



CUDA Toolkit 5.0 CUBLAS Library

PG-05326-050_v01 | April 2012



Contents

1	Introduction	4
1.1	Data layout	4
1.2	New and Legacy CUBLAS API	5
1.3	Example code	6
2	Using the CUBLAS API	9
2.1	Error status	9
2.2	CUBLAS context	9
2.3	Thread Safety	10
2.4	Scalar Parameters	10
2.5	Parallelism with Streams	10
2.6	Batching Kernels	11
2.7	Device API	11
3	CUBLAS Datatypes Reference	12
3.1	cublasHandle_t	12
3.2	cublasStatus_t	12
3.3	cublasOperation_t	13
3.4	cublasFillMode_t	14
3.5	cublasDiagType_t	14
3.6	cublasSideMode_t	14
3.7	cublasPointerMode_t	14
3.8	cublasAtomsicsMode_t	15
4	CUBLAS Helper Function Reference	16
4.1	cublasCreate()	16
4.2	cublasDestroy()	16
4.3	cublasGetVersion()	16
4.4	cublasSetStream()	17
4.5	cublasGetStream()	17
4.6	cublasGetPointerMode()	17
4.7	cublasSetPointerMode()	18
4.8	cublasSetVector()	18

4.9	<code>cublasGetVector()</code>	18
4.10	<code>cublasSetMatrix()</code>	19
4.11	<code>cublasGetMatrix()</code>	19
4.12	<code>cublasSetVectorAsync()</code>	20
4.13	<code>cublasGetVectorAsync()</code>	20
4.14	<code>cublasSetMatrixAsync()</code>	21
4.15	<code>cublasGetMatrixAsync()</code>	21
4.16	<code>cublasGetMatrixAsync()</code>	21
4.17	<code>cublasSetAtomicsMode()</code>	22
4.18	<code>cublasGetAtomicsMode()</code>	22
5	CUBLAS Level-1 Function Reference	23
5.1	<code>cublasI<t>amax()</code>	23
5.2	<code>cublasI<t>amin()</code>	24
5.3	<code>cublas<t>asum()</code>	25
5.4	<code>cublas<t>axpy()</code>	26
5.5	<code>cublas<t>copy()</code>	27
5.6	<code>cublas<t>dot()</code>	27
5.7	<code>cublas<t>nrm2()</code>	29
5.8	<code>cublas<t>rot()</code>	29
5.9	<code>cublas<t>rotg()</code>	31
5.10	<code>cublas<t>rotm()</code>	32
5.11	<code>cublas<t>rotmg()</code>	33
5.12	<code>cublas<t>scal()</code>	34
5.13	<code>cublas<t>swap()</code>	35
6	CUBLAS Level-2 Function Reference	37
6.1	<code>cublas<t>gbmv()</code>	37
6.2	<code>cublas<t>gemv()</code>	39
6.3	<code>cublas<t>ger()</code>	40
6.4	<code>cublas<t>sbmv()</code>	42
6.5	<code>cublas<t>spmv()</code>	44
6.6	<code>cublas<t>spr()</code>	45
6.7	<code>cublas<t>spr2()</code>	46
6.8	<code>cublas<t>symv()</code>	47
6.9	<code>cublas<t>syr()</code>	49
6.10	<code>cublas<t>syr2()</code>	50
6.11	<code>cublas<t>tbbmv()</code>	51
6.12	<code>cublas<t>tbsv()</code>	53
6.13	<code>cublas<t>tpmv()</code>	55
6.14	<code>cublas<t>tpsv()</code>	56
6.15	<code>cublas<t>trmv()</code>	58
6.16	<code>cublas<t>trsv()</code>	59
6.17	<code>cublas<t>hemv()</code>	60

