

In [1]: `!usr/local/cuda/bin/nvcc --version`

```
nvcc: NVIDIA (R) Cuda compiler driver
Copyright (c) 2005-2022 NVIDIA Corporation
Built on Wed_Sep_21_10:33:58_PDT_2022
Cuda compilation tools, release 11.8, V11.8.89
Build cuda_11.8.r11.8/compiler.31833905_0
```

In [2]: `!pip install git+https://github.com/andreinechaev/nvcc4jupyter.git`

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Collecting git+https://github.com/andreinechaev/nvcc4jupyter.git
  Cloning https://github.com/andreinechaev/nvcc4jupyter.git to /tmp/pip-req-build-xx11fcf
    Running command git clone --filter=blob:none --quiet https://github.com/andreinechaev/nvcc4jupyter.git /tmp/pip-req-build-xx11fcf
    Resolved https://github.com/andreinechaev/nvcc4jupyter.git to commit aac710a35f52bb78ab34d2e52517237941399eff
    Preparing metadata (setup.py) ... done
Building wheels for collected packages: NVCCPlugin
  Building wheel for NVCCPlugin (setup.py) ... done
  Created wheel for NVCCPlugin: filename=NVCCPlugin-0.0.2-py3-none-any.whl size=4305 sha256=9ac0d74cfea2102bb7cc31936181f5e05a8109d6dd004495e636cb6b0f809d95
    Stored in directory: /tmp/pip-ephem-wheel-cache-x7n3ram1/wheels/a8/b9/18/23f8ef71ceb0f63297dd1903aedd067e6243a68ea756d6feea
Successfully built NVCCPlugin
Installing collected packages: NVCCPlugin
Successfully installed NVCCPlugin-0.0.2
```

In [3]: `%load_ext nvcc_plugin`

```
created output directory at /content/src
Out bin /content/result.out
```

In [4]:

```
%cuda --name vector_add.cu
#include<iostream>
#include<bits/stdc++.h>
#include<cuda.h>

using namespace std;

void fill_array(int *arr,int size){
    for(int i = 0;i < size; i++){
        arr[i] = rand() % 100;
    }
}

void add_cpu(int *arr1, int *arr2, int *result, int size){
    for(int i = 0;i < size; i++){
        result[i] = arr1[i] + arr2[i];
    }
}

void print_matrix(int *arr, int size){
    for(int i = 0; i < size; i++){
        cout << arr[i] << " ";
    }
    cout << endl;
}

__global__ void add(int *arr1, int *arr2, int *arr3){
    int block_id = blockIdx.x;
    arr3[block_id] = arr1[block_id] + arr2[block_id];
}

int main(){
    int *arr1_cpu,*arr2_cpu,*result_cpu;
    int size;
    cout << "Enter size of vector: ";
    cin >> size;

    arr1_cpu = new int[size];
    arr2_cpu = new int[size];
    result_cpu = new int[size];

    fill_array(arr1_cpu,size);
    cout << "Array 1: ";
    print_matrix(arr1_cpu,size);
    fill_array(arr2_cpu,size);
    cout << "Array 2: ";
    print_matrix(arr2_cpu,size);

    int *arr1_gpu,*arr2_gpu,*result_gpu;

    cudaMallocManaged(&arr1_gpu, size * sizeof(int));
    cudaMallocManaged(&arr2_gpu, size * sizeof(int));
    cudaMallocManaged(&result_gpu, size * sizeof(int));

    cudaMemcpy(arr1_gpu,arr1_cpu,size * sizeof(int),cudaMemcpyHostToDevice);
    cudaMemcpy(arr2_gpu,arr2_cpu,size * sizeof(int),cudaMemcpyHostToDevice);
    cudaEvent_t start,stop;
    float elapsedTime;

    cudaEventCreate(&start);
    cudaEventCreate(&stop);
    cudaEventRecord(start,0);

    add<<<size,1>>>(arr1_gpu,arr2_gpu,result_gpu);
    cudaEventRecord(stop,0);
    cudaEventSynchronize(stop);
```

