# Parallel Programming with CUDA

*Streltsova O., Zuev M.*

Laboratory of Information Technologies

Joint Institute for Nuclear Research

HETEROGENEOUS COMPUTATIONS GROUP, HybriLIT

# cuBLAS library: cublas<t>scal()

The cuBLAS library is an implementation of BLAS (Basic Linear Algebra Subprograms) on top of the NVIDIA®CUDA runtime.
**Site:**

> **http://docs.nvidia.com/cuda/cublas/**

- cublasStatus_t cublasSscal(cublasHandle_t   handle, int n,
  const float *alpha, float *x, int incx)
- cublasStatus_t cublasDscal(cublasHandle_t handle, int n,
  const double *alpha, double *x, int incx)
- cublasStatus_t cublasCscal(cublasHandle_t handle, int n,
  const cuComplex *alpha, cuComplex *x, int incx)
- cublasStatus_t cublasCsscal(cublasHandle_t handle, int n,
  const float *alpha, cuComplex *x, int incx)
- cublasStatus_t cublasZscal(cublasHandle_t handle, int n,
  const cuDoubleComplex *alpha,
  cuDoubleComplex *x, int incx)
- cublasStatus_t cublasZdscal(cublasHandle_t handle, int n,
  const double *alpha, cuDoubleComplex *x, int incx)

# Program example

## Location

**cd Tutorial_HSchool2014/CUDA**

## File

**cuda_blas.cu**

# Program example

```
1.  #include <stdio.h>
2.  #include <stdlib.h>
3.  #include <cublas_v2.h>  // API v2 from CUDA 4.0
4.
5.  #include <math.h>
6.  #include <cuda.h>
7.
8.  #define NX 16384 // 64
9.
10. //----------------------- Error Handling  ----------------------------//
11. #define CUDA_CALL(x) do { cudaError_t err = x; if (( err ) != cudaSuccess ){ \
12. printf ("Error \"%s\" at %s :%d \n" , cudaGetErrorString(err), \
13. __FILE__ , __LINE__ ) ; return -1;\
14. }} while (0);
15.
16. #define CUBLAS_CALL(x) do { if (( x ) != CUBLAS_STATUS_SUCCESS ) {\
17. printf ("Error at %s :%d \n" , __FILE__ , __LINE__ ) ;\
18. return -1; }} while (0);
```

# Program example

```
19. int main() {
20.    cublasHandle_t handleNorm;
21.    cublasCreate(&handleNorm) ;
22.
23.    printf(" ==== Element-wise  vector  multiplication  by  a  scalar ==== \n" );
24.    printf(" Array  size  NX  =  %d\n", NX);
25.
26.    //-------------------- Input data: allocate memory on CPU --------------------//
27.    double* A_input = new double[NX];
28.    double* A_out    = new double[NX];
29.
30.    for (int ix= 0; ix< NX; ix++){
31.        A_input[ix]= (double)(ix);
32.    }
33.
34.    double normScalar= 3.0*M_PI;
```

# Program example

```
35.   //----------------- Allocate  memory on device -------------------//
36.   double* dev_Ainput;
37.   CUDA_CALL( cudaMalloc((void**)&dev_Ainput, NX*sizeof(double)) );
```

```
cudaMalloc ( void ** devPtr, // Pointer to allocated memory
             size_t  size )   // Requested allocation size in bytes
```

```
...
Error "out of memory" at cuda_blas.cu :37
```

```
38.   //------------------- Timing a CUDA application using events -------------------//
39.   cudaEvent_t start, stop_ini,  stop_cublas, stop_cpu;
40.   cudaEventCreate(&start);
41.   cudaEventCreate(&stop_ini);
42.   cudaEventCreate(&stop_cublas);
43.   cudaEventCreate(&stop_cpu);
44.
45.   cudaEventRecord(start);
```

# Program example

```
46.  //------------------- Copy from Host to Device Input data -------------------//
47.  CUDA_CALL( cudaMemcpy(dev_Ainput , A_input, NX*sizeof(double),
                                cudaMemcpyHostToDevice) );
```

```
cudaMemcpy ( void *              dst,   // Destination
             const void *        src,   // Source
             size_t              count, // Size
             enum cudaMemcpyKind kind ) // Type of transfer

cudaMemcpyKind: cudaMemcpyHostToHost, cudaMemcpyHostToDevice,
cudaMemcpyDeviceToHost, cudaMemcpyDeviceToDevice
```