6.18	cublas<t>hbm _v ()	61
6.19	cublas<t>hpm _v ()	63
6.20	cublas<t>her()	65
6.21	cublas<t>her ₂ ()	66
6.22	cublas<t>hpr()	67
6.23	cublas<t>hpr ₂ ()	68
7	CUBLAS Level-3 Function Reference	70
7.1	cublas<t>gemm()	70
7.2	cublas<t>gemmBatched()	72
7.3	cublas<t>symm()	74
7.4	cublas<t>syrk()	76
7.5	cublas<t>syr ₂ k()	78
7.6	cublas<t>trmm()	80
7.7	cublas<t>trsm()	83
7.8	cublas<t>hemm()	84
7.9	cublas<t>herk()	87
7.10	cublas<t>her ₂ k()	88
8	BLAS-like Extension	91
8.1	cublas<t>geam()	91
8.2	cublas<t>dgmm()	93
9	Appendix A: Using the CUBLAS Legacy API	95
9.1	Error Status	95
9.2	Initialization and Shutdown	95
9.3	Thread Safety	96
9.4	Memory Management	96
9.5	Scalar Parameters	96
9.6	Helper Functions	96
9.7	Level-1,2,3 Functions	97
9.8	Converting Legacy to the CUBLAS API	97
9.9	Examples	98
10	Appendix B: CUBLAS Fortran Bindings	101

Published by
NVIDIA Corporation
2701 San Tomas Expressway
Santa Clara, CA 95050

Notice

ALL NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, "MATERIALS") ARE BEING PROVIDED "AS IS". NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE.

Information furnished is believed to be accurate and reliable. However, NVIDIA Corporation assumes no responsibility for the consequences of use of such information or for any infringement of patents or other rights of third parties that may result from its use. No license is granted by implication or otherwise under any patent or patent rights of NVIDIA Corporation. Specifications mentioned in this publication are subject to change without notice. This publication supersedes and replaces all information previously supplied. NVIDIA Corporation products are not authorized for use as critical components in life support devices or systems without express written approval of NVIDIA Corporation.

Trademarks

NVIDIA, CUDA, and the NVIDIA logo are trademarks or registered trademarks of NVIDIA Corporation in the United States and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

Copyright

Copyright ©2005-2012 by NVIDIA Corporation. All rights reserved.

Portions of the SGEMM, DGEMM, CGEMM and ZGEMM library routines were written by Vasily Volkov and are subject to the Modified Berkeley Software Distribution License as follows:

Copyright ©2007-2009, Regents of the University of California

All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- ▶ Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- ▶ Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.

- Neither the name of the University of California, Berkeley nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE AUTHOR "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE AUTHOR BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

Portions of the SGEMM, DGEMM and ZGEMM library routines were written by Davide Barbieri and are subject to the Modified Berkeley Software Distribution License as follows:

Copyright ©2008-2009 Davide Barbieri @ University of Rome Tor Vergata.

All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- The name of the author may not be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE AUTHOR "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE AUTHOR BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

Portions of the DGEMM and SGEMM library routines optimized for Fermi architecture were developed by the University of Tennessee. Subsequently, several other routines that are optimized for the Fermi architecture have been derived from these initial DGEMM and SGEMM implementations. Those portions of the source code are thus subject to the Modified Berkeley Software Distribution License as follows:

Copyright ©2010 The University of Tennessee. All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- ▶ Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- ▶ Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer listed in this license in the documentation and/or other materials provided with the distribution.
- ▶ Neither the name of the copyright holders nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

Chapter 1

Introduction

The CUBLAS library is an implementation of BLAS (Basic Linear Algebra Subprograms) on top of the NVIDIA®CUDA™ runtime. It allows the user to access the computational resources of NVIDIA Graphics Processing Unit (GPU), but does not auto-parallelize across multiple GPUs.

To use the CUBLAS library, the application must allocate the required matrices and vectors in the GPU memory space, fill them with data, call the sequence of desired CUBLAS functions, and then upload the results from the GPU memory space back to the host. The CUBLAS library also provides helper functions for writing and retrieving data from the GPU.