```

    cudaEventElapsedTime(&elapsedTime,start,stop);
    cudaEventDestroy(start);
    cudaEventDestroy(stop);
    cudaMemcpy(result_cpu,result_gpu,size * sizeof(int),cudaMemcpyDeviceToHost);
    cout << "GPU result:\n";
    print_matrix(result_cpu,size);
    cout<<"Elapsed Time = "<<elapsedTime<<" milliseconds" << endl;
    cudaFree(arr1_gpu);
    cudaFree(arr2_gpu);
    cudaFree(result_gpu);

    cudaEventCreate(&start);
    cudaEventCreate(&stop);
    cudaEventRecord(start,0);

    add_cpu(arr1_cpu,arr2_cpu,result_cpu,size);
    cudaEventRecord(stop,0);
    cudaEventSynchronize(stop);
    cudaEventElapsedTime(&elapsedTime,start,stop);
    cudaEventDestroy(start);
    cudaEventDestroy(stop);
    cout << "CPU result:\n";
    print_matrix(result_cpu,size);
    cout<<"Elapsed Time = "<<elapsedTime<<" milliseconds" << endl;

    return 0;
}

```

Out[4]: 'File written in /content/src/vector_add.cu'

In [5]: `!nvcc /content/src/vector_add.cu -o /content/src/vector_add`

In []: `!/content/src/vector_add`

```

%%cuda --name matrix_multiply.cu
#include<iostream>
#include<bits/stdc++.h>
#include<cuda.h>

using namespace std;

void initialize_matrix(int *array, int rows, int cols){
    for(int i = 0 ; i < rows; i++){
        for(int j = 0; j < cols; j++){
            array[i*cols + j] = rand() % 10;
        }
    }
}

void print_matrix(int *array, int rows, int cols){
    for(int i = 0 ; i < rows; i++){
        for(int j = 0; j < cols; j++){
            cout << array[i*cols + j] << " ";
        }
        cout << endl;
    }
}

void matrix_multiplication_cpu(int *a, int *b, int *c, int common, int c_rows,int c_cols){
    for(int i = 0; i < c_rows; i++){
        for(int j = 0; j < c_cols; j++){
            int sum = 0;
            for(int k = 0; k < common; k++){
                sum += a[i*common + k] * b[k*c_cols + j];
            }
            c[i*c_cols + j] = sum;
        }
    }
}

__global__
void matrix_multiply(int *a, int *b, int *c, int m, int n, int k)
{
    int row = blockIdx.y*blockDim.y + threadIdx.y;
    int col = blockIdx.x*blockDim.x + threadIdx.x;
    int sum=0;

    if(col<k && row<m) {
        for(int j=0;j<n;j++){
            {
                sum += a[row*n+j] * b[j*k+col];
            }
            c[k*row+col]=sum;
        }
    }
}

int main(){

    int A_rows, A_cols, B_rows, B_cols, C_rows, C_cols;

```

```

    cout << "Dimensions of matrix 1:\n";
    cout << "Rows: ";
    cin >> A_rows;
    cout << "Columns: ";
    cin >> A_cols;
    cout << "Dimensions of matrix 2:\n";
    cout << "Rows: " << A_cols << endl << "Columns: ";
    cin >> B_cols;
    B_rows = A_cols;
    C_rows = A_rows;
    C_cols = B_cols;

    int A_size = A_rows * A_cols;
    int B_size = B_rows * B_cols;
    int C_size = C_rows * C_cols;

    int *A, *B, *C;
    int *m1,*m2,*result;

    A = new int[A_size];
    B = new int[B_size];
    C = new int[C_size];

    initialize_matrix(A,A_rows,A_cols);
    cout << "Matrix 1\n";
    print_matrix(A,A_rows,A_cols);
    initialize_matrix(B,B_rows,B_cols);
    cout << "Matrix 2\n";
    print_matrix(B,B_rows,B_cols);

    cudaMallocManaged(&m1, A_size * sizeof(int));
    cudaMallocManaged(&m2, B_size * sizeof(int));
    cudaMallocManaged(&result, C_size * sizeof(int));

    cudaMemcpy(m1,A,A_size * sizeof(int), cudaMemcpyHostToDevice);
    cudaMemcpy(m2,B,B_size * sizeof(int), cudaMemcpyHostToDevice);

    dim3 dimGrid(2,2);
    dim3 dimBlock(C_rows,C_cols);

    float gpu_elapsed_time;
    cudaEvent_t gpu_start,gpu_stop;