# Program example

```
49.    //------------------ Calculate time ------------------//
50.    cudaEventRecord(stop_ini, 0);
51.    float time_ini = 0.0;
52.    cudaEventSynchronize(stop_ini);
53.    cudaEventElapsedTime(&time_ini, start, stop_ini);
54.    printf("GPU compute time_load data(msec): %.5f\n", time_ini);
55.
56.    cudaError_t error = cudaGetLastError();
57.    if(error != cudaSuccess){
58.       printf("CUDA error, CUBLAS:  %s\n", cudaGetErrorString(error));
59.       exit(-1);
60.    }
61.
62.    //------------------ Check error ------------------//
63.    cudaError_t error = cudaGetLastError();
64.    if(error != cudaSuccess){
65.       printf("CUDA error, CUFFT:  %s\n", cudaGetErrorString(error));
66.       exit(-1);
67.    }
```

# Program example

```
68.  //---------------- For multiplication use function of CUBLAS library ------------------//
69.  CUBLAS_CALL( cublasDscal( handleNorm, NX, &normScalar, &dev_Ainput[0], 1) );
```

```
cublasDscal( cublasHandle_t  handle,   // Handle to the cuBLAS library context
             int             n,        // Number of elements in the vector x
             const double *  alpha,    // Scalar used for multiplication
             double *        x,        // Elements in the vector x
             int             incx )    // Stride between elements of x
```

```
70.  //------------------ Timing -------------------//
71.  cudaEventRecord(stop_cublas);
72.  float time_cublas = 0.0;
73.  cudaEventSynchronize(stop_cublas);
74.  cudaEventElapsedTime(&time_cublas, stop_ini, stop_cublas);
75.  printf("GPU compute time_cublas data(msec): %.5f\n", time_cublas);
```

# Program example

```
76.   //------------------ CPU computation ------------------//
77.   for (int ix= 0; ix< NX; ix++)
78.       A_out[ix]= A_input[ix]*normScalar;
79.
80.   cudaEventRecord(stop_cpu);
81.   float time_cpu = 0.0;
82.   cudaEventSynchronize(stop_cpu);
83.   cudaEventElapsedTime(&time_cpu, stop_cublas,stop_cpu);
84.   printf("CPU compute time_cpu (msec): %.5f\n", time_cpu);
85.
86.   //------------------ Copy from Device to Host Output data ------------------//
87.   CUDA_CALL( cudaMemcpy(A_out, dev_Ainput, NX*sizeof(double),
                                     cudaMemcpyDeviceToHost) );
87.   cudaDeviceSynchronize();
```

# Program example

```
88.    //------------------- Output result  on screen and file -------------------//
89.    char* file_name = "Rezult_GPU.dat" ;
90.    printf("============ Test result Vector for [0] and [1] elements \n");
91.    printf(  "ix=0, iy =0, A_out = %.16e, A_input = %.16e  \n", A_out[0], A_input[0]);
92.    printf(  "ix=1, iy =0, A_out = %.16e, A_input = %.16e  \n", A_out[1], A_input[1]);
93.
94.    FILE *fp4;
95.    fp4=fopen(file_name, "w");
96.    fprintf(fp4, "ix   ,   A_out[ix]),   A_input[ix]" );
97.    for (int ix = 0; ix < NX; ix++)
98.       fprintf(fp4, " %d    %.16e   %.16e  \n", ix, A_out[ix], A_input[ix]);
99.
100.   fclose(fp4);
101.   cublasDestroy(handleNorm) ;
102.   cudaFree(dev_Ainput);
103.   delete[] A_input;
104.   delete[] A_out;
105.   return 0;
106.}
```

# Program example. Compilation

## Add module

```
> module add cuda-6.0-x86_64
```

## Compilation with libraries

```
> nvcc -gencode=arch=compute_35,code=sm_35
-lcublas -O3 cublas_ex.cu -o cublas
```

# Program example. Running

Listing of **_script_cuda_**

```
#!/bin/sh
#SBATCH -p gpu
./test1
```

Running in batch

```
> sbatch script_cuda
```

# Program example. Output on the screen

```
 ==== Element-wise vector multiplication by a scalar =====
Array size  NX = 1048576
GPU compute time_load data (msec): 5.27184
GPU compute time_cublas data (msec): 0.22720
CPU compute time_cpu (msec): 6.57942
=============  Test result done on file Rezult_GPU.dat
=============  Test result Vector for [0] and [1] elements
ix=0, iy =0, A_out = 0.0000000000000000e+00,
A_input = 0.0000000000000000e+00
ix=1, iy =0, A_out = 9.4247779607693793e+00,
A_input = 1.0000000000000000e+00
```