1.1 Data layout

For maximum compatibility with existing Fortran environments, the CUBLAS library uses column-major storage, and 1-based indexing. Since C and C++ use row-major storage, applications written in these languages can not use the native array semantics for two-dimensional arrays. Instead, macros or inline functions should be defined to implement matrices on top of one-dimensional arrays. For Fortran code ported to C in mechanical fashion, one may choose to retain 1-based indexing to avoid the need to transform loops. In this case, the array index of a matrix element in row “i” and column “j” can be computed via the following macro

```
#define IDX2F(i , j , ld) (((j)-1)*(ld))+((i)-1))
```

Here, ld refers to the leading dimension of the matrix, which in the case of column-major storage is the number of rows of the allocated matrix (even if only a submatrix of it is being used). For natively written C and C++ code, one would most likely choose 0-based indexing, in which case the array index of a matrix element in row “i” and column “j” can be computed via the following macro

```
#define IDX2C(i , j , ld) (((j)*(ld))+i))
```


1.2 New and Legacy CUBLAS API

Starting with version 4.0, the CUBLAS Library provides a new updated API, in addition to the existing legacy API. This section discusses why a new API is provided, the advantages of using it, and the differences with the existing legacy API.

The new CUBLAS library API can be used by including the header file “cublas_v2.h”. It has the following features that the legacy CUBLAS API does not have:

- ▶ the `handle` to the CUBLAS library context is initialized using the `cublasCreate` function and is explicitly passed to every subsequent library function call. This allows the user to have more control over the library setup when using multiple host threads and multiple GPUs. This also allows the CUBLAS APIs to be reentrant.
- ▶ the scalars α and β can be passed by reference on the host or the device, instead of only being allowed to be passed by value on the host. This change allows library functions to execute asynchronously using streams even when α and β are generated by a previous kernel.
- ▶ when a library routine returns a scalar result, it can be returned by reference on the host or the device, instead of only being allowed to be returned by value only on the host. This change allows library routines to be called asynchronously when the scalar result is generated and returned by reference on the device resulting in maximum parallelism.
- ▶ the error status `cublasStatus_t` is returned by all CUBLAS library function calls. This change facilitates debugging and simplifies software development. Note that `cublasStatus` was renamed `cublasStatus_t` to be more consistent with other types in the CUBLAS library.
- ▶ the `cublasAlloc()` and `cublasFree()` functions have been deprecated. This change removes these unnecessary wrappers around `cudaMalloc()` and `cudaFree()`, respectively.
- ▶ the function `cublasSetKernelStream()` was renamed `cublasSetStream()` to be more consistent with the other CUDA libraries.

The legacy CUBLAS API, explained in more detail in the Appendix A, can be used by including the header file “cublas.h”. Since the legacy API is identical to the previously released CUBLAS library API, existing applications will work out of the box and automatically use this legacy API without any source code changes. In general, new applications should not use the legacy CUBLAS API, and existing applications should convert to using the new API if it requires sophisticated and optimal stream parallelism or if it calls CUBLAS routines concurrently from multiple threads. For the rest of the document, the new CUBLAS Library API will simply be referred to as the CUBLAS Library API.

As mentioned earlier the interfaces to the legacy and the CUBLAS library APIs are the header file “cublas.h” and “cublas_v2.h”, respectively. In addition, applications using the

CUBLAS library need to link against the DSO cublas.so (Linux), the DLL cublas.dll (Windows), or the dynamic library cublas.dylib (Mac OS X). Note: the same dynamic library implements both the new and legacy CUBLAS APIs.

1.3 Example code

For sample code references please see the two examples below. They show an application written in C using the CUBLAS library API with two indexing styles (Example 1.

