

# Compiler Design Notes

Ayush Raina

April 4, 2025

## Iteration Space and Iteration Vector

Consider the code block below:

```
1 for(i = 0; i < N; ++i) {
2     for(j = 0; j < N; ++j) {
3         /* Do Some Work */
4     }
5 }
```

$I_s = \{(i, j) : 0 \leq i, j < N\}$  and each point  $(i, j)$  in this space is called an iteration vector.  $i, j$  are called **iteration variables** or **loop induction values**. Iteration vectors tuples contains loop induction values starting from outermost loop to innermost loop.

Consider the following code block:

```
1 for(i = 0; i < N; ++i) {
2     for(j = 0; j < i; ++j) {
3         for(k = 0; k < j; ++k) {
4             /* Do Some Work */
5         }
6     }
7 }
```

In this case  $I_s = \{(i, j, k) : 0 \leq i < N, 0 \leq j < i, 0 \leq k < j\}$ , and each point  $(i, j, k)$  in this space is called an iteration vector.  $i, j, k \in \mathbb{Z}$  are called **iteration variables** or **loop induction values**.

## Lexicographical Order

$(i, j) < (i', j')$  if  $i < i'$  or  $i = i'$  and  $j < j'$ . In General for  $n$  dimensional space,  $(i_1, i_2, \dots, i_n) < (i'_1, i'_2, \dots, i'_n)$  if  $\exists k \in \{1, 2, \dots, n\}$  such that  $i_k < i'_k$  and  $i_i = i'_i$  for  $i < k$ .

## Lexicographical Ordering of Iteration Vectors

Consider the following code block:

```
1 for(i = 0; i < N; ++i) {
2     for(j = 0; j < N; ++j) {
3         /* Do Some Work */
4     }
5 }
```

The lexicographical ordering of iteration vectors is as follows:  $(0, 0), (0, 1), (0, 2), \dots, (0, N-1), (1, 0), (1, 1), \dots, (1, N-1), \dots, (N-1, 0), (N-1, 1), \dots, (N-1, N-1)$ . In this course we will plot  $i$  on y-axis and  $j$  on x-axis.

## Lexicographically $> 0$

Iteration vectors  $(a_1, a_2, \dots, a_n)$  are said to be lexicographically  $> 0$  if first non zero loop induction value is positive. For example,  $(0, 0, 0)$  is lexicographical  $\geq 0$ ,  $(0, 0, 1)$  is lexicographical  $> 0$ .

Consider the following code block:

```
1 for(i = 0; i < N; ++i) { /* Loop Header */
2     for(j = 0; j < i; ++j) { /* Loop Header */
3         A[i][j] = A[i-1][j] + A[i][j-1]; /* Loop Body */
4     }
5 }
```

Above code blocks has  $N^2$  iteration vectors. Hence  $N^2!$  possible ways to arrange these vectors. How to find the valid orderings of these vectors?. A valid ordering means an ordering which does not change the semantics of the program.

## Dependence Analysis

RAW - Read After Write  
WAW - Write After Write  
WAR - Write After Read

These are some of the dependencies that can occur in a program. Using this information we can take out the invalid orderings of the iteration vectors. In above program there are 2 reads  $R_1, R_2$ , one write  $W_1$ . Two instructions are dependent in one of them is write and other is read, same in the case of load and store. In this we have Write( $W_1$ ) after Read  $R_1, R_2$ . In this program to do  $W_1$  at  $(i, j)$  using  $R_1$  at  $(i', j')$ , the following conditions must be satisfied:

- $i = i' - 1$
- $j = j'$

Similarly we have dependence between  $W_1$  and  $R_2$ . To do  $W_1$  at  $(i, j)$  using  $R_2$  at  $(i', j')$ , the following conditions must be satisfied:

- $i = i'$
- $j = j' - 1$

## Distance Vectors

Consider the following code block:

```
1 for(i = 0; i < N; ++i) {
2     for(j = 0; j < i; ++j) {
3         A[j][i] = B[j][i]; /* Copying Array Elements */
4     }
5 }
```

In this case we can do loop interchange to get the following code block:

```
1 for(j = 0; j < N; ++j) {
2     for(i = 0; i < j; ++i) {
3         A[j][i] = B[j][i]; /* Copying Array Elements */
4     }
5 }
```

But this will change the execution order, but will this loop interchange will lead to same results? If we do dependence analysis here we will find that there is no dependence between the two loops because they both are working on different arrays. Hence any ordering of iterations will lead to same results.

