

# Tiling/Blocking MHA, GEMM, LN, CONV

P.D.

## Abstract

How a conversation with Zhongxlive-latex-base, texlive-fonts-recommendedand texlive-latex-recommendedyi Lin about Multi-head attention became the starting point for a stand alone tiling tool for the search, description and computation of blocked algorithms for AIE overlays.

We start with a formal description of MHA and we jot down several ideas how to describe, code, and interpret tiling strategies. This is our first attempt to describe formally and optimally the tiling process. We started with pure tiling as suggestions to an FPGA IP (XDNN), with layer fusion and explicit double buffering by software pipelining (DPUV3INT8). We introduce explicit code and compute parallelism by symmetric computations for sparse convolutions (Weight block sparsity). Here, we finally explore all valid solutions and optimum by any custom function with the ability to represent complex computation using double buffering at all level of the memory hierarchy without being hold back by how to implement each instructions (but having done so in the past).

## Index Terms

AI, FPGA, Performance, LLM BS and Tools

## I. MHA, ITS NOTATIONS AND AN INTRODUCTION

This is a HEAD

$$\mathbf{R} = \mathcal{S}_{max}(\mathbf{Q} * \mathbf{K}^t) * \mathbf{V} \quad (1)$$

Where  $\mathbf{K}, \mathbf{Q}, \mathbf{V} \in \mathbb{R}^{L \times d}$ , where  $d \in \{32, 64\}$  is small and  $L \in \{4096\}$  is large

$$\mathcal{S}_{max}(\mathbf{M}) = \mathbf{S}, \text{ with } \mathbf{S}, \mathbf{M} \in \mathbb{R}^{m \times n} \quad (2)$$

$$\mathcal{S}_{max}(\mathbf{M})_{i,j} = \frac{\exp \mathbf{M}_{i,j}}{\sum_j \exp \mathbf{M}_{i,j}} \quad (3)$$

Where the denominator is a sum by row and the result is literally a vector normalizing the row of hte numerator. Another way is to consider it a unitary matrix multiplying on the left

$$\Gamma = \frac{1}{\sum_j \exp \mathbf{M}_{i,j}} = \left[ \frac{1}{\sum_j \exp \mathbf{M}_{0,j}}, \dots, \frac{1}{\sum_j \exp \mathbf{M}_{m,j}} \right] * \mathbf{I} \quad (4)$$

$$\mathcal{S}_{max}(\mathbf{M}) = \Gamma \mathbf{M} = \mathbf{S} \quad (5)$$

$$\mathbf{Q} * \mathbf{K}^t = \mathbf{T} \in \mathbb{R}^{L \times L} \quad (6)$$

The submatrix of  $\mathbf{K}^t$  are actually transposed, so I do not need to carry the transpose sign everywhere. Assume that  $\mathbf{K}_i, \mathbf{Q}_j \in \mathbb{R}^{d \times d}$

$$\begin{pmatrix} \mathbf{Q}_0 \\ \mathbf{Q}_1 \\ \dots \\ \mathbf{Q}_n \end{pmatrix} * \begin{pmatrix} \mathbf{K}_0 & \mathbf{K}_1 & \dots & \mathbf{K}_n \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_0 \mathbf{K}_0 & \mathbf{Q}_0 \mathbf{K}_1 & \dots & \mathbf{Q}_0 \mathbf{K}_n \\ \mathbf{Q}_1 \mathbf{K}_0 & \mathbf{Q}_1 \mathbf{K}_1 & \dots & \mathbf{Q}_1 \mathbf{K}_n \\ \dots & \dots & \dots & \dots \\ \mathbf{Q}_n \mathbf{K}_0 & \mathbf{Q}_n \mathbf{K}_1 & \dots & \mathbf{Q}_n \mathbf{K}_n \end{pmatrix}$$

$$\mathbf{R}_0 = \mathcal{S}_{max}((\mathbf{Q}_0 \mathbf{K}_0 \quad \mathbf{Q}_0 \mathbf{K}_1 \quad \dots \quad \mathbf{Q}_0 \mathbf{K}_n)) \begin{pmatrix} \mathbf{V}_0 \\ \mathbf{V}_1 \\ \dots \\ \mathbf{V}_n \end{pmatrix}$$

To be precise,

$$\mathbf{R}_0 = \frac{1}{\sum_j \exp \mathbf{Q}_0 \mathbf{K}_j} (\exp \mathbf{Q}_0 \mathbf{K}_0 \quad \dots \quad \exp \mathbf{Q}_0 \mathbf{K}_n) \begin{pmatrix} \mathbf{V}_0 \\ \mathbf{V}_1 \\ \dots \\ \mathbf{V}_n \end{pmatrix} \quad (7)$$

$$\mathbf{R}_0 = \frac{\sum_j (\exp \mathbf{Q}_0 \mathbf{K}_j) \mathbf{V}_j}{\sum_j \exp \mathbf{Q}_0 \mathbf{K}_j} = \Gamma_0 \sum_j (\exp \mathbf{Q}_0 \mathbf{K}_j) \mathbf{V}_j \quad (8)$$

The sum at the denominator is by row, the result is a vector and it can be represented also as a left diagonal matrix multiplication using Equation 4. The sum at the numerator is a matrix addition (or a left matrix multiplication and thanks to associative). Thus, Equations 7 and 8 are proper and equivalent. We can apply to the remaining  $\mathbf{Q}_i$  and  $\mathbf{R}_i$ .