"Application Using C and CUBLAS: 1-based indexing" and Example 2. "Application Using C and CUBLAS: 0-based Indexing").

```
//Example 1. Application Using C and CUBLAS: 1-based indexing
//-----
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#include <cuda_runtime.h>
#include "cublas_v2.h"
#define M 6
#define N 5
#define IDX2F(i,j,ld) (((j)-1)*(ld))+((i)-1))

static __inline__ void modify (cublasHandle_t handle, float *m, int ldm, int n, int ←
    p, int q, float alpha, float beta){
    cublasSscal (handle, n-p+1, &alpha, &m[IDX2F(p,q,ldm)], ldm);
    cublasSscal (handle, ldm-p+1, &beta, &m[IDX2F(p,q,ldm)], 1);
}

int main (void){
    cudaError_t cudaStat;
    cublasStatus_t stat;
    cublasHandle_t handle;
    int i, j;
    float * devPtrA;
    float * a = 0;
    a = (float *)malloc (M * N * sizeof (*a));
    if (!a) {
        printf ("host memory allocation failed");
        return EXIT_FAILURE;
    }
    for (j = 1; j <= N; j++) {
        for (i = 1; i <= M; i++) {
            a[IDX2F(i,j,M)] = (float)((i-1) * M + j);
        }
    }
    cudaStat = cudaMalloc ((void**)&devPtrA, M*N*sizeof(*a));
    if (cudaStat != cudaSuccess) {
        printf ("device memory allocation failed");
        return EXIT_FAILURE;
    }
    stat = cublasCreate(&handle);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("CUBLAS initialization failed\n");
        return EXIT_FAILURE;
    }
}
```

```

    stat = cublasSetMatrix (M, N, sizeof(*a), a, M, devPtrA, M);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("data download failed");
        cudaFree (devPtrA);
        cublasDestroy(handle);
        return EXIT_FAILURE;
    }
    modify (handle, devPtrA, M, N, 2, 3, 16.0f, 12.0f);
    stat = cublasGetMatrix (M, N, sizeof(*a), devPtrA, M, a, M);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("data upload failed");
        cudaFree (devPtrA);
        cublasDestroy(handle);
        return EXIT_FAILURE;
    }
    cudaFree (devPtrA);
    cublasDestroy(handle);
    for (j = 1; j <= N; j++) {
        for (i = 1; i <= M; i++) {
            printf ("%7.0f", a[IDX2F(i,j,M)]);
        }
        printf ("\n");
    }
    return EXIT_SUCCESS;
}

```

```

//Example 2. Application Using C and CUBLAS: 0-based indexing
//-----
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#include <cuda_runtime.h>
#include "cublas_v2.h"
#define M 6
#define N 5
#define IDX2C(i,j,ld) (((j)*(ld))+(i))

static __inline__ void modify (cublasHandle_t handle, float *m, int ldm, int n, int←
    p, int q, float alpha, float beta){
    cublasSscal (handle, n-p, &alpha, &m[IDX2C(p,q,ldm)], ldm);
    cublasSscal (handle, ldm-p, &beta, &m[IDX2C(p,q,ldm)], 1);
}

int main (void){
    cudaError_t cudaStat;
    cublasStatus_t stat;
    cublasHandle_t handle;
    int i, j;
    float* devPtrA;
    float* a = 0;
    a = (float *)malloc (M * N * sizeof (*a));
    if (!a) {
        printf ("host memory allocation failed");
        return EXIT_FAILURE;
    }
    for (j = 0; j < N; j++) {
        for (i = 0; i < M; i++) {
            a[IDX2C(i,j,M)] = (float)(i * M + j + 1);
        }
    }
}

```

```

    cudaStat = cudaMalloc ((void**)&devPtrA, M*N*sizeof(*a));
    if (cudaStat != cudaSuccess) {
        printf ("device memory allocation failed");
        return EXIT_FAILURE;
    }
    stat = cublasCreate(&handle);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("CUBLAS initialization failed\n");
        return EXIT_FAILURE;
    }
    stat = cublasSetMatrix (M, N, sizeof(*a), a, M, devPtrA, M);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("data download failed");
        cudaFree (devPtrA);
        cublasDestroy(handle);
        return EXIT_FAILURE;
    }
    modify (handle, devPtrA, M, N, 1, 2, 16.0f, 12.0f);
    stat = cublasGetMatrix (M, N, sizeof(*a), devPtrA, M, a, M);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("data upload failed");
        cudaFree (devPtrA);
        cublasDestroy(handle);
        return EXIT_FAILURE;
    }
    cudaFree (devPtrA);
    cublasDestroy(handle);
    for (j = 0; j < N; j++) {
        for (i = 0; i < M; i++) {
            printf ("%7.0f", a[IDX2C(i,j,M)]);
        }
        printf ("\n");
    }
    return EXIT_SUCCESS;
}

```

Chapter 2

Using the CUBLAS API

This section describes how to use the CUBLAS library API. It does not contain a detailed reference for all API datatypes and functions—those are provided in subsequent chapters. The Legacy CUBLAS API is also not covered in this section—that is handled in an Appendix.

