an implementation choice for each elemental data type, but it must follow the following constraint. Given a vector length (**VLEN**) and a *matrix element width* (**MEW**), the vector length must match the element width times the tile size ($\lambda \times \lambda$ elements) times the number of tiles in the vector register. In other words, the following equality must hold:

$$\text{VLEN} = \text{MEW} \cdot \lambda^2 \cdot L.$$

Table 1 shows valid values of the $\langle \lambda, L \rangle$ pair for different combinations of **VLEN** and **MEW**. We will show below how to generate binary code that is agnostic to the choice of **VLEN**, λ , and L , and therefore portable, in both functionality and performance, across any implementation that follows the specification outlined in this document.

3 Matrix instructions

All matrix instructions in the common-type variant operate on matrix tiles of a common scalar data type. Each architected vector register stores a vector of matrix tiles, and instructions take as arguments either one, two, or three vector register identifiers. Tiles in a vector register **A** are identified as $\mathbf{A}[0], \mathbf{A}[1], \dots, \mathbf{A}[L-1]$. Additional arguments to matrix instructions can include *index registers*, containing unsigned integer values, and/or vector masks.

$\langle \lambda, L \rangle$				
VLEN	MEW			
	8	16	32	64
32	$\langle 2, 1 \rangle$	–	–	–
64	$\langle 2, 2 \rangle$	$\langle 2, 1 \rangle$	–	–
128	$\langle 2, 4 \rangle, \langle 4, 1 \rangle$	$\langle 2, 2 \rangle$	$\langle 2, 1 \rangle$	–
256	$\langle 2, 8 \rangle, \langle 4, 2 \rangle$	$\langle 2, 4 \rangle, \langle 4, 1 \rangle$	$\langle 2, 2 \rangle$	$\langle 2, 1 \rangle$
512	$\langle 2, 16 \rangle, \langle 4, 4 \rangle, \langle 8, 1 \rangle$	$\langle 2, 8 \rangle, \langle 4, 2 \rangle$	$\langle 2, 4 \rangle, \langle 4, 1 \rangle$	$\langle 2, 2 \rangle$
1024	$\langle 2, 32 \rangle, \langle 4, 8 \rangle, \langle 8, 2 \rangle$	$\langle 2, 16 \rangle, \langle 4, 4 \rangle, \langle 8, 1 \rangle$	$\langle 2, 8 \rangle, \langle 4, 2 \rangle$	$\langle 2, 4 \rangle, \langle 4, 1 \rangle$
2048	$\langle 2, 64 \rangle, \langle 4, 16 \rangle, \langle 8, 4 \rangle, \langle 16, 1 \rangle$	$\langle 2, 32 \rangle, \langle 4, 8 \rangle, \langle 8, 2 \rangle$	$\langle 2, 16 \rangle, \langle 4, 4 \rangle, \langle 8, 1 \rangle$	$\langle 2, 8 \rangle, \langle 4, 2 \rangle$

Table 1: Valid values of the $\langle \lambda, L \rangle$ pair for different combinations of matrix element width (MEW) and vector length (VLEN). An implementation must choose one of the valid pairs. For example, when processing single-precision floating-point data (MEW = 32), 512-bit vectors (VLEN = 512) can contain either 4 tiles of 2×2 elements ($\langle 2, 4 \rangle$), or 1 tile of 4×4 elements ($\langle 4, 1 \rangle$).

3.1 Matrix load instructions

Matrix load instructions have the general form

$$\text{mload}\{\text{variant}\}\langle T, \dots \rangle(\text{vd}, \mathbf{A}(i : +h, j : +w))$$

where **mload** identifies this as a matrix load instruction and **variant** is a specific variant. T is the matrix element data type, which specifies MEW. $\mathbf{A}(i : +h, j : +w)$ is the section of matrix \mathbf{A} to be loaded in the destination register group beginning with **vd**. Matrix load instructions support different forms of group multipliers, to be specified for each variant.

3.1.1 mload – matrix load (vanilla version)

Syntax: $\text{mload}\langle T, \text{RMUL}, \text{maxrows}, \text{CMUL}, \text{maxcols} \rangle(\text{vd}, \mathbf{A}(i, j), \text{lda})$

Arguments: $\mathbf{A}(i, j)$ is the address of the first element of the section of matrix \mathbf{A} to be loaded in the register group beginning with register **vd**. The leading dimension of \mathbf{A} , necessary to compute the address of the other elements in the section, is passed as **lda**.

Semantics: This instruction takes the $\text{LMUL} = \text{RMUL} \cdot \text{CMUL}$ vector registers starting at register **vd** and organizes them as a $\text{RMUL} \times \text{CMUL}$ two-dimensional array of registers, capable of holding up to $\text{RMUL} \cdot \text{CMUL} \cdot \lambda^2 L$ elements of type T . It then loads the array section $\mathbf{A}(i : +\max(\text{maxrows}, \text{RMUL} \cdot \lambda), j : +\max(\text{maxcols}, \text{CMUL} \cdot \lambda L))$ into those registers. Matrix \mathbf{A} is interpreted as being stored in row-major order, with a leading dimension of **lda**. Any elements of the target register group that are not loaded from \mathbf{A} are set to 0 (or, in the general case, a specified identity value).

Example: Let VLEN = 512 and MEW = 64. Consider the instruction

$$\text{mload}\langle \text{fp64}, 2, \infty, 2, \infty \rangle(\text{v00}, \mathbf{A}(0, 0), \text{lda}).$$

In this case, $\langle \lambda, L \rangle$ must be $\langle 2, 2 \rangle$ and vector registers **v00**, **v01**, **v02**, **v03** are loaded with a 4×8 section of matrix \mathbf{A} as follows:

$$\begin{aligned} \text{v00} &= \left[\begin{bmatrix} a_0^0 & a_0^1 \\ a_1^0 & a_1^1 \end{bmatrix} \quad \begin{bmatrix} a_0^2 & a_0^3 \\ a_1^2 & a_1^3 \end{bmatrix} \right] & \text{v01} &= \left[\begin{bmatrix} a_0^4 & a_0^5 \\ a_1^4 & a_1^5 \end{bmatrix} \quad \begin{bmatrix} a_0^6 & a_0^7 \\ a_1^6 & a_1^7 \end{bmatrix} \right] \\ \text{v02} &= \left[\begin{bmatrix} a_2^0 & a_2^1 \\ a_3^0 & a_3^1 \end{bmatrix} \quad \begin{bmatrix} a_2^2 & a_2^3 \\ a_3^2 & a_3^3 \end{bmatrix} \right] & \text{v03} &= \left[\begin{bmatrix} a_2^4 & a_2^5 \\ a_3^4 & a_3^5 \end{bmatrix} \quad \begin{bmatrix} a_2^6 & a_2^7 \\ a_3^6 & a_3^7 \end{bmatrix} \right] \end{aligned}$$

3.2 Matrix computation instructions

Matrix computation instructions have the general form

$$\mathbf{m}\{\text{comp}\}\langle T, \otimes, \oplus, \text{LMUL} \rangle(\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots)$$

where \mathbf{m} identifies this as a matrix instruction and comp identifies the specific computation, usually following the nomenclature established by BLAS. T is the matrix element data type, which specifies MEW. The implementation selects a valid $\langle \lambda, L \rangle$ pair compatible with its VLEN. The code will be agnostic to this choice. The multiplication/addition pair (\otimes, \oplus) forms a semiring on type T . (The precision of the operations may be different than the precision of the data type. For example, addition and multiplication for data type **bf16** can be performed in IEEE single-precision.) \mathbf{A} , \mathbf{B} , and \mathbf{C} are vector register identifiers that are used to store the tiles from matrices \mathbf{A} , \mathbf{B} , and \mathbf{C} used in the operation, respectively. \mathbf{A} and \mathbf{B} are input matrices. \mathbf{C} is the result matrix of the computation. A matrix instruction specifies a single group multiplier LMUL. The effective multipliers for each matrix ($\text{EMUL}(\mathbf{A})$, $\text{EMUL}(\mathbf{B})$, and $\text{EMUL}(\mathbf{C})$) are derived from LMUL. For the common-type variant $\text{EMUL}(\mathbf{A}) = \text{EMUL}(\mathbf{B}) = \text{EMUL}(\mathbf{C}) = \text{LMUL}$. By default, $\text{LMUL} = 1$.