## II. MHA, BLOCKED, ITERATIVE METHOD BY COLUMNS, AND CACHING

Consider only the contribution by  $Q_m$  and take  $N_{m,0} = D_{m,0} = 0$ ,  $M_{m,0} = \mathbf{0}$ , assume we have done the computation until step  $t-1$  and compute  $t$  to  $t < n$  then

$$\begin{aligned} M^1 &= \max_{\text{row}}(Q_m K_t, M_{m,t-1}) \\ E &= \exp(Q_m K_t - M^1) \\ S &= \exp(M^1 - M_{m,t-1}) \\ M_{m,t} &= M^1 \end{aligned}$$

$$N_{m,t} = \left( \frac{N_{m,t-1}}{S} + E V_t \right) \quad (9)$$

$$D_{m,t} = \frac{D_{m,t-1}}{S + E} \quad (10)$$

Finally, for  $t = n$  and for every element of  $Q_m$  we can repeat the process and compute the final result

$$\forall i \in [0, n], R_i = \frac{N_{i,n}}{D_{i,n}} \quad (11)$$

### A. Caching

Assume, we have  $K, Q, V \in \mathbb{R}^{L \times d}$ , we split them into partitions and we have  $n$  of them. We compute the result  $R$  as above Equation 11 and we remember the result  $R$ , the max vector  $M = \{M_i\}$ , the scaling diagonal  $D = \{D_i\}$ , and the numerators  $\{N_i\}$ . The last two we remember just the last *column computation*.

Now, consider the case that we read  $K_{n+1}, Q_{n+1}, V_{n+1}$  and we want to compute the update version  $R$ . We have two observations: first, we need to compute a fresh  $R_{n+1}$  from scratch:

$$R_{n+1}, M_{n+1}, N_{n+1}, D_{n+1} = \mathcal{S}_{\max}(Q_{n+1} * K^t) * V. \quad (12)$$

Second, the update of  $R_i$  with  $i < n+1$  is exactly the same Equation 9, 10, and 11.

### B. Range of $\exp(Q_0 K_j)$

The function  $e^x$  has the property that  $\frac{\partial e^x}{\partial x} = e^x$ . The function  $e^x$  is always positive, is always increasing, and  $e^0 = 1$  and also the derivative. For  $x > 0$  and with numerical representation that could be limited in range, we want to stay on the left side of  $x = 0$  or at a very least in a bound region to avoid saturation and overflow. Clearly, for  $x > 0$  we have  $\frac{\partial e^x}{\partial x} > 1$  and it really highlights small variation of  $x$ .

If we keep the range of  $Q_0 K_j$  being smaller than one, then the  $\exp()$  is bound to  $e$ . Due to the nature of the exponential function, the range will be positive and increasing (well exponentially). For example, assume we compute the product  $QK$  and we determine the maximum by row. This is a column vector  $M = \max_{\text{row}} QK$  then we can compute the final result.

$$R_0^t = \frac{\sum_j (\exp(Q_0 K_j - M)) V_j}{\sum_j \exp(Q_0 K_j - M)} = R_0 \frac{e^M}{e^M} = R_0 \quad (13)$$

The exponential range is under control, the summation at the denominator is increasing but bound to the number of additions. The denominator can be applied at the end of the computation but the normalization of the range as to be applied as soon as possible in such a way to have undesired overflow. In Section II, we combine everything into an iterative method that can be split into independent computations.

### C. The identity element in $\mathcal{S}_{\max}(M)$

In the ring based on the operation addition  $+$ , there is a identity element 0 (zero) such as  $\forall a, a + 0 = 0 + a = a$ . For multiplication  $*$  the identity element is 1 (one). If the operands are matrices, we can extend the zero and one elements to zero matrices and to diagonal one matrices.

Take  $\frac{\exp M_{i,j}}{\sum_j \exp M_{i,j}}$ , if for any reason we need to pad or extend  $M$  to  $N$ , what can we do so that  $\frac{\exp M_{i,j}}{\sum_j \exp M_{i,j}} = \mathcal{S}_{\max}(N)$ ? How do we pad it?

- **Padding the inputs before matrix multiplication will not work.** For example, take the last  $K_n \in d \times m$  and for the computation of the matrix multiplication  $Q_0 K_n$  we must increase  $m$  to  $m+1$  (say, alignment properties). This may affect also  $V_n$ . This will change the shape of the results and in practice, we would introduce  $d$  zeros into the matrix results.

This will introduce 1 into the matrix  $\exp(Q_0 K_n)$  and when applied to  $S_{max}(Q_0 K_n)$ , then we change the shape of the numerator and we increase the denominator (by one in this case). If we need to compute the max for range purpose as in Equation 13, we are changing also the values of the numerator matrix and given the max is not bijective, we cannot reverse the effects. Padding the operands and in combination with range constraints, we cannot have a clean computation.

- **Padding the output of  $\exp(M)$  with zeros.** In practice, the unitary element of the denominator sum is zero  $\exp(-\infty)$ , the unitary element for max is  $-\infty$ . Either pad the output with zero or the input with  $-\infty$ . In finite precision and for max, any min instead of  $-\infty$  value will work, however, the  $\exp(\min)$  will not translate to zero in floating point representation unless we fix the result to zero.

How do we introduce  $-\infty$  ? or how we introduce zeros to a specific output pattern?

#### D. Implementation of the head

We discuss here the case where we perform the operations with *int16*. Consider  $L = 512$ , we split the matrices  $Q$  and  $V$  into horizontal parts and we split  $K^t$  into vertical parts. We choose to split  $K^t$  and  $V$  equally into four core so each core will receive  $2^6 \times \lceil \frac{2^9}{2^2} \rceil = 2^7$  for a space of  $2^{14}$  Bytes and we need to store the same amount of  $V$  for a total space of  $2^{15}$  this is equivalent to four banks. We can split the computation so that  $V_i, K_i \in 64 \times 64$  and thus we send into each core  $2 V_i$  and  $2 K_i$ .

Let consider what should be the size of  $Q_0 = m \times 64$  so that we can fit a basic computation in a core.

$T = Q_0 K_i = m \times 64$ ,  $T = \exp T = m \times 64$ ,  $D+ = T$ ,  $T = TV_i$ , and  $N+ = T$ . We need to have space for 4 matrices of size  $m \times 64$ . Thus  $m2^3 2^6$  must fit into 2 banks  $2^{14}$ . We have  $m = 2^5 = 32$ .

Core  $A_i$  compute  $N_{i*2, (i+1)*2-1}$  and  $D_{i*2, (i+1)*2-1}$ . With three reductions per column we have the computation of  $R_0$ , we repeat the process. The subvolume for  $Q_i = 32 \times 64$ . We need to broad cast  $Q_i$  16 times.

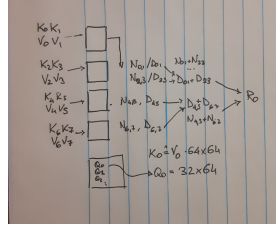


Fig. 1. HEAD for int16

In this section, we show the algorithm that can be implemented on a  $2 \times 4$  having the problem size  $d = 64$  and  $L = 512$

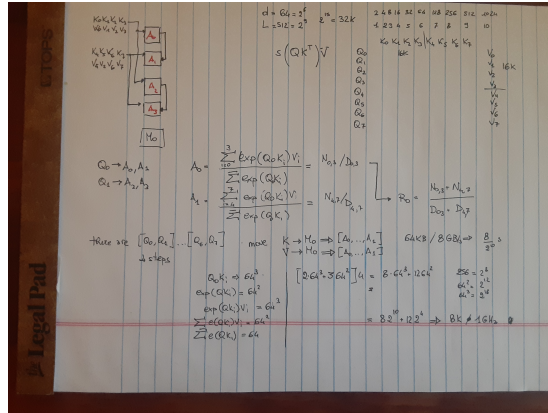


Fig. 2. HEAD for int 8

#### E. The curious case of the problem size $l = 77$ and $L = 768$

This is the place holder for the curious case and may be a code generator.

### III. GEMM FOR ASYMMETRIC SYSTOLIC ARRAYS

We start a conversation about  $C = A * B$  where we imply that the right operand is larger than the left operand. If it may happen the opposite, then we may want to compute  $C^t = B^t * A^t$ , the shapes may be different but the sizes stay the same. The algorithm in a nut shell is a systolic algorithm where the matrix C is spatially separated and computed completely before written. Any matrices we are going to present in this section will have the usual notation

$$M = \begin{pmatrix} M_{0,0} & M_{0,1} & M_{0,2} & M_{0,3} & \dots \\ M_{1,0} & M_{1,1} & M_{1,2} & M_{1,3} & \dots \\ M_{2,0} & M_{2,1} & M_{2,2} & M_{2,3} & \dots \\ M_{3,0} & M_{3,1} & M_{3,2} & M_{3,3} & \dots \\ \dots & \dots & \dots & \dots & \dots \end{pmatrix}$$

Where each element is a sub matrix. We will refer these as subvolumes but they can be reduce to scalar if necessary. We assume an array of AIEs in a formation  $4 \times 2$ . Two columns and four rows.

$$\begin{pmatrix} C_{0,0} = \sum_k A_{0,k} B_{k,0} & C_{1,0} = \sum_k A_{1,k} B_{k,0} \\ C_{0,1} = \sum_k A_{0,k} B_{k,1} & C_{1,1} = \sum_k A_{1,k} B_{k,1} \\ C_{0,2} = \sum_k A_{0,k} B_{k,2} & C_{1,2} = \sum_k A_{1,k} B_{k,2} \\ C_{0,3} = \sum_k A_{0,k} B_{k,3} & C_{1,3} = \sum_k A_{1,k} B_{k,3} \end{pmatrix} \quad (14)$$

$A_{0,k}$  is broad cast on column zero and  $A_{1,k}$  to column one, one subvolume at a time.  $B_{k,0}$  is broad cast on row zero and  $B_{k,3}$  on row three, one subvolume at a time. This is a systolic communication.

#### A. Sub-Volume

The basic computation for one core is  $C_{i,j} = \sum_k A_{i,k} B_{k,j}$ , the basic subvolume for this computation is  $C_{i,j}, A_{i,k}, B_{k,j}$ . They are not same-size matrices. They need to be in the AIE cores. This has eight banks of 8 KB each ( $2^{13}$  elements). We assume that the computation and communication to the core is double buffered.

---

#### Algorithm 1 Core streaming computation for $C_{0,0}$

---

- 1:  $T_A = A_{0,0}, T_B = B_{0,0}$
  - 2:  $\forall i[1, K-1], TT_A = A_{0,i}, TT_B = B_{i,0}, C_{0,0} += T_A T_B, T_A = TT_A, T_B = TT_B$
  - 3:  $C_{0,0} += T_A T_B$
- 

We decide sub-volume  $C_{0,0}$  will fit a bank  $2^{13}$ , to keep the most balance computation we choose  $(128 = 2^7) \times (64 = 2^6)$  (8 KB). This means that  $A_{0,0}$  will be  $128 \times 64$  (8 KB) and  $B_{0,0}$  will be  $64 \times 64$  (4 KB).

One step (tick) of the overlay computation is an *outer* product

$$\begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \end{pmatrix} = \begin{pmatrix} A_{0,k} \\ A_{1,k} \end{pmatrix} (B_{k,0} \ B_{k,1} \ B_{k,2} \ B_{k,3}) \quad (15)$$

We read 16 KB of  $A$  ( $256 \times 64$ ), we read 16 KB of  $B$  ( $64 \times 256$ ), and we compute  $256 \times 256$  elements of  $C$ . In practice, the full computation is formally, where we introduce a step in the outer loop:

$$\sum_{i=0,2}^M \sum_{j=0,4}^N \begin{pmatrix} C_{i,j} & C_{i,j+1} & C_{i,j+2} & C_{i,j+3} \\ C_{i+1,j} & C_{i+1,j+1} & C_{i+1,j+2} & C_{i+1,j+3} \end{pmatrix} = \sum_k^K \begin{pmatrix} A_{i,k} \\ A_{i+1,k} \end{pmatrix} (B_{k,j} \ B_{k,j+1} \ B_{k,j+2} \ B_{k,j+3})$$

We have a ratio compute over communication per core

$$\mathbb{F} = \frac{2^7 * 2^6 * 2^7}{2 * 2^7 * 2^6} = \frac{2^{20}}{2^{14}} = 2^6 \quad (16)$$

This ratio is useful to estimate the smallest subvolume we can afford in order to hide communication latency by double buffering or figure out the ratio between the computation speed and bandwidth.

For example, assume we can send each operand by different channels at 8 GB/s =  $2^3 * 2^{30}$  and 256 operations can be done in one cycle at 1 GHz. The ratio above will become

$$\mathbb{F}_c = \frac{2^7 * 2^6 * 2^7 * 2^{33}}{2 * 2^7 * 2^6 * 2^8 * 10^9} \sim \frac{2^{53}}{2^{14+8+28=50}} > 2^3$$

For any values larger than one, the execution time is compute bound. In the same way we can choose the operand sizes:

$$\frac{2^3 * 2^3 * 2^3 * 2^{33}}{2 * 2^3 * 2^3 * 2^8 * 10^9} \sim \frac{2^{42}}{2^{7+8+28=43}} > 1/2$$

any sub volume smaller than  $8 \times 8$  will be communication bound and double buffering will not help. This Ratio can be used for the investigation of double buffering and latency of reading operands in Mem-Tiles.

1) *Transpose Algorithm*: Because the sizes and shapes of the operands does not need to be the same, we may need to consider also the transpose algorithm

$$\begin{pmatrix} C_{0,0} & C_{0,1} \\ C_{1,0} & C_{1,1} \\ C_{2,0} & C_{2,1} \\ C_{3,0} & C_{3,1} \end{pmatrix} = \begin{pmatrix} A_{0,k} \\ A_{1,k} \\ A_{2,k} \\ A_{3,k} \end{pmatrix} (B_{k,0} \quad B_{k,1}) \quad (17)$$

We read 16 KB of  $A$  ( $64 \times 256$ ), we read 16 KB of  $B$  ( $256 \times 64$ ), and we compute  $256 \times 256$  elements of  $C$ . In practice, the full computation is formally, where we introduce a step in the outer loop:

$$\sum_{i=0,4}^M \sum_{j=0,2}^N \begin{pmatrix} C_{i,j} & C_{i,j+1} \\ C_{i+1,j} & C_{i+1,j+1} \\ C_{i+2,j} & C_{i+2,j+1} \\ C_{i+3,j} & C_{i+3,j+1} \end{pmatrix} = \sum_k^K \begin{pmatrix} A_{i,k} \\ A_{i+1,k} \\ A_{i+2,k} \\ A_{i+3,k} \end{pmatrix} (B_{k,j} \quad B_{k,j+1})$$

The main difference and utility to have two algorithms (for memtile-core) is the better way to have a different ratio of the computation. As a practical point, Equation 15 leads to a reuse of  $A$  in memtile (if  $A$  smaller space) and Equation 17 leads to a reuse of  $B$  (if  $B$  is smaller). This will be evident when we write and present a code generator based on this document. For

2) *Different sizes: int16*: In the case, we are working with an element 2 Bytes (int 16) and we can keep the same subvolumes, the fraction  $\mathbb{F} = \frac{2^6}{2}$ . The number of computation does not change, the communication increases linearly.

The sub-volume computations has to change. We take  $VcC_{0,0} = 64 \times 64 = 2^6 \times 2^6$ , this means of size  $2 \times 2^{12}$ , this is 8 KB (again).  $A_{0,k}, B_{k,0} = 64 \times 64$ , one bank again. The subvolume computation is  $64 \times 64 \times 64$ . Changing the subvolume does not affect the Bytes transfer per tick (16 KB for each operand). This is still balanced.

*B. Mem-tiles and the case  $A = 512 \times 768$   $B = 768 \times 768$*

The matrix  $A$  has size  $(512 = 2^9) \times (768 = 3 \times 2^8)$  that is  $3 \times 2^{17}$  and this fits nicely in a memtile of size 512 KB, that is  $2^{19}$  elements. The matrix  $B$  has size  $768 \times 768$  that is  $9 \times 2^{16} > 2^{19}$ . The nice property is that we can split the matrix  $B$  by column into three parts exactly and we can do double buffering

$$\sum_{j=0,4}^{12} \left\{ \sum_{i=0,2}^4 \begin{pmatrix} C_{i,j} & C_{i,j+1} & C_{i,j+2} & C_{i,j+3} \\ C_{i+1,j} & C_{i+1,j+1} & C_{i+1,j+2} & C_{i+1,j+3} \end{pmatrix} \right\} = \sum_k^{12} \begin{pmatrix} A_{i,k} \\ A_{i+1,k} \end{pmatrix} (B_{k,j} \quad B_{k,j+1} \quad B_{k,j+2} \quad B_{k,j+3})$$

The curly brackets show the part of the computation when the sub matrix of  $B$  of size  $768 \times 256$  is double buffered.

*C. Mem-tiles and the case int16  $A = 512 \times 768$   $B = 768 \times 768$*

The matrix  $A$  has size  $(512 = 2^9) \times (768 = 3 \times 2^8)$  that is  $2 \times 3 \times 2^{17}$ , We need to tile and double buffering  $A$ . We need need to feed two columns:  $A_{0,*}$  and  $A_{1,*}$  at a time. One ping per column will be  $2 \times (64)$ : If we split  $A$  into 4 tiles of 128 rows (4 subvolumes), We can have 4 iterations where we double buffer  $A_{i,*}$  and  $A_{i+1,*}$ . We need to keep a space of 384 KB in memtile 0.

The matrix  $B$  has size  $768 \times 768$  that is  $2 \times 9 \times 2^{16} > 2^{19}$ . If we keep the same sub-volume, we need  $B_{*,0} \dots B_{*,3}$  thus  $2 \times 256 \times 768$ , four double buffering we need twice as much as that. To make this possible we need to change the subvolume of  $B$  and of  $C$ .

To allow double buffering of the  $B$  in memtile, The sub-volume computations has to change. We take  $VcC_{0,0} = 64 \times 32 = 2^6 \times 2^5$ , this means of size  $2 \times 2^{11}$ , this is 4 KB (again).  $A_{0,k} = 64 \times 64$  one bank,  $B_{k,0} = 64 \times 32$ , 4 KB. The subvolume computation is  $64 \times 64 \times 32$ . Changing the subvolume does affect the Bytes transfer per tick (16 KB for A and 8KB for B).

$$\sum_{j=0,4}^{24} \left\{ \sum_{i=0,2}^8 \begin{pmatrix} C_{i,j} & C_{i,j+1} & C_{i,j+2} & C_{i,j+3} \\ C_{i+1,j} & C_{i+1,j+1} & C_{i+1,j+2} & C_{i+1,j+3} \end{pmatrix} \right\} = \sum_k^{12} \left\{ \begin{pmatrix} A_{i,k} \\ A_{i+1,k} \end{pmatrix} (B_{k,j} \quad B_{k,j+1} \quad B_{k,j+2} \quad B_{k,j+3}) \right\}$$

---

**Algorithm 2** Small Small still need tiling

---

```

1: for i=0 to 8 step 2 do
2:   Ping/Pong  $A_{i,*} \dots A_{i+1,*}$ 
3:   # Reuse A but hide latency
4:   for j=0 to 24 step 4 do
5:     Ping/Pong  $B_{*,j} \dots B_{*,j+3}$ 
6:      $\begin{pmatrix} C_{i,j} & C_{i,j+1} & C_{i,j+2} & C_{i,j+3} \\ C_{i+1,j} & C_{i+1,j+1} & C_{i+1,j+2} & C_{i+1,j+3} \end{pmatrix} = \sum_k^{12} \begin{pmatrix} A_{i,k} \\ A_{i+1,k} \end{pmatrix} (B_{k,j} \quad B_{k,j+1} \quad B_{k,j+2} \quad B_{k,j+3})$ 
7:   end for
8: end for

```

---

$$\mathbb{F}_{cm} = \frac{2^6 * 2^6 * 2^5 / 128 / 10^9}{2(2^{12} + 2^{11}) / 8 / 2^{30}} = 0.75$$

This is communication bound. If we consider the computation of a core and the time to read a ping of  $B$  from DDR to memtile (considering we can use 2 channels 8 GBs). We can reach break even with the computation if we use 4 channels (copy to both Memtiles)

$$\mathbb{F}_{md} = \frac{2^6 * 3 * 2^8 * 2^5 / 128 / 10^9}{2 * 2^7 * 32^8 / (2 * 8 * 2^{30})} = 0.53$$

If we assume we can use 4 channels for DDR. The execution time will be  $6 * 8 * 12 * 2(2^{12} + 2^{11}) / 8 / 2^{30}$  which is 411 us

#### D. Mem-tiles and the case $A = 512 \times 768$ $B = 768 \times 768 * 4$

The matrix  $A$  has size  $(512 = 2^9) \times (768 = 3 * 2^8)$  that is  $3 * 2^{17}$  and this fits nicely in a memtile of size 512 KB, that is  $2^{19}$  elements. The matrix  $B$  has size  $768 \times 768 * 4$  that is  $9 * 2^{18} > 2^{19}$ . The nice property is that we can split the matrix  $B$  by column into 12 parts exactly and we can do double buffering

$$\sum_{j=0,4}^{48} \left\{ \sum_{i=0,2}^4 \begin{pmatrix} C_{i,j} & C_{i,j+1} & C_{i,j+2} & C_{i,j+3} \\ C_{i+1,j} & C_{i+1,j+1} & C_{i+1,j+2} & C_{i+1,j+3} \end{pmatrix} = \sum_k^{12} \begin{pmatrix} A_{i,k} \\ A_{i+1,k} \end{pmatrix} (B_{k,j} \quad B_{k,j+1} \quad B_{k,j+2} \quad B_{k,j+3}) \right\}$$

#### E. Mem-tiles and the case $A = 512 \times 768 * 4$ $B = 768 * 4 \times 768$

The matrix  $A$  has size  $(512 = 2^9) \times (768 = 3 * 2^{10})$  that is  $3 * 2^{19}$ . The matrix  $B$  has size  $768 * 4 \times 768$  that is  $9 * 2^{18} > 2^{19}$ .

We present two algorithms. The first algorithm change the subvolume so that we can double buffer A and reuse a slice of A multiple times.

Consider first that a matrix of size  $512 \times 768$  fits a Memtile. A slice of  $A$  of size  $64 \times 768 * 4$  allows a double buffering. This means that a column will solve a subvolume of size  $m = 32$  (64 divided by 2). Equivalently, a slice of  $B$  of size  $768 * 4 \times 64$  will allow double buffering and thus a sub volume of  $n = 16$  (64 divided by 4).

$C_{i,j}$  has size  $32 \times 16$ .  $A_{i,k}$  has size  $32 \times 64$ .  $B_{k,j}$  has size  $64 \times 16$ .

$$\sum_{i=0,2}^8 \left\{ \sum_{j=0,4}^{96} \begin{pmatrix} C_{i,j} & C_{i,j+1} & C_{i,j+2} & C_{i,j+3} \\ C_{i+1,j} & C_{i+1,j+1} & C_{i+1,j+2} & C_{i+1,j+3} \end{pmatrix} = \sum_k^{192} \left\{ \begin{pmatrix} A_{i,k} \\ A_{i+1,k} \end{pmatrix} (B_{k,j} \quad B_{k,j+1} \quad B_{k,j+2} \quad B_{k,j+3}) \right\} \right\}$$

---

#### Algorithm 3 Large Large

---

```

1: for i=0 to 8 step 2 do
2:   Ping/Pong  $A_{i,*} \dots A_{i+1,*}$ 
3:   # Reuse A but hide latency
4:   for j=0 to 48 step 4 do
5:     Ping/Pong  $B_{*,j} \dots B_{*,j+3}$ 
6:      $\begin{pmatrix} C_{i,j} & C_{i,j+1} & C_{i,j+2} & C_{i,j+3} \\ C_{i+1,j} & C_{i+1,j+1} & C_{i+1,j+2} & C_{i+1,j+3} \end{pmatrix} = \sum_k^{192} \begin{pmatrix} A_{i,k} \\ A_{i+1,k} \end{pmatrix} (B_{k,j} \quad B_{k,j+1} \quad B_{k,j+2} \quad B_{k,j+3})$ 
7:   end for
8: end for

```

---

#### IV. GEMM TILES GENERATOR

In general, the core computation needs to be a multiple aligned to 8, say, for the channels (columns of the operands). Subvolumes being a power of two are nice and not necessary. What if we can explore the sub-volume solution space with some type of properties? In Equation 16, the ratio  $\mathbb{F}$  is a function of the subvolume ( $\mathbb{F}(m, n, k)$ ). This sub-volume can be for a core (streaming) or the representation of the cluster computation (all cores with all memtiles). We argue for three requirements: we should have a subvolume that fit into a core (basic subvolume), we want to have largest number of operations per subvolume ( $\mathbb{C}(m, n, k)$ ), and among them we prefer the one that have the largest ratio  $\mathbb{F} = \frac{\mathbb{C}}{\mathbb{S}}$  (where  $\mathbb{S}$  stands for space and 24 KB represent the space for 6 banks in a core). We can call  $\mathbb{T}$  the set of valid solutions and we call  $\mathbb{W}$  the sorted solution space as follows:

$$\mathbb{T} = \{[m, n, k] \text{ s.t. } \mathbb{S}(m, n, k) < 24KB\} \quad (18)$$

$$\mathbb{W} = \text{sort}(\mathbb{T}, \lambda x : (-\mathbb{C}(x), -\mathbb{F}(x))) \quad (19)$$

We call  $\mathbb{Q}_0$  the sum over  $k$  of Equation 15 and  $\mathbb{Q}_1$  of Equation 17. These represent two algorithms and also subproblems. The cluster of cores work together reading from Mem-Tile and in the context of the larger problem we can double buffer either **A** or **B**.

If we chose  $(m, n, k) = \mathbb{W}[0]$  as starting block for the core computation. Then  $\mathbb{Q}_0$  will have a memtile subvolume  $M_0 = (2 * m, 4 * n, K)$  and  $\mathbb{Q}_1$   $M_1 = (4 * m, 2 * n, K)$ . Both computations will be the reduction of  $K/k$  iterations. And we can choose the algorithm so that  $\mathbb{F}(M_0)$  and  $\mathbb{F}(M_1)$  is the largest. The core subvolume is determined by space and  $\mathbb{F}$ , so it is the Mem-tile subvolume and thus the algorithm itself.

At this time, we will not deal with padding and alignment per core and per memtile. We summarily say that the core subvolume  $(m, n, k)$  will determine a minimum zero padding so that a core will produce the correct result. Padding at cluster level does not need to be, we can accept either the turn off or columns/rows or the computation of garbage and not writing back to DDR.

##### A. Double Buffering aka DB

Rarely **A** and **B** will fit the Memtile and some double buffering may be required. Given a core subvolume, If we can do double buffering for both operands, there is no problem.

$\mathbb{Q}_1$  is more efficient for large **B**, we tend to reuse **A**, and at a minimum we need to double buffer **B**. Clearly the core subvolume is related to the memtile subvolume and whether we can double buffer. If we must tile the  $K$  dimension, we can design the algorithm and estimate the execution time. Then we can restart the search space asking for a problem size that we can do DB of **B** and compare.

$\mathbb{Q}_0$  is more efficient for large **A**, we tend to reuse **B**, and at a minimum we need to double buffer **A**. Clearly the core subvolume is related to the memtile subvolume and whether we can double buffer. If we must tile the  $K$  dimension, we can design the algorithm and estimate the execution time. Then we can restart the search space asking for a problem size that we can do DB of **A** and compare.

This search is almost complete. This is simply a set of heuristics for the selection of the solution space, the algorithms, and the tiling. In the, following section we shall present the ideas for a automatic selection and exploration of the solution space and its optimal solution.

#### V. SOLUTION SPACE AND ITS EXPLORATION

We are working on a software tool that will explore a solution space for a given problem. For GEMM this is triplet  $(M, N, K)$  (i.e.,  $C \in M \times N$ ,  $A \in M \times K$ ,  $B \in K \times N$ ). Briefly, tiling  $(M, N, K)$  to memtile is another triplet  $(M_m, N_m, K_m)$ . This can describe the (largest) problem we solve using data in memtile. We can go down further and describe the tiling at core level as  $(M_c, N_c, K_c)$ . In practice,  $(M, N, K) = P$ ,  $(M_m, N_m, K_m) = m$ , and  $(M_c, N_c, K_c) = c$  are enough to describe a complete computation.

The solution space is the composition of all valid  $\mathcal{S} = \{(P, m, c) \text{ s.t. valid}\}$ . A tiling tools have two main goals: First, we need to describe a complete set of solutions. In general, this can be quite large; however hardware requirements and designs strategies make the solution space a limited one. Second, we can address cost functions and their effect on the overall solution space simply by translating the solution  $P, m, c$  either into actual code or other realistic and consistent measures.