2.1 Error status

All CUBLAS library function calls return the error status `cublasStatus_t`.

2.2 CUBLAS context

The application must initialize the `handle` to the CUBLAS library context by calling the `cublasCreate()` function. Then, the `handle` is explicitly passed to every subsequent library function call. Once the application finishes using the library, it must call the `cublasDestroy()` function to release the resources associated with the CUBLAS library context.

This approach allows the user to explicitly control the library setup when using multiple host threads and multiple GPUs. For example, the application can use `cudaSetDevice()` to associate different devices with different host threads and in each of those host threads it can initialize a unique `handle` to the CUBLAS library context, which will use the particular device associated with that host thread. Then, the CUBLAS library function calls made with different `handle` will automatically dispatch the computation to different devices.

The device associated with a particular CUBLAS context is assumed to remain unchanged between the corresponding `cublasCreate()` and `cublasDestroy()` calls. In order for the CUBLAS library to use a different device in the same host thread, the application must set the new device to be used by calling `cudaSetDevice()` and then create another CUBLAS context, which will be associated with the new device, by calling `cublasCreate()`.

2.3 Thread Safety

The library is thread safe and its functions can be called from multiple host threads, even with the same `handle`.

2.4 Scalar Parameters

In the CUBLAS API the scalar parameters α and β can be passed by reference on the host or the device

Also, the few functions that return a scalar result, such as `amax()`, `amin()`, `asum()`, `rotg()`, `rotmg()`, `dot()` and `nrm2()`, return the resulting value by reference on the host or the device. Notice that even though these functions return immediately, similarly to matrix and vector results, the scalar result is ready only when execution of the routine on the GPU completes. This requires proper synchronization in order to read the result from the host.

These changes allow the library functions to execute completely asynchronously using streams even when α and β are generated by a previous kernel. For example, this situation can arise when iterative methods for solution of linear systems and eigenvalue problems are implemented using the CUBLAS library.

2.5 Parallelism with Streams

If the application uses the results computed by multiple independent tasks, CUDA™ streams can be used to overlap the computation performed in these tasks.

The application can conceptually associate each stream with each task. In order to achieve the overlap of computation between the tasks, the user should create CUDA™ streams using the function `cudaStreamCreate()` and set the stream to be used by each individual CUBLAS library routine by calling `cublasSetStream()` just before calling the actual CUBLAS routine. Then, the computation performed in separate streams would be overlapped automatically when possible on the GPU. This approach is especially useful when the computation performed by a single task is relatively small and is not enough to fill the GPU with work.

We recommend using the new CUBLAS API with scalar parameters and results passed by reference in the device memory to achieve maximum overlap of the computation when using streams.

A particular application of streams, batching of multiple small kernels, is described below.

2.6 Batching Kernels

In this section we will explain how to use streams to batch the execution of small kernels. For instance, suppose that we have an application where we need to make many small independent matrix-matrix multiplications with dense matrices.

It is clear that even with millions of small independent matrices we will not be able to achieve the same *GFLOPS* rate as with a one large matrix. For example, a single $n \times n$ large matrix-matrix multiplication performs n^3 operations for n^2 input size, while $1024 \times \frac{n}{32} \times \frac{n}{32}$ small matrix-matrix multiplications perform $1024 \left(\frac{n}{32}\right)^3 = \frac{n^3}{32}$ operations for the same input size. However, it is also clear that we can achieve a significantly better performance with many small independent matrices compared with a single small matrix.