Now change B to A in above code to get the following code block:

```
1 for(i = 0; i < N; ++i) {
2     for(j = 0; j < i; ++j) {
3         A[j][i] = A[i][j]; /* Copying Array Elements */
4     }
5 }
```

We will test for Read After Write dependence between the two loops. Suppose Read  $R_1$  is at  $(i, j)$  and Write  $W_1$  is at  $(i', j')$ . To do  $R_1$  at  $(i, j)$  using  $W_1$  at  $(i', j')$ , the one of the following two conditions must be satisfied:

- $i' \geq i + 1, j' = j, i' = j, j' = i$
- $i' = i, j' \geq j + 1, i' = j, j' = i$

where  $0 \leq i, j, i', j' \leq N$ . Similarly we can do the same for Write After Read dependence analysis. Suppose Write  $W_1$  is at  $(j, i)$  and Read  $R_1$  is at  $(i', j')$ . To do  $W_1$  at  $(j, i)$  using  $R_1$  at  $(i', j')$ , the one of the following two conditions must be satisfied:

- $i' \geq i + 1, j' = j, i' = j, j' = i$
- $i' = i, j' \geq j + 1, i' = j, j' = i$

where  $0 \leq i, j, i', j' \leq N$ .

Consider the  $n$  dimensional space  $I_s = \{(i_1, i_2, \dots, i_n) : 0 \leq i_1, i_2, \dots, i_n < N\}$ . Let  $\vec{s}, \vec{t}$  be given. The distance vector between  $\vec{s}$  and  $\vec{t}$  is defined as  $\vec{d} = \vec{t} - \vec{s}$ . Then  $\forall (\vec{s}, \vec{t})$  st there is a dependence between  $\vec{s}$  to  $\vec{t}$ , then  $T(\vec{t}) - T(\vec{s}) > 0$ , where T is a linear transformation.

Further if T is a linear transformation then  $T(\vec{t}) - T(\vec{s}) = T(\vec{t} - \vec{s}) = T(\vec{d}) > 0$ . Hence  $\vec{d}$  is lexicographically  $> 0$  if  $(\vec{s}, \vec{t})$  is a constant distance vector.

## Fully Permutable Loops

If all the Permutations of distance vectors are constant i.e all the components are positive then the loops are said to be fully permutable. All the Loop Orderings will lead to same results.

## Optimizations : Locality

Locality means being close together, we want to focus on data locality which means how close are we to data that we are using. We will talk about Latency and Bandwidth. Latency: Time taken to access a single data element. Bandwidth: Number of data elements that can be accessed in a given time, say per cycle.

Further we have two types of locality:

1. **Locality in Space (Spatial Locality)** : If we are accessing a data element at some particular address, we are likely to access the different address close to it.

2. **Locality in Time (Temporal Locality)**: If we are accessing a data element at some particular address, we are likely to access the same address in near future.

Compiler transformations can enhance locality. Loop Interchange can enhance the Spatial Locality. For example

```
1  for(int i = 0; i < 100; i++) {
2      for(int j = 0; j < 100; j++) {
3          /* Do Some work using f(j,i) */
4      }
5  }
```

In the above code, even when who cache line is loaded, we are not accessing the data in the order it is stored in the cache line. Hence we can interchange the loops to get the following code block:

```
1  for(int j = 0; j < 100; j++) {
2      for(int i = 0; i < 100; i++) {
3          /* Do Some work using f(j,i) */
4      }
5  }
```

This will enhance the spatial locality. Following is a example of Temporal Reuse: (Reuse in Time)

```
1  for(int i = 0; i < 100; i++) {
2      for(int j = 0; j < 100; j++) {
3          /* Do Some work using f(i) */
4      }
5  }
```

In General there are two types of reuse for both spatial and temporal locality - **Self Reuse, Group Reuse**

## Group Temporal Reuse

Consider the following code block:

```
1  load a[i][j];
2  load a[i][j+1];
3  load a[i][j+2];
```

These kind of access is called Group Temporal Reuse because  $a[i][2], a[i][3], \dots, a[i][N - 3]$  are accessed 3 times. This is called Group Temporal Reuse. These kind of access can also be called as Group Spatial Reuse since we are accessing nearby (**contiguous**) memory locations.

```
1  load a[i][j];
```

This of access is called Self Temporal Reuse because  $a[i][j]$  is again accessed if put in a loop not depending on i and j.