In all previous work, we did compute algorithm choices (such as tiling) by back of the envelop cost estimates, code generation cost estimate, simulations, and actual time measures. Here, we reduce to a minimum the description of the algorithm and hardware. We compute tiling for the application of double buffering at core and memtile level. We aim at symmetric computations (so that we can have uniform and predictable data movement).

In these subsections, we introduce the main idea for tiling and solution space we use for each computation. We assume in all cases that the computation start from DDR and we end in DDR. Problem  $P$  will start in DDR and the result will be in

DDR. The operand memory allocation in DDR is given and the memtile is considered clean and clear. If the P fit memtile, then  $(P, P, *)$  may be a valid solution. For this tool, we do not use addresses and we assume that allocation is always possible if space is available. We are aware this last statement is not always true.

#### A. GEMM

The current architecture is  $4 \times 2$  architecture and thus the computation is rectangular by nature. We have two main algorithms Equation 15 and 17. We search first valid core solutions  $(M_c, N_c, K_c) = c$  using property of alignment specific to the HW, then to the space requirements for the operands (with double buffering).

For Equation 15, we need to read from memtile  $(2 * M_c, 4 * N_c, K) = m$  and for 17 we need  $(4 * M_c, 2 * N_c, K) = m$ . If  $K$  is too large, we will split  $K$  accordingly  $((4 * M_c, 2 * N_c, K/4) = m, \text{ say})$  and stream the computation.

First, the data in memtile is split perfectly across the cores. We aim at padding at core level and if we can avoid better. If we can double buffer A or B or both we try, if we split by  $K$ , we double buffer portions of A and B on the  $K$  dimension. Second we can split the computation between DDR and memtile  $m$  and thus we have the final solution  $P, m, c$ . This describes completely our solution space from core to original problem size. We use a cost function, time, space, compute and ratio to explore the solution space and the Pareto curve.

We provide two different classes and approaches: we provide a solution space stemming from a single problem and from a set. A single problem create a list of solutions but we do not take advantage of any repetition. For a list of problem we create dictionary of solutions. This can help us to explore solutions that would not be available and consider in case extra constraints are introduced.

#### B. MHA

To exploit performance is to read into cores the matrices  $K$  and  $V$  so that to perform the largest block inner product  $\sum_j \exp(Q_0 K_j) * V_j$ . The idea is to compute the space requirement in order to compute locally in a core the numerator Equation 9 and the denominator Equation 10. Then do a reduction by a single core for the final computation Equation 11.

The minimum number of cores is one, if it is possible we rather use a single column. The reduction is simpler done by a single column, and the next column can work on  $Q_1$  and  $R_1$ .

We start with a problem of size  $P = (d, L_0, L_1, r)$  so that  $Q \in L_0 \times d$ ,  $K \in d \times L_1$ , and  $V \in L_1 \times r$ . We find all the problem  $c = (d, l_c, \ell_c, r)$  that fit in a core and we can estimate the execution time.

#### C. Layer Norm

The solution space and cost estimation follow the same ideas developed for MHA in the previous sections.

#### D. Convolution

The way we build the solution space is by observing the output tensor from the smallest to the largest and compute the projection to the input. Projection is a function that given an output tensor of size  $(1, H, W, C_{out})$  and a convolution with stride  $(0, s_h, s_w, 0)$  and kernel  $(C_{out}, k_h, k_w, C_{in})$ , the projection is  $(1, (H - 1)s_h + k_h, (W - 1)s_w + k_w, C_{in})$ . We have an starting problem  $P$ .

Consider now that the output tensor will be created by two (or more) columns of cores. We consider the computation split by width so that each column compute may compute any  $v \in [0, w/2]$  and each core compute  $c_o \in [0, C_{out}/4]$  with a proper granularity (alignment). We assume we are streaming the computation by  $H$  (this is important for the computation description and its time estimates)

The core will use a space for  $c$

- output  $(1, h, v, c_o)$
- input  $(1, (h - 1)s_h * k_h, (v - 1)s_w + k + w, C_{in})$
- weights  $(c, k_h, k_w, C_{in}) + (1, 1, 1, c_o)$

For each core subvolume, we consider in memtile only multiple of the core subvolume  $m$ ; that is,  $H = m * h$  and  $V = n * v$ .

- output  $(1, H, V, c_o)$
- input  $(1, (H - 1)s_h * k_h, (H - 1)s_w + k + w, C_{in})$
- weights  $(c, k_h, k_w, C_{in}) + (1, 1, 1, c_o)$

So for the same  $c$  we may have multiple  $m$  that can fit in memtile with double buffering.

We have then a complete representation of the solution by our usual triplet  $(P, m, c)$  and we can estimate the performance by a simple cost function.

Notice, If we want to split the computation by input channel,  $C_{in}$ , we may just transform the convolution into two or more convolutions followed by element wise additions. If we want to split the computation by  $C_{out}$ , we can transform the convolution into two convolutions and a concatenation.