    cudaEventCreate(&gpu_start);
    cudaEventCreate(&gpu_stop);
    cudaEventRecord(gpu_start);
    matrix_multiply<<dimGrid,dimBlock>>>(m1,m2,result,C_rows,A_cols,C_cols);
    cudaEventRecord(gpu_stop);
    cudaEventSynchronize(gpu_stop);
    cudaEventElapsedTime(&gpu_elapsed_time, gpu_start, gpu_stop);
    cudaEventDestroy(gpu_start);
    cudaEventDestroy(gpu_stop);

    cudaMemcpy(C,result,C_size*sizeof(int),cudaMemcpyDeviceToHost);
    cout << "GPU result:\n";
    print_matrix(C,C_rows,C_cols);
    cout<<"GPU Elapsed time is: "<<gpu_elapsed_time<<" milliseconds"<<endl;

    cudaEventCreate(&gpu_start);
    cudaEventCreate(&gpu_stop);
    cudaEventRecord(gpu_start);
    matrix_multiplication_cpu(A,B,C,A_cols,C_rows,C_cols);
    cudaEventRecord(gpu_stop);
    cudaEventSynchronize(gpu_stop);
    cudaEventElapsedTime(&gpu_elapsed_time, gpu_start, gpu_stop);
    cudaEventDestroy(gpu_start);
    cudaEventDestroy(gpu_stop);

    cout << "CPU result:\n";
    print_matrix(C,C_rows,C_cols);
    cout<<"CPU Elapsed time is: "<<gpu_elapsed_time<<" milliseconds"<<endl;

    cudaFree(m1);
    cudaFree(m2);
    cudaFree(result);

    return 0;
}

```

Out[7]: 'File written in /content/src/matrix_multiply.cu'

In [8]: `!nvcc /content/src/matrix_multiply.cu -o /content/src/matrix_multiply`

In [10]: `!./content/src/matrix_multiply`

```
Dimensions of matrix 1:
Rows: 10
Columns: 8
Dimensions of matrix 2:
Rows: 8
Columns: 12
Matrix 1
3 6 7 5 3 5 6 2
9 1 2 7 0 9 3 6
0 6 2 6 1 8 7 9
2 0 2 3 7 5 9 2
2 8 9 7 3 6 1 2
9 3 1 9 4 7 8 4
5 0 3 6 1 0 6 3
2 0 6 1 5 5 4 7
6 5 6 9 3 7 4 5
2 5 4 7 4 4 3 0
Matrix 2
7 8 6 8 8 4 3 1 4 9 2 0
6 8 9 2 6 6 4 9 5 0 4 8
7 1 7 2 7 2 2 6 1 0 6 1
5 9 4 9 0 9 1 7 7 1 1 5
9 7 7 6 7 3 6 5 6 3 9 4
8 1 2 9 3 9 0 8 8 5 0 9
6 3 8 5 6 1 1 5 9 8 4 8
1 0 3 0 4 4 4 4 7 6 3 1
GPU result:
236 168 226 188 189 175 84 227 210 126 134 187
214 163 165 237 161 217 69 190 233 193 71 156
204 140 198 183 159 212 83 248 270 159 106 221
202 126 175 179 156 122 72 163 208 151 124 164
243 182 222 190 180 209 90 253 195 84 134 184
277 237 242 288 210 236 98 237 295 225 125 217
134 122 139 136 116 101 51 112 146 120 76 88
177 83 156 132 160 121 81 158 176 133 123 115
271 215 241 252 206 243 101 261 263 169 135 203
193 164 173 172 132 160 68 187 169 81 103 155
GPU Elapsed time is: 0.01104 milliseconds
CPU result:
236 168 226 188 189 175 84 227 210 126 134 187
214 163 165 237 161 217 69 190 233 193 71 156
204 140 198 183 159 212 83 248 270 159 106 221
202 126 175 179 156 122 72 163 208 151 124 164
243 182 222 190 180 209 90 253 195 84 134 184
277 237 242 288 210 236 98 237 295 225 125 217
134 122 139 136 116 101 51 112 146 120 76 88
177 83 156 132 160 121 81 158 176 133 123 115
271 215 241 252 206 243 101 261 263 169 135 203
193 164 173 172 132 160 68 187 169 81 103 155
CPU Elapsed time is: 0.007648 milliseconds
```

In [9]: