



Example of CUDA processing flow

1. Copy data from main memory to GPU memory
2. CPU initiates the GPU **compute kernel**
3. GPU's CUDA cores execute the kernel in parallel
4. Copy the resulting data from GPU memory to main memory

```

1 #include<stdio.h>
2 #include<stdlib.h>
3
4 // CUDA kernel: compute vector sum z_d = x_d + y_d using a grid of threads on GPU
5 __global__ void vecAddKernel(float *x_d, float *y_d, float *z_d, unsigned int n){
6     unsigned int i = blockDim.x * blockIdx.x + threadIdx.x;
7     if(i < n){
8         z_d[i] = x_d[i]+y_d[i];
9     }
10 }
11
12 int main(int argc, char** argv){
13     unsigned int n = 1024;
14
15     // Allocate host memory for arrays x_h, y_h, and z_h; and initialize arrays x_h and y_h.
16     float *x_h = (float*) malloc(sizeof(float)*n);
17     for(unsigned int i=0; i<n; i++) x_h[i] = (float)rand()/(float)(RAND_MAX);
18     float *y_h = (float*) malloc(sizeof(float)*n);
19     for(unsigned int i=0; i<n; i++) y_h[i] = (float)rand()/(float)(RAND_MAX);
20     float *z_h = (float*) calloc(n, sizeof(float));
21
22     // (1) Allocate device memory for arrays x_d, y_d, and z_d.
23     float *x_d, *y_d, *z_d;
24     cudaMalloc((void**)&x_d, sizeof(float)*n);
25     cudaMalloc((void**)&y_d, sizeof(float)*n);
26     cudaMalloc((void**)&z_d, sizeof(float)*n);
27
28     // (2) Copy arrays x_h and z_h to device memory x_d and y_d, respectively.
29     cudaMemcpy(x_d, x_h, sizeof(float)*n, cudaMemcpyHostToDevice);
30     cudaMemcpy(y_d, y_h, sizeof(float)*n, cudaMemcpyHostToDevice);
31
32     // (3) Call kernel to launch a grid of threads to perform the vector addition on GPU.
33     vecAddKernel<<<<<<(n/256,0), 256>>>>(x_d,y_d,z_d,n);
34
35     // (4) Copy the result data from the device memory of array z_d to the host memory of array z_h.
36     cudaMemcpy(z_h, z_d, sizeof(float)*n, cudaMemcpyDeviceToHost);
37
38     // (5) Free device memory of arrays x_d, y_d, and z_d
39     cudaFree(x_d);
40     cudaFree(y_d);
41     cudaFree(z_d);
42
43     // Free host memory of arrays x_h, y_h, and z_h
44     free(x_h);
45     free(y_h);
46     free(z_h);
47
48     return 0;
49 }

```

Define a CUDA Kernel Function