3.2.1 mgemmm – matrix multiplication

Syntax: $\text{mgemmm}\langle T, \otimes, \oplus \rangle(\mathbf{A}, \mathbf{B}, \mathbf{C})$

Arguments: \mathbf{A} , \mathbf{B} , and \mathbf{C} are vector registers, each with L matrix tiles of shape $\lambda \times \lambda$.

Semantics: This instruction computes $\mathbf{C}[i] \leftarrow \mathbf{C}[i] \oplus \mathbf{A}[i] \otimes \mathbf{B}[i], \forall i \in [0, L)$.

Example: Let $\text{VLEN} = 512$ and $\text{MEW} = 64$. In this case, $\langle \lambda, L \rangle$ must be $\langle 2, 2 \rangle$ and vector registers \mathbf{A} , \mathbf{B} , and \mathbf{C} each contain two 2×2 tiles of 64-bit elements. The computation performed is illustrated as follows:

$\mathbf{A}[0] : \begin{bmatrix} a_0^0 & a_0^1 \\ a_1^0 & a_1^1 \end{bmatrix}$	$\mathbf{B}[0] : \begin{bmatrix} b_0^0 & b_0^1 \\ b_1^0 & b_1^1 \end{bmatrix} \quad \mathbf{B}[1] : \begin{bmatrix} b_0^2 & b_0^3 \\ b_1^2 & b_1^3 \end{bmatrix}$
$\mathbf{A}[1] : \begin{bmatrix} a_0^2 & a_0^3 \\ a_1^2 & a_1^3 \end{bmatrix}$	$\mathbf{C}[0] : \begin{bmatrix} c_0^0 + a_0^0 b_0^0 + a_1^0 b_1^0 & c_0^1 + a_0^0 b_0^1 + a_1^0 b_1^1 \\ c_1^0 + a_1^0 b_0^0 + a_1^1 b_1^0 & c_1^1 + a_1^0 b_0^1 + a_1^1 b_1^1 \end{bmatrix}$
$\mathbf{C}[1] : \begin{bmatrix} c_0^2 + a_0^2 b_0^2 + a_1^2 b_1^2 & c_0^3 + a_0^2 b_0^3 + a_1^2 b_1^3 \\ c_1^2 + a_1^2 b_0^2 + a_1^3 b_1^2 & c_1^3 + a_1^2 b_0^3 + a_1^3 b_1^3 \end{bmatrix}$	

3.2.2 mgemmm0 – matrix multiplication with $\mathbf{A}[0]$

Syntax: $\text{mgemmm0}\langle T, \otimes, \oplus \rangle(\mathbf{A}, \mathbf{B}, \mathbf{C})$

Arguments: \mathbf{A} , \mathbf{B} , and \mathbf{C} are vector registers, each with L matrix tiles of shape $\lambda \times \lambda$.

Semantics: This instruction computes $\mathbf{C}[i] \leftarrow \mathbf{C}[i] \oplus \mathbf{A}[0] \otimes \mathbf{B}[i], \forall i \in [0, L)$.

3.2.3 mgemmmx – matrix multiplication with $\mathbf{A}[x]$

Syntax: $\text{mgemmmx}\langle T, \otimes, \oplus \rangle(\mathbf{A}, \mathbf{B}, \mathbf{C}, x)$

Arguments: \mathbf{A} , \mathbf{B} , and \mathbf{C} are vector registers, each with L matrix tiles of shape $\lambda \times \lambda$. The index register x identifies the matrix tile $\mathbf{A}[x], x \in [0, L)$.

Semantics: This instruction computes $\mathbf{C}[i] \leftarrow \mathbf{C}[i] \oplus \mathbf{A}[x] \otimes \mathbf{B}[i], \forall i \in [0, L)$.

4 Case study: dgemm

The BLAS `dgemm` routine computes $\mathbf{C} \leftarrow \alpha \mathbf{A} \mathbf{B} + \beta \mathbf{C}$, where \mathbf{C} is a $M \times N$ matrix, \mathbf{A} is a $M \times K$ matrix, \mathbf{B} is a $K \times N$ matrix, α and β are scalars. All data types are double-precision floating-point numbers. (We are ignoring the matrix transpose variants.)

Pseudo code for $\langle \text{VLEN}, \lambda, L \rangle$ agnostic implementation of `dgemm` is shown in Figure 1. It consists of a double-nested parallel loop across the rows and columns of result matrix \mathbf{C} . Each iteration of the loop nest computes an $m \times n$ panel of \mathbf{C} in the *micro-kernel* `μdgemm`, with $m = 4\lambda$ and $n = 4\lambda L$.

```

01 procedure dgemm( $M, N, K, \alpha, \mathbf{A}, \mathbf{B}, \beta, \mathbf{C}$ )
02   for( $I \leftarrow 0; I < M; I \leftarrow I + \text{ml}$ )
03     ml  $\leftarrow \min(M - I, 4\lambda)$ 
04     for( $J \leftarrow 0; J < N; J \leftarrow J + \text{nl}$ )
05       nl  $\leftarrow \min(N - J, 4\lambda L)$ 
06        $\mu\text{dgemm}(K, \alpha, \mathbf{A}(I : 4\lambda, :), \mathbf{B}(:, J : 4\lambda L), \beta, \mathbf{C}(I : 4\lambda, J : 4\lambda L))$ 
07     end
08   end
09 end

```

Figure 1: The (M, N, K) BLAS routine `dgemm` can be implemented as a two-dimensional loop, each iteration computing a $4\lambda \times 4\lambda L$ panel of the result matrix \mathbf{C} . The computation of each panel is performed by a micro-kernel that takes as input a block of 4λ rows of \mathbf{A} and a block of $4\lambda L$ columns of \mathbf{B} .

Pseudo code for the micro-kernel is shown in Figure 2. It corresponds to the mapping of matrices to vector registers shown in Figure 3. The code, including the corresponding binary code, produces the correct result independent of the values of `VLEN`, λ , and L chosen by an implementation, as long as they are consistent. We call this code *geometry agnostic*.

The micro-kernel works as follows. We first zero the contents of the $m \times n$ panel of \mathbf{C} in registers `v16–v31`, using standard RISC-V vector instructions. We then iterate over the inner dimension. The outer loop proceeds in chunks of λL , while the inner loop proceeds in chunks of λ . At each iteration of the outer (k) loop, we load $4L$ tiles of matrix \mathbf{A} in registers `v08–v11`. This is accomplished with the matrix load instruction discussed above. At each iteration of the inner (x) loop, we process one tile of \mathbf{A} from each of those registers. At the beginning of each x iteration we load a $\lambda \times 4\lambda L$ panel of \mathbf{B} . We then multiply each of the 4 tiles of shape $\lambda \times \lambda$ in position x of registers `v08–v11` by the $4L$ tiles of shape $\lambda \times \lambda$ in registers `v12–v15`. This updates the $16L$ tiles of \mathbf{C} in registers `v16–v31`. At the end of the outer loop we will have computed $\mathbf{A} \times \mathbf{B}$. The scaling by scalar α can be accomplished with existing RISC-V vector instructions, as can the addition of $\beta \mathbf{C}$. The loading and storing of the panel of \mathbf{C} can again use matrix load and store (not discussed, but similar to load) instructions.

The binary code in Figure 2 is agnostic to the values of `VLEN`, λ , and L , producing the correct result in every case. The computational intensity (η), defined as the ratio of floating-point operations by the number of elements transferred, can be computed as follows. Each iteration of the outer loop of Figure 2 loads $4\lambda^2 L$ elements of matrix \mathbf{A} . Each iteration of the inner loop of Figure 2 loads $4\lambda^2 L$ elements of matrix \mathbf{B} . Therefore, the total number of elements loaded in each iteration of the outer loop is $4\lambda^2 L(1 + L)$. Each `mgemmx` instruction in the inner loop performs λ^3 multiply-adds for each tile of matrix \mathbf{C} produced, or $\lambda^3 L$ multiply-adds per instruction. Since there are 16 instructions in the inner loop body and that body is executed L time per iteration of the outer loop, we perform $16\lambda^3 L^2$ multiply-adds per iteration of the outer loop. Therefore, the computational intensity is

$$\eta = \frac{16\lambda^3 L^2}{4\lambda^2 L(1 + L)} = 4 \frac{\lambda L}{1 + L} = \begin{cases} 2\lambda & (L = 1) \\ \frac{8}{3}\lambda & (L = 2) \end{cases}.$$

The computational intensity scales with λ , which is proportional to the square root of `VLEN`. Furthermore, the computational intensity is highest $L = 1$ or $L = 2$, which are the optimal choices for a given `VLEN`.

```

01  procedure  $\mu$ dgemm( $K, \alpha, \mathbf{A}, \mathbf{B}, \beta, \mathbf{C}$ )
02       $\begin{bmatrix} \text{v16} & \text{v17} & \text{v18} & \text{v19} \\ \text{v20} & \text{v21} & \text{v22} & \text{v23} \\ \text{v24} & \text{v25} & \text{v26} & \text{v27} \\ \text{v28} & \text{v29} & \text{v30} & \text{v31} \end{bmatrix} \leftarrow 0$ 
03      for( $k \leftarrow 0; k < K; k \leftarrow k + \lambda L$ )
04           $\text{kl} \leftarrow \min(K - k, \lambda L)$ 
05          mload( $\langle \text{fp64}, 4, \text{ml}, 1, \text{kl} \rangle(\text{v08}, \mathbf{A}(0, k), \text{lda})$ )
06          for( $x \leftarrow 0; x < L; x \leftarrow x + 1$ )
07              mload( $\langle \text{fp64}, 1, \text{kl}, 4, \text{nl} \rangle(\text{v12}, \mathbf{B}(k, 0), \text{ldb})$ )
08              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v08}, \text{v12}, \text{v16}, x)$ )
09              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v08}, \text{v13}, \text{v17}, x)$ )
10              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v08}, \text{v14}, \text{v18}, x)$ )
11              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v08}, \text{v15}, \text{v19}, x)$ )
12              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v09}, \text{v12}, \text{v20}, x)$ )
13              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v09}, \text{v13}, \text{v21}, x)$ )
14              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v09}, \text{v14}, \text{v22}, x)$ )
15              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v09}, \text{v15}, \text{v23}, x)$ )
16              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v10}, \text{v12}, \text{v24}, x)$ )
17              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v10}, \text{v13}, \text{v25}, x)$ )
18              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v10}, \text{v14}, \text{v26}, x)$ )
19              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v10}, \text{v15}, \text{v27}, x)$ )
20              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v11}, \text{v12}, \text{v28}, x)$ )
21              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v11}, \text{v13}, \text{v29}, x)$ )
22              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v11}, \text{v14}, \text{v30}, x)$ )
23              mgemmx( $\langle \text{fp64}, \times, + \rangle(\text{v11}, \text{v15}, \text{v31}, x)$ )
24               $\text{kl} \leftarrow \text{kl} - \lambda$ 
25          end
26      end
27       $(\text{v16}, \text{v17}, \dots, \text{v31}) \leftarrow \alpha \times (\text{v16}, \text{v17}, \dots, \text{v31})$ 
28       $\mathbf{C} \leftarrow \begin{bmatrix} \text{v16} & \text{v17} & \text{v18} & \text{v19} \\ \text{v20} & \text{v21} & \text{v22} & \text{v23} \\ \text{v24} & \text{v25} & \text{v26} & \text{v27} \\ \text{v28} & \text{v29} & \text{v30} & \text{v31} \end{bmatrix} + \beta \times \mathbf{C}$ 
29  end

```

$$\begin{pmatrix} \text{v08} \leftarrow \mathbf{A}(0 \times \lambda : +\lambda, k : +\lambda L) \\ \text{v09} \leftarrow \mathbf{A}(1 \times \lambda : +\lambda, k : +\lambda L) \\ \text{v10} \leftarrow \mathbf{A}(2 \times \lambda : +\lambda, k : +\lambda L) \\ \text{v11} \leftarrow \mathbf{A}(3 \times \lambda : +\lambda, k : +\lambda L) \end{pmatrix}$$

$$\begin{pmatrix} \text{v12} \leftarrow \mathbf{B}(k : +\lambda, 0 \times \lambda L : +\lambda L) \\ \text{v13} \leftarrow \mathbf{B}(k : +\lambda, 1 \times \lambda L : +\lambda L) \\ \text{v14} \leftarrow \mathbf{B}(k : +\lambda, 2 \times \lambda L : +\lambda L) \\ \text{v15} \leftarrow \mathbf{B}(k : +\lambda, 3 \times \lambda L : +\lambda L) \end{pmatrix}$$

Figure 2: The micro-kernel of the BLAS routine `dgemm` computes a $4\lambda \times 4\lambda L$ panel of the result matrix \mathbf{C} . The body of the innermost loop is executed once for each tile of the \mathbf{A} registers. In the body, it updates $16 \times L$ tiles of \mathbf{C} , each with a product of a tile of \mathbf{A} and a tile of \mathbf{B} . The code is agnostic to the values of λ and L .

	v12 [B ₀ ⁰ B ₀ ¹]	v13 [B ₀ ² B ₀ ³]	v14 [B ₀ ⁴ B ₀ ⁵]	v15 [B ₀ ⁶ B ₀ ⁷]
v08 [A ₀ ⁰ A ₀ ¹]	v16 [0 0]	v17 [0 0]	v18 [0 0]	v19 [0 0]
v09 [A ₁ ⁰ A ₁ ¹]	v20 [0 0]	v21 [0 0]	v22 [0 0]	v23 [0 0]
v10 [A ₂ ⁰ A ₂ ¹]	v24 [0 0]	v25 [0 0]	v26 [0 0]	v27 [0 0]
v11 [A ₃ ⁰ A ₃ ¹]	v28 [0 0]	v29 [0 0]	v30 [0 0]	v31 [0 0]

$$\begin{aligned}
\mathbf{v08} = & \begin{bmatrix} \begin{bmatrix} a_0^0 & a_0^1 & a_0^2 & a_0^3 \\ a_1^0 & a_1^1 & a_1^2 & a_1^3 \\ a_2^0 & a_2^1 & a_2^2 & a_2^3 \\ a_3^0 & a_3^1 & a_3^2 & a_3^3 \end{bmatrix} & \begin{bmatrix} a_0^4 & a_0^5 & a_0^6 & a_0^7 \\ a_1^4 & a_1^5 & a_1^6 & a_1^7 \\ a_2^4 & a_2^5 & a_2^6 & a_2^7 \\ a_3^4 & a_3^5 & a_3^6 & a_3^7 \end{bmatrix} \end{bmatrix} \quad \mathbf{v12} = \begin{bmatrix} \begin{bmatrix} b_0^0 & b_0^1 & b_0^2 & b_0^3 \\ b_1^0 & b_1^1 & b_1^2 & b_1^3 \\ b_2^0 & b_2^1 & b_2^2 & b_2^3 \\ b_3^0 & b_3^1 & b_3^2 & b_3^3 \end{bmatrix} & \begin{bmatrix} b_0^4 & b_0^5 & b_0^6 & b_0^7 \\ b_1^4 & b_1^5 & b_1^6 & b_1^7 \\ b_2^4 & b_2^5 & b_2^6 & b_2^7 \\ b_3^4 & b_3^5 & b_3^6 & b_3^7 \end{bmatrix} \end{bmatrix}
\end{aligned}$$

Figure 3: State of the computation at the beginning (line 08) of iteration $\langle k = 0, x = 0 \rangle$, with $\text{VLEN} = 2048$. Each of the registers contain 32 fp64 elements of either matrix **A** (v08-v11), **B** (v12-v15), or **C** (v16-v31). The 32 elements are organized as two 4×4 tiles in each register. The contents of v08 and v12 are shown explicitly. This particular configuration of $\lambda = 4$ and $L = 2$ is only one of the valid configurations for $\text{VLEN} = 2048$ shown in Table 1.

	v12 [B ₀ ⁰ B ₀ ¹]	v13 [B ₀ ² B ₀ ³]	v14 [B ₀ ⁴ B ₀ ⁵]	v15 [B ₀ ⁶ B ₀ ⁷]
v08 [A ₀ ⁰ A ₀ ¹]	v16 [A ₀ ⁰ B ₀ ⁰ A ₀ ⁰ B ₀ ¹]	v17 [A ₀ ⁰ B ₀ ² A ₀ ⁰ B ₀ ³]	v18 [A ₀ ⁰ B ₀ ⁴ A ₀ ⁰ B ₀ ⁵]	v19 [A ₀ ⁰ B ₀ ⁶ A ₀ ⁰ B ₀ ⁷]
v09 [A ₁ ⁰ A ₁ ¹]	v20 [A ₁ ⁰ B ₀ ⁰ A ₁ ⁰ B ₀ ¹]	v21 [A ₁ ⁰ B ₀ ² A ₁ ⁰ B ₀ ³]	v22 [A ₁ ⁰ B ₀ ⁴ A ₁ ⁰ B ₀ ⁵]	v23 [A ₁ ⁰ B ₀ ⁶ A ₁ ⁰ B ₀ ⁷]
v10 [A ₂ ⁰ A ₂ ¹]	v24 [A ₂ ⁰ B ₀ ⁰ A ₂ ⁰ B ₀ ¹]	v25 [A ₂ ⁰ B ₀ ² A ₂ ⁰ B ₀ ³]	v26 [A ₂ ⁰ B ₀ ⁴ A ₂ ⁰ B ₀ ⁵]	v27 [A ₂ ⁰ B ₀ ⁶ A ₂ ⁰ B ₀ ⁷]
v11 [A ₃ ⁰ A ₃ ¹]	v28 [A ₃ ⁰ B ₀ ⁰ A ₃ ⁰ B ₀ ¹]	v29 [A ₃ ⁰ B ₀ ² A ₃ ⁰ B ₀ ³]	v30 [A ₃ ⁰ B ₀ ⁴ A ₃ ⁰ B ₀ ⁵]	v31 [A ₃ ⁰ B ₀ ⁶ A ₃ ⁰ B ₀ ⁷]

Figure 4: State of the computation at the end (line 24) of iteration $\langle k = 0, x = 0 \rangle$.

	v12 [B_1^0 B_1^1]	v13 [B_1^2 B_1^3]	v14 [B_1^4 B_1^5]	v15 [B_1^6 B_1^7]
v08 [A_0^0 A_0^1]	v16 [$A_0^0 B_0^0$ $A_0^0 B_0^1$] + [$A_0^1 B_0^0$ $A_0^1 B_0^1$]	v17 [$A_0^0 B_0^2$ $A_0^0 B_0^3$] + [$A_0^1 B_0^2$ $A_0^1 B_0^3$]	v18 [$A_0^0 B_0^4$ $A_0^0 B_0^5$] + [$A_0^1 B_0^4$ $A_0^1 B_0^5$]	v19 [$A_0^0 B_0^6$ $A_0^0 B_0^7$] + [$A_0^1 B_0^6$ $A_0^1 B_0^7$]
v09 [A_1^0 A_1^1]	v20 [$A_1^0 B_0^0$ $A_1^0 B_0^1$] + [$A_1^1 B_0^0$ $A_1^1 B_0^1$]	v21 [$A_1^0 B_0^2$ $A_1^0 B_0^3$] + [$A_1^1 B_0^2$ $A_1^1 B_0^3$]	v22 [$A_1^0 B_0^4$ $A_1^0 B_0^5$] + [$A_1^1 B_0^4$ $A_1^1 B_0^5$]	v23 [$A_1^0 B_0^6$ $A_1^0 B_0^7$] + [$A_1^1 B_0^6$ $A_1^1 B_0^7$]
v10 [A_2^0 A_2^1]	v24 [$A_2^0 B_0^0$ $A_2^0 B_0^1$] + [$A_2^1 B_0^0$ $A_2^1 B_0^1$]	v25 [$A_2^0 B_0^2$ $A_2^0 B_0^3$] + [$A_2^1 B_0^2$ $A_2^1 B_0^3$]	v26 [$A_2^0 B_0^4$ $A_2^0 B_0^5$] + [$A_2^1 B_0^4$ $A_2^1 B_0^5$]	v27 [$A_2^0 B_0^6$ $A_2^0 B_0^7$] + [$A_2^1 B_0^6$ $A_2^1 B_0^7$]
v11 [A_3^0 A_3^1]	v28 [$A_3^0 B_0^0$ $A_3^0 B_0^1$] + [$A_3^1 B_0^0$ $A_3^1 B_0^1$]	v29 [$A_3^0 B_0^2$ $A_3^0 B_0^3$] + [$A_3^1 B_0^2$ $A_3^1 B_0^3$]	v30 [$A_3^0 B_0^4$ $A_3^0 B_0^5$] + [$A_3^1 B_0^4$ $A_3^1 B_0^5$]	v31 [$A_3^0 B_0^6$ $A_3^0 B_0^7$] + [$A_3^1 B_0^6$ $A_3^1 B_0^7$]

Figure 5: State of the computation at the end (line 24) of iteration $\langle k = 0, x = 1 \rangle$.

	v12 [B_2^0 B_2^1]	v13 [B_2^2 B_2^3]	v14 [B_2^4 B_2^5]	v15 [B_2^6 B_2^7]
v08 [A_0^2 A_0^3]	v16 $\begin{bmatrix} A_0^0 B_0^0 & A_0^0 B_0^1 \\ + & + \\ A_0^1 B_0^0 & A_0^1 B_0^1 \\ + & + \\ A_0^2 B_0^0 & A_0^2 B_0^1 \end{bmatrix}$	v17 $\begin{bmatrix} A_0^0 B_0^2 & A_0^0 B_0^3 \\ + & + \\ A_0^1 B_0^2 & A_0^1 B_0^3 \\ + & + \\ A_0^2 B_0^2 & A_0^2 B_0^3 \end{bmatrix}$	v18 $\begin{bmatrix} A_0^0 B_0^4 & A_0^0 B_0^5 \\ + & + \\ A_0^1 B_0^4 & A_0^1 B_0^5 \\ + & + \\ A_0^2 B_0^4 & A_0^2 B_0^5 \end{bmatrix}$	v19 $\begin{bmatrix} A_0^0 B_0^6 & A_0^0 B_0^7 \\ + & + \\ A_0^1 B_0^6 & A_0^1 B_0^7 \\ + & + \\ A_0^2 B_0^6 & A_0^2 B_0^7 \end{bmatrix}$
v09 [A_1^2 A_1^3]	v20 $\begin{bmatrix} A_1^0 B_0^0 & A_1^0 B_0^1 \\ + & + \\ A_1^1 B_0^0 & A_1^1 B_0^1 \\ + & + \\ A_1^2 B_0^0 & A_1^2 B_0^1 \end{bmatrix}$	v21 $\begin{bmatrix} A_1^0 B_0^2 & A_1^0 B_0^3 \\ + & + \\ A_1^1 B_0^2 & A_1^1 B_0^3 \\ + & + \\ A_1^2 B_0^2 & A_1^2 B_0^3 \end{bmatrix}$	v22 $\begin{bmatrix} A_1^0 B_0^4 & A_1^0 B_0^5 \\ + & + \\ A_1^1 B_0^4 & A_1^1 B_0^5 \\ + & + \\ A_1^2 B_0^4 & A_1^2 B_0^5 \end{bmatrix}$	v23 $\begin{bmatrix} A_1^0 B_0^6 & A_1^0 B_0^7 \\ + & + \\ A_1^1 B_0^6 & A_1^1 B_0^7 \\ + & + \\ A_1^2 B_0^6 & A_1^2 B_0^7 \end{bmatrix}$
v10 [A_2^2 A_2^3]	v24 $\begin{bmatrix} A_2^0 B_0^0 & A_2^0 B_0^1 \\ + & + \\ A_2^1 B_0^0 & A_2^1 B_0^1 \\ + & + \\ A_2^2 B_0^0 & A_2^2 B_0^1 \end{bmatrix}$	v25 $\begin{bmatrix} A_2^0 B_0^2 & A_2^0 B_0^3 \\ + & + \\ A_2^1 B_0^2 & A_2^1 B_0^3 \\ + & + \\ A_2^2 B_0^2 & A_2^2 B_0^3 \end{bmatrix}$	v26 $\begin{bmatrix} A_2^0 B_0^4 & A_2^0 B_0^5 \\ + & + \\ A_2^1 B_0^4 & A_2^1 B_0^5 \\ + & + \\ A_2^2 B_0^4 & A_2^2 B_0^5 \end{bmatrix}$	v27 $\begin{bmatrix} A_2^0 B_0^6 & A_2^0 B_0^7 \\ + & + \\ A_2^1 B_0^6 & A_2^1 B_0^7 \\ + & + \\ A_2^2 B_0^6 & A_2^2 B_0^7 \end{bmatrix}$
v11 [A_3^2 A_3^3]	v28 $\begin{bmatrix} A_3^0 B_0^0 & A_3^0 B_0^1 \\ + & + \\ A_3^1 B_0^0 & A_3^1 B_0^1 \\ + & + \\ A_3^2 B_0^0 & A_3^2 B_0^1 \end{bmatrix}$	v29 $\begin{bmatrix} A_3^0 B_0^2 & A_3^0 B_0^3 \\ + & + \\ A_3^1 B_0^2 & A_3^1 B_0^3 \\ + & + \\ A_3^2 B_0^2 & A_3^2 B_0^3 \end{bmatrix}$	v30 $\begin{bmatrix} A_3^0 B_0^4 & A_3^0 B_0^5 \\ + & + \\ A_3^1 B_0^4 & A_3^1 B_0^5 \\ + & + \\ A_3^2 B_0^4 & A_3^2 B_0^5 \end{bmatrix}$	v31 $\begin{bmatrix} A_3^0 B_0^6 & A_3^0 B_0^7 \\ + & + \\ A_3^1 B_0^6 & A_3^1 B_0^7 \\ + & + \\ A_3^2 B_0^6 & A_3^2 B_0^7 \end{bmatrix}$

Figure 6: State of the computation at the end (line 24) of iteration $\langle k = \lambda L, x = 0 \rangle$.

	$\begin{matrix} \text{v12} \\ [B_3^0 \ B_3^1] \end{matrix}$	$\begin{matrix} \text{v13} \\ [B_3^2 \ B_3^3] \end{matrix}$	$\begin{matrix} \text{v14} \\ [B_3^4 \ B_3^5] \end{matrix}$	$\begin{matrix} \text{v15} \\ [B_3^6 \ B_3^7] \end{matrix}$
$\begin{matrix} \text{v08} \\ [A_0^2 \ A_0^3] \end{matrix}$	$\begin{matrix} \text{v16} \\ \begin{bmatrix} A_0^0 B_0^0 & A_0^0 B_0^1 \\ + & + \\ A_0^1 B_0^0 & A_0^1 B_0^1 \\ + & + \\ A_0^2 B_0^0 & A_0^2 B_0^1 \\ + & + \\ A_0^3 B_0^0 & A_0^3 B_0^1 \end{bmatrix} \end{matrix}$	$\begin{matrix} \text{v17} \\ \begin{bmatrix} A_0^0 B_0^2 & A_0^0 B_0^3 \\ + & + \\ A_0^1 B_0^2 & A_0^1 B_0^3 \\ + & + \\ A_0^2 B_0^2 & A_0^2 B_0^3 \\ + & + \\ A_0^3 B_0^2 & A_0^3 B_0^3 \end{bmatrix} \end{matrix}$	$\begin{matrix} \text{v18} \\ \begin{bmatrix} A_0^0 B_0^4 & A_0^0 B_0^5 \\ + & + \\ A_0^1 B_0^4 & A_0^1 B_0^5 \\ + & + \\ A_0^2 B_0^4 & A_0^2 B_0^5 \\ + & + \\ A_0^3 B_0^4 & A_0^3 B_0^5 \end{bmatrix} \end{matrix}$	$\begin{matrix} \text{v19} \\ \begin{bmatrix} A_0^0 B_0^6 & A_0^0 B_0^7 \\ + & + \\ A_0^1 B_0^6 & A_0^1 B_0^7 \\ + & + \\ A_0^2 B_0^6 & A_0^2 B_0^7 \\ + & + \\ A_0^3 B_0^6 & A_0^3 B_0^7 \end{bmatrix} \end{matrix}$
$\begin{matrix} \text{v09} \\ [A_1^2 \ A_1^3] \end{matrix}$	$\begin{matrix} \text{v20} \\ \begin{bmatrix} A_1^0 B_1^0 & A_1^0 B_1^1 \\ + & + \\ A_1^1 B_1^0 & A_1^1 B_1^1 \\ + & + \\ A_1^2 B_1^0 & A_1^2 B_1^1 \\ + & + \\ A_1^3 B_1^0 & A_1^3 B_1^1 \end{bmatrix} \end{matrix}$	$\begin{matrix} \text{v21} \\ \begin{bmatrix} A_1^0 B_1^2 & A_1^0 B_1^3 \\ + & + \\ A_1^1 B_1^2 & A_1^1 B_1^3 \\ + & + \\ A_1^2 B_1^2 & A_1^2 B_1^3 \\ + & + \\ A_1^3 B_1^2 & A_1^3 B_1^3 \end{bmatrix} \end{matrix}$	$\begin{matrix} \text{v22} \\ \begin{bmatrix} A_1^0 B_1^4 & A_1^0 B_1^5 \\ + & + \\ A_1^1 B_1^4 & A_1^1 B_1^5 \\ + & + \\ A_1^2 B_1^4 & A_1^2 B_1^5 \\ + & + \\ A_1^3 B_1^4 & A_1^3 B_1^5 \end{bmatrix} \end{matrix}$	$\begin{matrix} \text{v23} \\ \begin{bmatrix} A_1^0 B_1^6 & A_1^0 B_1^7 \\ + & + \\ A_1^1 B_1^6 & A_1^1 B_1^7 \\ + & + \\ A_1^2 B_1^6 & A_1^2 B_1^7 \\ + & + \\ A_1^3 B_1^6 & A_1^3 B_1^7 \end{bmatrix} \end{matrix}$
$\begin{matrix} \text{v10} \\ [A_2^2 \ A_2^3] \end{matrix}$	$\begin{matrix} \text{v24} \\ \begin{bmatrix} A_2^0 B_2^0 & A_2^0 B_2^1 \\ + & + \\ A_2^1 B_2^0 & A_2^1 B_2^1 \\ + & + \\ A_2^2 B_2^0 & A_2^2 B_2^1 \\ + & + \\ A_2^3 B_2^0 & A_2^3 B_2^1 \end{bmatrix} \end{matrix}$	$\begin{matrix} \text{v25} \\ \begin{bmatrix} A_2^0 B_2^2 & A_2^0 B_2^3 \\ + & + \\ A_2^1 B_2^2 & A_2^1 B_2^3 \\ + & + \\ A_2^2 B_2^2 & A_2^2 B_2^3 \\ + & + \\ A_2^3 B_2^2 & A_2^3 B_2^3 \end{bmatrix} \end{matrix}$	$\begin{matrix} \text{v26} \\ \begin{bmatrix} A_2^0 B_2^4 & A_2^0 B_2^5 \\ + & + \\ A_2^1 B_2^4 & A_2^1 B_2^5 \\ + & + \\ A_2^2 B_2^4 & A_2^2 B_2^5 \\ + & + \\ A_2^3 B_2^4 & A_2^3 B_2^5 \end{bmatrix} \end{matrix}$	$\begin{matrix} \text{v27} \\ \begin{bmatrix} A_2^0 B_2^6 & A_2^0 B_2^7 \\ + & + \\ A_2^1 B_2^6 & A_2^1 B_2^7 \\ + & + \\ A_2^2 B_2^6 & A_2^2 B_2^7 \\ + & + \\ A_2^3 B_2^6 & A_2^3 B_2^7 \end{bmatrix} \end{matrix}$
$\begin{matrix} \text{v11} \\ [A_3^2 \ A_3^3] \end{matrix}$	$\begin{matrix} \text{v28} \\ \begin{bmatrix} A_3^0 B_3^0 & A_3^0 B_3^1 \\ + & + \\ A_3^1 B_3^0 & A_3^1 B_3^1 \\ + & + \\ A_3^2 B_3^0 & A_3^2 B_3^1 \\ + & + \\ A_3^3 B_3^0 & A_3^3 B_3^1 \end{bmatrix} \end{matrix}$	$\begin{matrix} \text{v29} \\ \begin{bmatrix} A_3^0 B_3^2 & A_3^0 B_3^3 \\ + & + \\ A_3^1 B_3^2 & A_3^1 B_3^3 \\ + & + \\ A_3^2 B_3^2 & A_3^2 B_3^3 \\ + & + \\ A_3^3 B_3^2 & A_3^3 B_3^3 \end{bmatrix} \end{matrix}$	$\begin{matrix} \text{v30} \\ \begin{bmatrix} A_3^0 B_3^4 & A_3^0 B_3^5 \\ + & + \\ A_3^1 B_3^4 & A_3^1 B_3^5 \\ + & + \\ A_3^2 B_3^4 & A_3^2 B_3^5 \\ + & + \\ A_3^3 B_3^4 & A_3^3 B_3^5 \end{bmatrix} \end{matrix}$	$\begin{matrix} \text{v31} \\ \begin{bmatrix} A_3^0 B_3^6 & A_3^0 B_3^7 \\ + & + \\ A_3^1 B_3^6 & A_3^1 B_3^7 \\ + & + \\ A_3^2 B_3^6 & A_3^2 B_3^7 \\ + & + \\ A_3^3 B_3^6 & A_3^3 B_3^7 \end{bmatrix} \end{matrix}$

Figure 7: State of the computation at the end (line 24) of iteration $\langle k = \lambda L, x = 1 \rangle$.

5 Trivial extension to mixed data types

Let $T(\mathbf{A})$, $T(\mathbf{B})$, and $T(\mathbf{C})$ denote the data types of matrices \mathbf{A} , \mathbf{B} , and \mathbf{C} , respectively. If $\text{sizeof}(T(\mathbf{A})) = \text{sizeof}(T(\mathbf{B})) = \text{sizeof}(T(\mathbf{C}))/n$, we can treat all three matrices as having $\text{MEW} = \text{sizeof}(T(\mathbf{C}))$ and define the elements of \mathbf{A} and \mathbf{B} to be n -vectors of their natural data types ($T(\mathbf{A})$ and $T(\mathbf{B})$, respectively). The \otimes operator becomes the dot-product of two n -vectors of type $T(\mathbf{A})$ and $T(\mathbf{B})$. The individual dot-products are summed up with the \oplus . With these definitions it is trivial to extend the common-type variant to treat mixed types for the three matrices, subject to the constraints outlined above.

Figure 8 illustrates the contents of sample 4×4 tiles of matrices \mathbf{A} , \mathbf{B} , and \mathbf{C} for $T(\mathbf{A}) = T(\mathbf{B}) = \text{fp16}$ and $T(\mathbf{C}) = \text{fp64}$. In this case $\text{MEW} = 64$ and all tiles take 1024 bits of storage. The elements of \mathbf{C} are 64-bit floating-point numbers, while the elements of \mathbf{A} and \mathbf{B} are 4-vectors of 16-bit floating-point numbers. The matrix multiplication is defined over a pair $\langle \otimes, \oplus \rangle$ of operators. The \otimes operator is the dot-product of two 4-vectors of half-precision floating-point numbers, producing a double-precision result which is then accumulated with the \oplus operator.

$$A_0^0 = \begin{bmatrix} [a_0^0(0) & a_0^0(1) & a_0^0(2) & a_0^0(3)] & [a_0^1(0) & a_0^1(1) & a_0^1(2) & a_0^1(3)] & [a_0^2(0) & a_0^2(1) & a_0^2(2) & a_0^2(3)] & [a_0^3(0) & a_0^3(1) & a_0^3(2) & a_0^3(3)] \\ [a_1^0(0) & a_1^0(1) & a_1^0(2) & a_1^0(3)] & [a_1^1(0) & a_1^1(1) & a_1^1(2) & a_1^1(3)] & [a_1^2(0) & a_1^2(1) & a_1^2(2) & a_1^2(3)] & [a_1^3(0) & a_1^3(1) & a_1^3(2) & a_1^3(3)] \\ [a_2^0(0) & a_2^0(1) & a_2^0(2) & a_2^0(3)] & [a_2^1(0) & a_2^1(1) & a_2^1(2) & a_2^1(3)] & [a_2^2(0) & a_2^2(1) & a_2^2(2) & a_2^2(3)] & [a_2^3(0) & a_2^3(1) & a_2^3(2) & a_2^3(3)] \\ [a_3^0(0) & a_3^0(1) & a_3^0(2) & a_3^0(3)] & [a_3^1(0) & a_3^1(1) & a_3^1(2) & a_3^1(3)] & [a_3^2(0) & a_3^2(1) & a_3^2(2) & a_3^2(3)] & [a_3^3(0) & a_3^3(1) & a_3^3(2) & a_3^3(3)] \end{bmatrix}$$

$$B_0^0 = \begin{bmatrix} [b_0^0(0) & b_0^0(1) & b_0^0(2) & b_0^0(3)] & [b_0^1(0) & b_0^1(1) & b_0^1(2) & b_0^1(3)] & [b_0^2(0) & b_0^2(1) & b_0^2(2) & b_0^2(3)] & [b_0^3(0) & b_0^3(1) & b_0^3(2) & b_0^3(3)] \\ [b_1^0(0) & b_1^0(1) & b_1^0(2) & b_1^0(3)] & [b_1^1(0) & b_1^1(1) & b_1^1(2) & b_1^1(3)] & [b_1^2(0) & b_1^2(1) & b_1^2(2) & b_1^2(3)] & [b_1^3(0) & b_1^3(1) & b_1^3(2) & b_1^3(3)] \\ [b_2^0(0) & b_2^0(1) & b_2^0(2) & b_2^0(3)] & [b_2^1(0) & b_2^1(1) & b_2^1(2) & b_2^1(3)] & [b_2^2(0) & b_2^2(1) & b_2^2(2) & b_2^2(3)] & [b_2^3(0) & b_2^3(1) & b_2^3(2) & b_2^3(3)] \\ [b_3^0(0) & b_3^0(1) & b_3^0(2) & b_3^0(3)] & [b_3^1(0) & b_3^1(1) & b_3^1(2) & b_3^1(3)] & [b_3^2(0) & b_3^2(1) & b_3^2(2) & b_3^2(3)] & [b_3^3(0) & b_3^3(1) & b_3^3(2) & b_3^3(3)] \end{bmatrix}$$

$$C_0^0 = \begin{bmatrix} c_0^0 & c_0^1 & c_0^2 & c_0^3 \\ c_1^0 & c_1^1 & c_1^2 & c_1^3 \\ c_2^0 & c_2^1 & c_2^2 & c_2^3 \\ c_3^0 & c_3^1 & c_3^2 & c_3^3 \end{bmatrix}$$

Figure 8: Contents of sample tiles of matrices \mathbf{A} , \mathbf{B} , and \mathbf{C} for $\lambda = 4$, $T(\mathbf{A}) = T(\mathbf{B}) = \text{fp16}$, and $T(\mathbf{C}) = \text{fp64}$. In this case, $\text{MEW} = 64$, the \otimes operator is the dot-product of 4-vectors of fp16 elements producing a fp64 result, and the \oplus operator is a fp64 addition.

6 General mixed data types

For the general case, let $\text{sizeof}(T(\mathbf{A})) \leq \text{sizeof}(T(\mathbf{B})) \leq \text{sizeof}(T(\mathbf{C}))$. We choose a shape $\lambda \times \lambda$ for the tiles of each matrix. Let $L(\mathbf{C})$ be the number of tiles of matrix \mathbf{C} in a vector register. Let $n(\mathbf{A}) = \frac{\text{sizeof}(T(\mathbf{C}))}{\text{sizeof}(T(\mathbf{A}))}$ and $n(\mathbf{B}) = \frac{\text{sizeof}(T(\mathbf{C}))}{\text{sizeof}(T(\mathbf{B}))}$. Then, $L(\mathbf{A}) = n(\mathbf{A})L(\mathbf{C})$ and $L(\mathbf{B}) = n(\mathbf{B})L(\mathbf{C})$ are the number of tiles of matrices \mathbf{A} and \mathbf{B} , respectively, in a vector register.

The `mgemmx` instruction is modified to multiply $n(\mathbf{B})$ tiles of \mathbf{A} by $n(\mathbf{B}) \times L(\mathbf{C})$ tiles of \mathbf{B} , producing $L(\mathbf{C})$ result tiles that are accumulated into \mathbf{C} . The innermost loop (line 06) of the `μ dgemm` kernel is modified to increment x by $n(\mathbf{B})$, with an upper bound of $L(\mathbf{A})$.

	$\mathbf{v12}$ $\begin{bmatrix} B_0^0 \\ B_1^0 \end{bmatrix}$	$\mathbf{v13}$ $\begin{bmatrix} B_0^1 \\ B_1^1 \end{bmatrix}$	$\mathbf{v14}$ $\begin{bmatrix} B_0^2 \\ B_1^2 \end{bmatrix}$	$\mathbf{v15}$ $\begin{bmatrix} B_0^3 \\ B_1^3 \end{bmatrix}$
$\mathbf{v08}$ $\begin{bmatrix} A_0^0 & A_0^1 & A_0^2 & A_0^3 \end{bmatrix}$	$\mathbf{v16}$ $\begin{bmatrix} 0 \end{bmatrix}$	$\mathbf{v17}$ $\begin{bmatrix} 0 \end{bmatrix}$	$\mathbf{v18}$ $\begin{bmatrix} 0 \end{bmatrix}$	$\mathbf{v19}$ $\begin{bmatrix} 0 \end{bmatrix}$
$\mathbf{v09}$ $\begin{bmatrix} A_1^0 & A_1^1 & A_1^2 & A_1^3 \end{bmatrix}$	$\mathbf{v20}$ $\begin{bmatrix} 0 \end{bmatrix}$	$\mathbf{v21}$ $\begin{bmatrix} 0 \end{bmatrix}$	$\mathbf{v22}$ $\begin{bmatrix} 0 \end{bmatrix}$	$\mathbf{v23}$ $\begin{bmatrix} 0 \end{bmatrix}$
$\mathbf{v10}$ $\begin{bmatrix} A_2^0 & A_2^1 & A_2^2 & A_2^3 \end{bmatrix}$	$\mathbf{v24}$ $\begin{bmatrix} 0 \end{bmatrix}$	$\mathbf{v25}$ $\begin{bmatrix} 0 \end{bmatrix}$	$\mathbf{v26}$ $\begin{bmatrix} 0 \end{bmatrix}$	$\mathbf{v27}$ $\begin{bmatrix} 0 \end{bmatrix}$
$\mathbf{v11}$ $\begin{bmatrix} A_3^0 & A_3^1 & A_3^2 & A_3^3 \end{bmatrix}$	$\mathbf{v28}$ $\begin{bmatrix} 0 \end{bmatrix}$	$\mathbf{v29}$ $\begin{bmatrix} 0 \end{bmatrix}$	$\mathbf{v30}$ $\begin{bmatrix} 0 \end{bmatrix}$	$\mathbf{v31}$ $\begin{bmatrix} 0 \end{bmatrix}$

$$\mathbf{v08} = \left[\begin{array}{c} \begin{bmatrix} a_0^0 & a_0^1 & a_0^2 & a_0^3 \\ a_1^0 & a_1^1 & a_1^2 & a_1^3 \\ a_2^0 & a_2^1 & a_2^2 & a_2^3 \\ a_3^0 & a_3^1 & a_3^2 & a_3^3 \end{bmatrix} \\ \begin{bmatrix} a_0^4 & a_0^5 & a_0^6 & a_0^7 \\ a_1^4 & a_1^5 & a_1^6 & a_1^7 \\ a_2^4 & a_2^5 & a_2^6 & a_2^7 \\ a_3^4 & a_3^5 & a_3^6 & a_3^7 \end{bmatrix} \\ \begin{bmatrix} a_0^8 & a_0^9 & a_0^{10} & a_0^{11} \\ a_1^8 & a_1^9 & a_1^{10} & a_1^{11} \\ a_2^8 & a_2^9 & a_2^{10} & a_2^{11} \\ a_3^8 & a_3^9 & a_3^{10} & a_3^{11} \end{bmatrix} \\ \begin{bmatrix} a_0^{12} & a_0^{13} & a_0^{14} & a_0^{15} \\ a_1^{12} & a_1^{13} & a_1^{14} & a_1^{15} \\ a_2^{12} & a_2^{13} & a_2^{14} & a_2^{15} \\ a_3^{12} & a_3^{13} & a_3^{14} & a_3^{15} \end{bmatrix} \end{array} \right]$$

$$\mathbf{v12} = \left[\begin{array}{c} \begin{bmatrix} b_0^0 & b_0^1 & b_0^2 & b_0^3 \\ b_1^0 & b_1^1 & b_1^2 & b_1^3 \\ b_2^0 & b_2^1 & b_2^2 & b_2^3 \\ b_3^0 & b_3^1 & b_3^2 & b_3^3 \end{bmatrix} \\ \begin{bmatrix} b_4^0 & b_4^1 & b_4^2 & b_4^3 \\ b_5^0 & b_5^1 & b_5^2 & b_5^3 \\ b_6^0 & b_6^1 & b_6^2 & b_6^3 \\ b_7^0 & b_7^1 & b_7^2 & b_7^3 \end{bmatrix} \end{array} \right]$$

Figure 9: State of the computation at the beginning of iteration $\langle k = 0, x = 0 \rangle$, with $\text{VLEN} = 1024$, $T(\mathbf{C}) = \text{fp64}$, $T(\mathbf{B}) = \text{fp32}$, and $T(\mathbf{A}) = \text{fp16}$. This particular configuration of $\lambda = 4$, $L(\mathbf{C}) = 1$, $L(\mathbf{B}) = 2$, and $L(\mathbf{A}) = 4$ is only one of the valid configurations for $\text{VLEN} = 1024$ shown in Table 1.

	v12 $\begin{bmatrix} B_0^0 \\ B_1^0 \end{bmatrix}$	v13 $\begin{bmatrix} B_0^1 \\ B_1^1 \end{bmatrix}$	v14 $\begin{bmatrix} B_0^2 \\ B_1^2 \end{bmatrix}$	v15 $\begin{bmatrix} B_0^3 \\ B_1^3 \end{bmatrix}$
v08 $\begin{bmatrix} A_0^0 & A_0^1 & A_0^2 & A_0^3 \end{bmatrix}$	v16 $\begin{bmatrix} A_0^0 B_0^0 + A_0^1 B_1^0 \end{bmatrix}$	v17 $\begin{bmatrix} A_0^0 B_0^1 + A_0^1 B_1^1 \end{bmatrix}$	v18 $\begin{bmatrix} A_0^0 B_0^2 + A_0^1 B_1^2 \end{bmatrix}$	v19 $\begin{bmatrix} A_0^0 B_0^3 + A_0^1 B_1^3 \end{bmatrix}$
v09 $\begin{bmatrix} A_1^0 & A_1^1 & A_1^2 & A_1^3 \end{bmatrix}$	v20 $\begin{bmatrix} A_1^0 B_0^0 + A_1^1 B_1^0 \end{bmatrix}$	v21 $\begin{bmatrix} A_1^0 B_0^1 + A_1^1 B_1^1 \end{bmatrix}$	v22 $\begin{bmatrix} A_1^0 B_0^2 + A_1^1 B_1^2 \end{bmatrix}$	v23 $\begin{bmatrix} A_1^0 B_0^3 + A_1^1 B_1^3 \end{bmatrix}$
v10 $\begin{bmatrix} A_2^0 & A_2^1 & A_2^2 & A_2^3 \end{bmatrix}$	v24 $\begin{bmatrix} A_2^0 B_0^0 + A_2^1 B_1^0 \end{bmatrix}$	v25 $\begin{bmatrix} A_2^0 B_0^1 + A_2^1 B_1^1 \end{bmatrix}$	v26 $\begin{bmatrix} A_2^0 B_0^2 + A_2^1 B_1^2 \end{bmatrix}$	v27 $\begin{bmatrix} A_2^0 B_0^3 + A_2^1 B_1^3 \end{bmatrix}$
v11 $\begin{bmatrix} A_3^0 & A_3^1 & A_3^2 & A_3^3 \end{bmatrix}$	v28 $\begin{bmatrix} A_3^0 B_0^0 + A_3^1 B_1^0 \end{bmatrix}$	v29 $\begin{bmatrix} A_3^0 B_0^1 + A_3^1 B_1^1 \end{bmatrix}$	v30 $\begin{bmatrix} A_3^0 B_0^2 + A_3^1 B_1^2 \end{bmatrix}$	v31 $\begin{bmatrix} A_3^0 B_0^3 + A_3^1 B_1^3 \end{bmatrix}$

Figure 10: State of the computation at the end (line 24) of iteration $\langle k = 0, x = 0 \rangle$.

	v12 $\begin{bmatrix} B_2^0 \\ B_3^0 \end{bmatrix}$	v13 $\begin{bmatrix} B_2^1 \\ B_3^1 \end{bmatrix}$	v14 $\begin{bmatrix} B_2^2 \\ B_3^2 \end{bmatrix}$	v15 $\begin{bmatrix} B_2^3 \\ B_3^3 \end{bmatrix}$
v08 $\begin{bmatrix} A_0^0 & A_0^1 & A_0^2 & A_0^3 \end{bmatrix}$	v16 $\begin{bmatrix} A_0^0 B_0^0 + A_0^1 B_1^0 \\ + \\ A_0^2 B_2^0 + A_0^3 B_3^0 \end{bmatrix}$	v17 $\begin{bmatrix} A_0^0 B_0^1 + A_0^1 B_1^1 \\ + \\ A_0^2 B_2^1 + A_0^3 B_3^1 \end{bmatrix}$	v18 $\begin{bmatrix} A_0^0 B_0^2 + A_0^1 B_1^2 \\ + \\ A_0^2 B_2^2 + A_0^3 B_3^2 \end{bmatrix}$	v19 $\begin{bmatrix} A_0^0 B_0^3 + A_0^1 B_1^3 \\ + \\ A_0^2 B_2^3 + A_0^3 B_3^3 \end{bmatrix}$
v09 $\begin{bmatrix} A_1^0 & A_1^1 & A_1^2 & A_1^3 \end{bmatrix}$	v20 $\begin{bmatrix} A_1^0 B_0^0 + A_1^1 B_1^0 \\ + \\ A_1^2 B_2^0 + A_1^3 B_3^0 \end{bmatrix}$	v21 $\begin{bmatrix} A_1^0 B_0^1 + A_1^1 B_1^1 \\ + \\ A_1^2 B_2^1 + A_1^3 B_3^1 \end{bmatrix}$	v22 $\begin{bmatrix} A_1^0 B_0^2 + A_1^1 B_1^2 \\ + \\ A_1^2 B_2^2 + A_1^3 B_3^2 \end{bmatrix}$	v23 $\begin{bmatrix} A_1^0 B_0^3 + A_1^1 B_1^3 \\ + \\ A_1^2 B_2^3 + A_1^3 B_3^3 \end{bmatrix}$
v10 $\begin{bmatrix} A_2^0 & A_2^1 & A_2^2 & A_2^3 \end{bmatrix}$	v24 $\begin{bmatrix} A_2^0 B_0^0 + A_2^1 B_1^0 \\ + \\ A_2^2 B_2^0 + A_2^3 B_3^0 \end{bmatrix}$	v25 $\begin{bmatrix} A_2^0 B_0^1 + A_2^1 B_1^1 \\ + \\ A_2^2 B_2^1 + A_2^3 B_3^1 \end{bmatrix}$	v26 $\begin{bmatrix} A_2^0 B_0^2 + A_2^1 B_1^2 \\ + \\ A_2^2 B_2^2 + A_2^3 B_3^2 \end{bmatrix}$	v27 $\begin{bmatrix} A_2^0 B_0^3 + A_2^1 B_1^3 \\ + \\ A_2^2 B_2^3 + A_2^3 B_3^3 \end{bmatrix}$
v11 $\begin{bmatrix} A_3^0 & A_3^1 & A_3^2 & A_3^3 \end{bmatrix}$	v28 $\begin{bmatrix} A_3^0 B_0^0 + A_3^1 B_1^0 \\ + \\ A_3^2 B_2^0 + A_3^3 B_3^0 \end{bmatrix}$	v29 $\begin{bmatrix} A_3^0 B_0^1 + A_3^1 B_1^1 \\ + \\ A_3^2 B_2^1 + A_3^3 B_3^1 \end{bmatrix}$	v30 $\begin{bmatrix} A_3^0 B_0^2 + A_3^1 B_1^2 \\ + \\ A_3^2 B_2^2 + A_3^3 B_3^2 \end{bmatrix}$	v31 $\begin{bmatrix} A_3^0 B_0^3 + A_3^1 B_1^3 \\ + \\ A_3^2 B_2^3 + A_3^3 B_3^3 \end{bmatrix}$

Figure 11: State of the computation at the end (line 24) of iteration $\langle k = 0, x = 2 \rangle$.