## Matrix Multiplication

Consider the following code block:

```
1  for(i = 0; i < N; ++i) {
2      for(j = 0; j < N; ++j) {
3          for(k = 0; k < N; ++k) {
4              C[i][j] += A[i][k] * B[k][j];
5          }
6      }
7  }
```

Assuming all matrices  $\in \mathcal{R}^{n \times n}$ . Consider Cache Line Size of 64 Bytes, which means it can fit 16 elements of size 4 bytes. Further consider that  $N^2 \gg C_s$  (Cache Size).

Suppose if there is no Cache, then there are  $4N^3$  memory accesses. Now Given Cache, lets try to compute the number of memory accesses which in this case means number of cache misses.

When we Load  $A[1][1]$  is accessed then  $A[1][1], A[1][2], A[1][3], \dots, A[1][15]$  are loaded in the cache. For next 15 iterations we will have cache hits. But when we access  $A[1][16]$  then  $A[1][16], A[1][17], \dots, A[1][31]$  are loaded in the cache and continue.

Overall we will have  $\frac{N^3}{L}$  cache misses for A because of spatial locality. But B does not have spatial locality. Hence we will have  $N^3$  cache misses for B. For  $N$  iterations of  $k$  loop,  $C[i][j]$  is accesses so we assume this is in cache. Hence we will have  $N^2$  loads and  $N^2$  stores for C.

Hence overall memory accesses =  $\frac{N^3}{L} + N^3 + 2N^2 \sim \mathcal{O}(N^3)$ . Look about how accesses to B can conflict with accesses to A and why we only use 4 way set associative cache.

Suppose now we change the loop order from  $(i, j, k) \rightarrow (k, j, i)$ , then we will have the following code block:

```

1 for(k = 0; k < N; ++k) {
2   for(j = 0; j < N; ++j) {
3     for(i = 0; i < N; ++i) {
4       C[i][j] += A[i][k] * B[k][j];
5     }
6   }
7 }

```

Then  $N^3$  accesses for A,  $N^2$  accesses for B and  $2N^3$  accesses for C. Hence overall memory accesses =  $3N^3 + N^2 \sim \mathcal{O}(N^3)$ . This order is worse than the previous order.

Now Lets consider table with all 6 possible orderings of loops. We will have the following table:

Loop Order	A	B	C
$(i, j, k)$	$N^3/L$	$N^3$	$2N^2$
$(k, j, i)$	$N^3$	$N^2$	$2N^3$
$(i, k, j)$	$N^2$	$N^3/L$	$2N^3/L$
$(j, i, k)$	$N^3/L$	$N^3$	$2N^2$
$(j, k, i)$	$N^3$	$2N^2$	$2N^3$
$(k, i, j)$	$2N^2$	$N^3/16$	$2N^3/16$

This is a example of fully permutable loops. In this case we can interchange the loops in any order and we will get the same results but performance will be different. From above table we can find that  $(i, k, j)$  and  $(k, i, j)$  are best orderings and they differ in  $N^2$  memory accesses.

## March 11, 2025

In the above example of matrix multiplication, with 3 loops.  $i$  goes from  $1 \rightarrow N$ ,  $j$  goes from  $1 \rightarrow N$  and  $k$  goes from  $1 \rightarrow N$ . There are  $N^3!$  possible orderings of these loops. We saw 6 of them in the previous example in which dependence is not violated. We want to find out how many are valid execution orders which do not violate the dependencies.

**Note:** Floating Point operations are not associative. On changing the order, overflow can occur.

```

1 for(i = 0; i < N; ++i) {
2   for(j = 0; j < N; ++j) {
3     for(k = 0; k < N; ++k) {
4       C[i][j] += A[i][k] * B[k][j];
5     }
6   }
7 }

```

If  $k$  is innermost loop then we have  $N^2!$  possible orderings which are all valid. But when Loop if of the form  $i, k, j$  then we some more which we need to count.

Given a nested loop and access conditions, we need to find about temporal and spatial locality. Consider the following code block:

```

1 for(i = 0; i < N; ++i) {
2   for(j = 0; j < N; ++j) {
3     for(k = 0; k < N; ++k) {
4       A[i][j + 2*k + k] /* is accessed */

```

```

5     }
6   }
7 }

```

We can write the access as  $\text{matirx}(A)$   $\text{vector}(V)$  product where:

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 1 \end{bmatrix} \text{ and } B = \begin{bmatrix} i \\ j \\ k \end{bmatrix}$$

## Identifying Spatial Reuse

In General for a iteration vector  $\vec{i} = (i_1, i_2, \dots, i_n)$ , there is spatial reuse for  $i_k$  if  $k^{th}$  column has a small non zero element which be possibly negative only in the last row. All other elements in  $k^{th}$  column are zero.

## Identifying Temporal Reuse

```

1 for(i = 0; i < N; ++i) {
2   for(j = 0; j < N; ++j) {
3     for(k = 0; k < N; ++k) {
4       A[i][k] /* is accessed */
5     }
6   }
7 }

```

In this loop, we have temporal reuse over  $k$  loop. (Check and add this)

Consider this code block:

```

1 for(i = 0; i < N; ++i) {
2   for(j = 0; j < N; ++j) {
3     A[i+j][2(i+j)] /* is accessed */
4   }
5 }

```

Is there a temporal reuse in this kind of memory access ? In this case around  $2N$  unique elements are accessed but there are  $N^2$  iteration vectors, so there is kind of pseudo temporal reuse (Do they access same element at any point ?)

Consider another code block:

```

1 for(i = 0; i < N; ++i) {
2   for(j = 0; j < N; ++j) {
3     for(k = 0; k < N; ++k) {
4       A[i][1+j+k][j+k] /* is accessed */
5     }
6   }
7 }

```

The access matrix for the above memory access is:  $A = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix}$  and loop vector is  $B = \begin{bmatrix} i \\ j \\ k \end{bmatrix}$

Access matrix may not be square matrix. Dimensions of the array can be more than the number of loops. Again if 3 dimensional iteration space containing iteration vectors are accessing a 2D matrix there has to be temporal reuse.

An Access has no temporal reuse if  $M$  has full column rank, where  $M$  is the access matrix. We can also say like this for  $\vec{v} \neq \vec{v}' \implies M\vec{v} \neq M\vec{v}' \implies M(\vec{v} - \vec{v}') \neq 0$ .

Hence for temporal reuse  $\exists \vec{v} \neq \vec{v}'$  such that  $M(\vec{v} - \vec{v}') = 0$ . This means  $(\vec{v} - \vec{v}') \in \text{NULL}(M)$ .

In the above definitions we were talking about self spatial and self temporal reuse. Consider the following code block:

```

1 for(int i = 0; i < N; i++) {
2   A[i] + A[i+1] + A[i+2] /* is accessed */
3 }

```

This is an example of  $\theta(1)$  Temporal Reuse. Most of the Deep Learning Computations have very high temporal reuse  $\sim \theta(N^2)$ .

## Degree of Temporal Reuse

Degree of Temporal Reuse, depends on the rank of the access matrix. If the rank of the access matrix is  $r$  then the degree of temporal reuse is  $N - r$ .

```
1 for(int i = 0; i < N; i++) {  
2     B[A[i]]; /* is accessed */  
3 }
```

These kinds of accesses are called indirect accesses. It is not easy to find about locality of these kind of accesses at compile time.

## March 13

If we have full column rank for the access matrix, then the transformation is one to one and in this case there will be no temporal reuse.

Today we will look at another loop transformation called Loop Tiling for exploiting Temporal Locality. This is also called Loop Blocking. Let us go back to our old matrix multiplication example.

```
1 for(i = 0; i < N; ++i) {  
2     for(j = 0; j < N; ++j) {  
3         for(k = 0; k < N; ++k) {  
4             C[i][j] += A[i][k] * B[k][j];  
5         }  
6     }  
7 }
```

In Loop Interchange  $(i, k, j)$  the best we could do is  $\theta(N^3/L)$  cache misses. Can we do better? In the above case we are not getting temporal reuse for  $B$  and this will slow us down.

When we need simultaneous reuse along multiple dimensions of the iteration space, here will be using loop tiling. We are given that  $c \ll N^2e$  where  $e$  is the number of elements in the cache line and  $c$  is cache size. If we carefully observe, to compute row  $C[i]$  the  $A[i]$  row is fixed and is multiplied with all the columns of  $B$  to get the row  $C[i]$ . Once  $C[i][j]$  is computed, now to compute  $C[i][j+1]$  the  $A[i][0]$  may not be present in the cache due to its small size.

Suppose we do tiling with block size  $B$  in horizontal direction then number of cache misses:

$$\begin{aligned} M(A) &= \frac{B}{L} \times \frac{N}{B} \times \frac{N}{B} \times N = \frac{N^3}{BL} \\ M(B) &= \frac{B^2}{L} \times \frac{N}{B} \times \frac{N}{B} \times N = \frac{N^3}{L} \\ M(C) &= \frac{B}{L} \times \frac{N}{B} \times N = \frac{N^2}{L}, \text{ since both read and write is there so } \frac{2N^2}{L} \end{aligned}$$

Since for  $B$  Temporal Reuse is along  $i$  direction and I am currently finishing whole  $j - k$  plane before moving to next  $i$  plane. Now we will take 3D Tile Block. Assumption is if  $B^2/L$  elements fits into cache then along  $i$  direction, we will have all hits once these elements are loaded. Hence  $B^2e \leq c$ .

$$\begin{aligned} M(B) &= \frac{B^2}{L} \times \frac{N}{B} \times \frac{N}{B} \times \frac{N}{B} = \frac{N^3}{BL} \\ M(A) &= \frac{B}{L} \times B \times \frac{N}{B} \times \frac{N}{B} \times \frac{N}{B} = \frac{N^3}{BL} \\ M(C) &= \frac{B^2}{L} \times \frac{N}{B} \times \frac{N}{B} = \frac{N^2}{L}, \text{ but there is read and write so } \frac{2N^2}{L} \end{aligned}$$

Hence  $e \times (B^2 + B^2 + B^2) \leq c \implies B \leq \sqrt{\frac{c}{3e}}$ . Add calculation for  $B$  here for different sizes of  $c$ , generally  $c = 256\text{kb}$ . From a old paper, the lower bound on cache misses is  $\theta(\frac{N^3}{\sqrt{c}})$  for matrix multiplication. In this case Tile Size  $T \propto \sqrt{c}$ . In this case we choosed same  $B$  along all dimensions, but can we choose different  $B$  along different dimensions and we need to do the analysis for this.

## Point of Diminishing Returns

As we increase the tile size, the increase in cache hit percentage will keep decreasing. Suppose for  $B = 2, 4, 8, 16, 32, 64, 128$  then cache hit percentages are 0.5, 0.75, 0.875, 0.9375, 0.96875, 0.984375, 0.9921875. Hence we can see that the increase in cache hit percentage is decreasing. This is called Point of Diminishing Returns.

Learn about 4 way associative cache and how conflicts will be handled if we use this cache.

March 18,2025

Code for Tiled Matrix Multiplication in 3 dimensions: Suppose we choose tile size  $B$ , we will have 3D Tile Block. This reduces cache misses by a factor of  $B$ , where  $B \sim \mathcal{O}(\sqrt{C})$ .

```
1 for(int i = 0; i < N; i += B) { /* Tile Control Loop 1 */
2   for(int j = 0; j < N; j += B) { /* Tile Control Loop 2 */
3     for(int k = 0; k < N; k += B) { /* Tile Control Loop 3 */
4
5       for(int ii = 0; ii < i+B; ii++) { /* Loop 1 */
6         for(int jj = 0; jj < j+B; jj++) { /* Loop 2 */
7           for(int kk = 0; kk < k+B; kk++) { /* Loop 3 */
8             C[ii][jj] += A[ii][kk] * B[kk][jj];
9           }
10        }
11      }
12    }
13  }
14 }
15 }
```

There can be the case in which  $B \nmid N$ . Then we need to change the inner loop bounds to get the following code

```
1 for(int i = 0; i < N; i += B) { /* Tile Control Loop 1 */
2   for(int j = 0; j < N; j += B) { /* Tile Control Loop 2 */
3     for(int k = 0; k < N; k += B) { /* Tile Control Loop 3 */
4
5       for(int ii = i; ii < min(i+B,N); ii++) { /* Loop 1 */
6         for(int jj = j; jj < min(j+B,N); jj++) { /* Loop 2 */
7           for(int kk = k; kk < min(k+B,N); kk++) { /* Loop 3 */
8             C[ii][jj] += A[ii][kk] * B[kk][jj];
9           }
10        }
11      }
12    }
13  }
14 }
15 }
```

We also do not need to compute min in every iteration, since  $i$  does not depend on  $ii$  and same for  $j, k$ . We can precompute these values and use them in the loop to get the following code block:

```
1 for(int i = 0; i < N; i += B) { /* Tile Control Loop 1 */
2   for(int j = 0; j < N; j += B) { /* Tile Control Loop 2 */
3     for(int k = 0; k < N; k += B) { /* Tile Control Loop 3 */
4
5       int ii_end = min(i+B,N); /* Precomputed Value 1 */
6       int jj_end = min(j+B,N); /* Precomputed Value 2 */
7       int kk_end = min(k+B,N); /* Precomputed Value 3 */
8
9       for(int ii = i; ii < ii_end; ii++) { /* Loop 1 */
10        for(int jj = j; jj < jj_end; jj++) { /* Loop 2 */
11          for(int kk = k; kk < kk_end; kk++) { /* Loop 3 */
12            C[ii][jj] += A[ii][kk] * B[kk][jj];
13          }
14        }
15      }
16    }
17  }
18 }
19 }
```

In General, Tiling results in  $2n$  loop nest from  $n$  loop nest. But we can also choose to tile only few loops and not all loops. Now we will look at the conditions for validity of Tiling, when it is okay to change the execution order like above and when it is not okay.

## Validity of Tiling

The following is the sufficient condition for validity of tiling: The tiling is valid if and only if there is no cyclic dependence between tiling and we should be able to come up with a total order to schedule the tiles. This is not a necessary condition as we can construct an example in which Tile  $(0, 1)$  depends on Tile  $(0, 0)$  and Tile  $(1, 0)$  depends on  $(0, 1)$  through a back edge, but I chose my Tile size in such a way that back edge goes to above tile to avoid the cycle.

For example  $(i, j)$  in iteration space and other points  $(i + 1, j), (i, j + 1), (i + 1, j - 1)$  all three depend on  $(i, j)$  then we cannot tile here, as it will create a cycle. In contrast to this we say tiling along a certain dimension is valid if all

the components of dependence vector are positive or negative. We can also say that dependence should be in the same direction along a dimension.

Can we apply some loop transformation to make all dependences in same direction so that tiling is now valid ?

## Affine Transformations

These transformations can enable tiling.

**Skew:**  $T(i, j) = (i, i + j)$ . As you go up, shift the  $j$  by factor of  $i$ . This transformation is called Skew Transformation. Add the example from uday sir's slides here.

## Different kinds of cache misses

Write about cold misses, conflict misses, capacity misses, compulsory misses.

Explicit Copying: Add about This THis benefits in eliminating conflict misses, TLB Misses and overheads gets amortized. Prefetching works smoothly and improves its performance.

## March 20, 2025

We covered the following topics till now:

- Loop Interchange
- Loop Permutation
- Loop Tiling
- Explicit Copying - Add about This. This is also called packing and unpacking. It improves locality by eliminating conflict misses and TLB capacity misses. Improves Prefetching.

Loop transformations consists of following steps:

- Dependence Analysis
- transformation based on cost model
- Code Generation - Transforming the IR.

**Example 1:**

## Loop Skewing

Consider the following 2D loop nest:

```
1 for(int i = 0; i < N; i++) {  
2     for(int j = 0; j < N; j++) {  
3         a[i][j] = a[i-1][j] + a[i-1][j-1] + a[i-1][j+1];  
4     }  
5 }
```

Assume that this nest has following dependences:  $\mathcal{D} = \{(1, 0), (1, 1), (1, -1)\}$ , Here skew transformation can be applied to make all dependences in the same direction. This transformation can be represented in matrix form as follows:

$$\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} i \\ j \end{bmatrix} = \begin{bmatrix} i \\ i + j \end{bmatrix}$$

New Loop variables:  $(t_1, t_2) = (i, i + j)$ , now we have to eliminate  $(i, j)$  from the following system of equations:

Original System

1.  $t_1 = i$
2.  $t_2 = i + j$
3.  $0 \leq i \leq N - 1$
4.  $0 \leq j \leq N - 1$

New System

1.  $0 \leq t_1 \leq N - 1$
2.  $0 \leq t_2 - t_1 \leq N - 1 \implies 0 \leq t_2 \leq N + t_1 - 1$



We have to first eliminate old loop variables then successively apply Fourier Motzkin Elimination one by one to get bounds for new loop variables starting from inner most loop variable to outermost loop variable. In the above system after applying Fourier Motzkin Elimination we get:  $N \geq 1, 0 \leq t_1 \leq N - 1$  which turns out be same as earlier in this case.

Hence our transformed dependences in this case are:  $\mathcal{D}' = \{(1, 1), (1, 2), (1, 0)\}$  after applying skew transformation.

**Example 2:** Now consider the following transformation:  $T(i, j) = (i + j, j)$ . This transformation can be represented in matrix form as follows:  $\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} i \\ j \end{bmatrix} = \begin{bmatrix} i + j \\ j \end{bmatrix}$

This transformation is not valid in above case, since after applying this there will still be a negative component in dependence vector.

For now consider the following code:

```
1 for(int i = 0; i < N; i++) {
2   for(int j = 0; j < N; j++) {
3     a[i][j] = a[i-1][j] + a[i][j-1];
4   }
5 }
```

This has  $\mathcal{D} = \{(0, 1), (1, 0)\}$ . We want to calculate new loop bounds for  $(t_1, t_2)$  after applying this transformation:

$$\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

Original System

1.  $t_1 = i + j$
2.  $t_2 = j$
3.  $0 \leq i \leq N - 1$
4.  $0 \leq j \leq N - 1$

New System

1. On Eliminating  $j$  from original system we get:
  - a.  $0 \leq t_2 \leq N - 1$
  - b.  $i = t_1 - t_2$
  - c.  $0 \leq i \leq N - 1$

2. On Eliminating  $i$  from original system we get:

- a.  $0 \leq t_2 \leq N - 1$  (Same as above)
- b.  $0 \leq t_1 - t_2 \leq N - 1$  (New)

Now we can write these 2 inequalities as:

- a.  $t_2 \geq 0$  and  $t_2 \geq t_1 - N + 1$  from first and second inequality respectively.
- b.  $t_2 \leq t_1$  and  $t_2 \leq N - 1$  from first and second inequality respectively.

Hence inner loop bounds are  $\boxed{\max(0, t_1 - N + 1) \leq t_2 \leq \min(t_1, N - 1)}$ . From above equations we can see that  $N \geq 1$ . Since  $t_2 \leq N - 1$  Hence  $0 \leq t_1 \leq 2N - 1$ .

Here is the final code:

```
1 for(int t = 0; t < 2*N - 2; t++) {
2   for(int u = max(0, t-N+1); u <= min(t, N-1); u++) {
3     a[t-u][u] = a[t-u-1][u] + a[t-u][u-1];
4   }
5 }
```

**Example 3:** Consider the following code block:

```
1 for(int i = 0; i < N; i++) {
2   for(int j = 0; j < i; j++) {
3     /* Do Some Work */
4   }
5 }
```

If we apply the transformation  $T(i, j) = (j, i)$  then after fourier motzkin elimination we get the following bounds for new loop variables:  $0 \leq t_1 \leq N - 1$  and  $t_1 \leq t_2 \leq N - 1$ . Hence the new code block will be:

```
1 for(int t = 0; t < N; t++) {
2   for(int u = t; u < N; u++) {
3     /* Do Some Work */
4   }
5 }
```

## Carrying or Satisfying a Dependence

Consider the following code block:

```
1 for(int i = 0; i < N; i++) {  
2     for(int j = 0; j < N; j++) {  
3         for(int k = 0; k < N; k++) {  
4             /* Do Some Work */  
5         }  
6     }  
7 }
```

Suppose we have dependence  $\mathcal{D} = \{(1, 0, 1), (0, 0, 1), (1, 1, 0)\}$ . We say that  $i^{th}$  loop carries the dependence if index of leftmost element of the dependence vector is  $i$ . In this case  $i^{th}$  loop carries the dependence.

We say that a loop is **parallel** if it does not carry any dependence. In above dependences dependence carrying loops are  $i, k, i$  respectively. So in this case for a given  $i$ , the  $j$  loop is parallel.

## March 25, 2025

### Affine Expressions in MLIR

$l(\vec{x})$  is a linear function if it can be written as  $l(\vec{x}) = \mathcal{A}\vec{x}$ .  $l_a(\vec{x})$  is affine expression if it can be written as  $l_a(\vec{x}) = \mathcal{A}\vec{x} + b$ .

An affine map maps zero or more identities to one or more result affine expressions, which can be multi dimensional also. For example:

1. **#map1 = (d0) → (d0-4)**
2. **#map2 = (d0,d1) → (d0+d1, d0-d1)**

## March 27, 2025

Remember that Types never change during execution, every operation strictly requires matching types. We can also define custom types in new dialects.

### Types Relevant for Dense Matrices / Tensors

Add about this later. Checkout slides for details. In particular add about:

1. tensor
2. memref

### Polyhedral Structures in IR

Add about this later. Checkout slides for details. In particular add about:

1. Affine Expressions
2. Affine Maps
3. Integer Sets
4. Affine apply operation
5. Affine If operation
6. Affine For operation

## Transformations for Parallelism

Techniques used for parallelization include **Interchange**: can give rise to outer parallelism which is sync free parallelism, **Loop Skewing**: can expose inner parallelism, **Loop Tiling**. There is something called **Tiling for Parallelism**.

April 3, 2025

Consider the dependence  $(1, 0, 0)$ , here  $j, k$  loops are parallel. We can do the transformation  $T = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$ , this will give us the following transformation when multiplied with the vector  $\begin{bmatrix} i \\ j \\ k \end{bmatrix} = \begin{bmatrix} k \\ j \\ i \end{bmatrix}$ . Similarly  $(j, k, i)$  is also a valid transformation. To decide which one to use, we need to consider other factors like locality.

Let  $m$  be the dimensionality of iteration space, Let  $\mathcal{D}$  be the dependence matrix, then  $m - \text{Rank}(\mathcal{D})$  gives degrees of synchronization free parallelism.

When Outer Loop is sequential and inner loop is parallel, then we call it **Parallelism with synchronization**. Check slide 9 for example.

Add about **do Across** and **do All Parallelism**. Check Slide 33 for reference.

**Tilings:** In General tiling functions map  $n$  dimensional iteration space to  $2n$  dimensional iteration space. Can we have a transformation  $T(i_1, i_2, \dots, i_n) = (i_1, i_2, \dots, i_{2n})$ . Yes, we can have this transformation. Here is a example of rectangular tiling of skewed iteration space:

$$T(i, j) = (\lfloor \frac{i}{64} \rfloor, \lfloor \frac{i+j}{64} \rfloor, i, i+j) . \text{ Check slide 45 for another example which has 3D space.}$$

**Loop Fusion:** Combining two loops into one loop. Opposite of this is called loop distribution. We need to discuss, what are its benefits and when it is valid to do so . Add about **producer-consumer reuse** which is a form of temporal reuse. Add about **Input reuse** which is in the case of sibling fusion. It reduces loop overheads. It also reduces code size.

**Downsides of Fusion:** Fusion can result in more dependences, which can result in loss of useful transformations that could lead to parallelism, tilings and vectorization.

It can lead to more cache misses (conflict misses). Add about this one. Fusion can be bad for capacity misses. Add about this also.

It increases register pressure also. Add more about this one.

It can also decrease the Prefetching performance. Add about this one.

**Validity of Fusion:** Check Slide 47 for fusion.

April 4, 2025