The *Fermi* architecture family of GPUs allows us to execute multiple kernels simultaneously. Hence, in order to batch the execution of independent kernels, we can run each of them in a separate stream. In particular, in the above example we could create 1024 CUDA™ streams using the function `cudaStreamCreate()`, then preface each call to `cublas<T>gemm()` with a call to `cublasSetStream()` with a different stream for each of the matrix-matrix multiplications. This will ensure that when possible the different computations will be executed concurrently. Although the user can create many streams, in practice it is not possible to have more than 16 concurrent kernels executing at the same time.

2.7 Device API

For devices with compute capability higher than 3.5, starting with CUDA Toolkit release 5.0, the CUBLAS library routines can be called from the device.

In order to use the CUBLAS library routines from the device the user must include the header file “cublas_v2.h” corresponding to the new CUBLAS API and link the static CUBLAS library with the device side routines to the application.

The CUBLAS library routines are called from the device in exactly the same way they are called from the host, with the following exceptions:

- ▶ The legacy CUBLAS API is not supported on the device.
- ▶ The pointer mode `CUBLAS_POINTER_MODE_HOST` is not supported on the device, in other words, scalar input and output parameters must be allocated in the device memory.
- ▶ Furthermore, the input and output scalar parameters must be allocated and released on the device using the `cudaMalloc()/malloc()` and `cudaFree()/free()` routines, respectively, in other words, they can not be passed by reference from the local memory to the routines.

Chapter 3

CUBLAS Datatypes Reference

3.1 cublasHandle_t

The `cublasHandle_t` type is a pointer type to an opaque structure holding the CUBLAS library context. The CUBLAS library context must be initialized using `cublasCreate()` and the returned handle must be passed to all subsequent library function calls. The context should be destroyed at the end using `cublasDestroy()`.

3.2 cublasStatus_t

The `cublasStatus_t` type is used for function status returns. All CUBLAS library functions return their status, which can have the following values.

CUBLAS_STATUS_SUCCESS	The operation completed successfully.
CUBLAS_STATUS_NOT_INITIALIZED	<p>The CUBLAS library was not initialized. This is usually caused by the lack of a prior <code>cublasCreate()</code> call, an error in the CUDA Runtime API called by the CUBLAS routine, or an error in the hardware setup.</p> <p>To correct: call <code>cublasCreate()</code> prior to the function call; and check that the hardware, an appropriate version of the driver, and the CUBLAS library are correctly installed.</p>
CUBLAS_STATUS_ALLOC_FAILED	<p>Resource allocation failed inside the CUBLAS library. This is usually caused by a <code>cudaMalloc()</code> failure.</p> <p>To correct: prior to the function call, deallocate previously allocated memory as much as possible.</p>
CUBLAS_STATUS_INVALID_VALUE	

An unsupported value or parameter was passed to the function (a negative vector size, for example).

To correct: ensure that all the parameters being passed have valid values.

CUBLAS_STATUS_ARCH_MISMATCH

The function requires a feature absent from the device architecture; usually caused by the lack of support for double precision.

To correct: compile and run the application on a device with appropriate compute capability, which is 1.3 for double precision.

CUBLAS_STATUS_MAPPING_ERROR

An access to GPU memory space failed, which is usually caused by a failure to bind a texture.

To correct: prior to the function call, unbind any previously bound textures.

CUBLAS_STATUS_EXECUTION_FAILED

The GPU program failed to execute. This is often caused by a launch failure of the kernel on the GPU, which can be caused by multiple reasons.

To correct: check that the hardware, an appropriate version of the driver, and the CUBLAS library are correctly installed.

CUBLAS_STATUS_INTERNAL_ERROR

An internal CUBLAS operation failed. This error is usually caused by a `cudaMemcpyAsync()` failure.

To correct: check that the hardware, an appropriate version of the driver, and the CUBLAS library are correctly installed. Also, check that the memory passed as a parameter to the routine is not being deallocated prior to the routine's completion.

3.3 cublasOperation_t

The `cublasOperation_t` type indicates which operation needs to be performed with the dense matrix. Its values correspond to Fortran characters 'N' or 'n' (non-transpose), 'T' or 't' (transpose) and 'C' or 'c' (conjugate transpose) that are often used as parameters to legacy BLAS implementations.

Value	Meaning
CUBLAS_OP_N	the non-transpose operation is selected
CUBLAS_OP_T	the transpose operation is selected
CUBLAS_OP_C	the conjugate transpose operation is selected

3.4 cublasFillMode_t

The `cublasFillMode_t` type indicates which part (lower or upper) of the dense matrix was filled and consequently should be used by the function. Its values correspond to Fortran characters ‘L’ or ‘l’ (lower) and ‘U’ or ‘u’ (upper) that are often used as parameters to legacy BLAS implementations.

Value	Meaning
<code>CUBLAS_FILL_MODE_LOWER</code>	the lower part of the matrix is filled
<code>CUBLAS_FILL_MODE_UPPER</code>	the upper part of the matrix is filled

3.5 cublasDiagType_t

The `cublasDiagType_t` type indicates whether the main diagonal of the dense matrix is unity and consequently should not be touched or modified by the function. Its values correspond to Fortran characters ‘N’ or ‘n’ (non-unit) and ‘U’ or ‘u’ (unit) that are often used as parameters to legacy BLAS implementations.

Value	Meaning
<code>CUBLAS_DIAG_NON_UNIT</code>	the matrix diagonal has non-unit elements
<code>CUBLAS_DIAG_UNIT</code>	the matrix diagonal has unit elements

3.6 cublasSideMode_t

The `cublasSideMode_t` type indicates whether the dense matrix is on the left or right side in the matrix equation solved by a particular function. Its values correspond to Fortran characters ‘L’ or ‘l’ (left) and ‘R’ or ‘r’ (right) that are often used as parameters to legacy BLAS implementations.

Value	Meaning
<code>CUBLAS_SIDE_LEFT</code>	the matrix is on the left side in the equation
<code>CUBLAS_SIDE_RIGHT</code>	the matrix is on the right side in the equation

3.7 cublasPointerMode_t

The `cublasPointerMode_t` type indicates whether the scalar values are passed by reference on the host or device. It is important to point out that if several scalar values are present in the function call, all of them must conform to the same single pointer mode. The pointer mode can be set and retrieved using `cublasSetPointerMode()` and `cublasGetPointerMode()` routines, respectively.

Value	Meaning
CUBLAS_POINTER_MODE_HOST	the scalars are passed by reference on the host
CUBLAS_POINTER_MODE_DEVICE	the scalars are passed by reference on the device

3.8 cublasAtomicMode_t

The `cublasAtomicMode_t` type indicates whether CUBLAS routines which has an alternate implementation using atomics can be used. The atomics mode can be set and queried using `cublasSetAtomicMode()` and `cublasGetAtomicMode()` routines, respectively.

Value	Meaning
CUBLAS_ATOMICS_NOT_ALLOWED	the usage of atomics is not allowed
CUBLAS_ATOMICS_ALLOWED	the usage of atomics is allowed

Chapter 4

CUBLAS Helper Function Reference

4.1 cublasCreate()

```
cublasStatus_t  
cublasCreate(cublasHandle_t *handle)
```

This function initializes the CUBLAS library and creates a handle to an opaque structure holding the CUBLAS library context. It allocates hardware resources on the host and device and must be called prior to making any other CUBLAS library calls.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the initialization succeeded
CUBLAS_STATUS_NOT_INITIALIZED	the CUDA™ Runtime initialization failed
CUBLAS_STATUS_ALLOC_FAILED	the resources could not be allocated

4.2 cublasDestroy()

```
cublasStatus_t  
cublasDestroy(cublasHandle_t handle)
```

This function releases hardware resources used by the CUBLAS library. The release of GPU resources may be deferred until the application exits. This function is usually the last call with a particular handle to the CUBLAS library.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the shut down succeeded
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized

4.3 cublasGetVersion()

```
cublasStatus_t  
cublasGetVersion(cublasHandle_t handle, int *version)
```

This function returns the version number of the CUBLAS library.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized

4.4 cublasSetStream()

```
cublasStatus_t
cublasSetStream(cublasHandle_t handle, cudaStream_t streamId)
```

This function sets the CUBLAS library stream, which will be used to execute all subsequent calls to the CUBLAS library functions. If the CUBLAS library stream is not set, all kernels use the *default* NULL stream. In particular, this routine can be used to change the stream between kernel launches and then to reset the CUBLAS library stream back to NULL.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the stream was set successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized

4.5 cublasGetStream()

```
cublasStatus_t
cublasGetStream(cublasHandle_t handle, cudaStream_t *streamId)
```

This function gets the CUBLAS library stream, which is being used to execute all calls to the CUBLAS library functions. If the CUBLAS library stream is not set, all kernels use the *default* NULL stream.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the stream was returned successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized

4.6 cublasGetPointerMode()

```
cublasStatus_t
cublasGetPointerMode(cublasHandle_t handle, cublasPointerMode_t *mode)
```

This function obtains the pointer mode used by the CUBLAS library. Please see the section on the `cublasPointerMode_t` type for more details.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the pointer mode was obtained successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized

4.7 cublasSetPointerMode()

```
cublasStatus_t
cublasSetPointerMode(cublasHandle_t handle, cublasPointerMode_t mode)
```

This function sets the pointer mode used by the CUBLAS library. The *default* is for the values to be passed by reference on the host. Please see the section on the `cublasPointerMode_t` type for more details.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the pointer mode was set successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized

4.8 cublasSetVector()

```
cublasStatus_t
cublasSetVector(int n, int elemSize,
                const void *x, int incx, void *devicePtr, int incy)
```

This function copies `n` elements from a vector `x` in host memory space to a vector `y` in GPU memory space. Elements in both vectors are assumed to have a size of `elemSize` bytes. The storage spacing between consecutive elements is given by `incx` for the source vector `x` and `incy` for the destination vector `y`.

In general, `y` points to an object, or part of an object, that was allocated via `cublasAlloc()`. Since column-major format for two-dimensional matrices is assumed, if a vector is part of a matrix, a vector increment equal to 1 accesses a (partial) column of that matrix. Similarly, using an increment equal to the leading dimension of the matrix results in accesses to a (partial) row of that matrix.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>incx, incy, elemSize <= 0</code>
CUBLAS_STATUS_MAPPING_ERROR	there was an error accessing GPU memory

4.9 cublasGetVector()

```
cublasStatus_t
cublasGetVector(int n, int elemSize,
                const void *x, int incx, void *y, int incy)
```

This function copies `n` elements from a vector `x` in GPU memory space to a vector `y` in host memory space. Elements in both vectors are assumed to have a size of `elemSize`

bytes. The storage spacing between consecutive elements is given by `incx` for the source vector `x` and `incy` for the destination vector `y`.

In general, `x` points to an object, or part of an object, that was allocated via `cublasAlloc()`. Since column-major format for two-dimensional matrices is assumed, if a vector is part of a matrix, a vector increment equal to 1 accesses a (partial) column of that matrix. Similarly, using an increment equal to the leading dimension of the matrix results in accesses to a (partial) row of that matrix.

Return Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	the operation completed successfully
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	the library was not initialized
<code>CUBLAS_STATUS_INVALID_VALUE</code>	the parameters <code>incx, incy, elemSize <= 0</code>
<code>CUBLAS_STATUS_MAPPING_ERROR</code>	there was an error accessing GPU memory

4.10 cublasSetMatrix()

```
cublasStatus_t
cublasSetMatrix(int rows, int cols, int elemSize,
                const void *A, int lda, void *B, int ldb)
```

This function copies a tile of `rows × cols` elements from a matrix `A` in host memory space to a matrix `B` in GPU memory space. It is assumed that each element requires storage of `elemSize` bytes and that both matrices are stored in column-major format, with the leading dimension of the source matrix `A` and destination matrix `B` given in `lda` and `ldb`, respectively. The leading dimension indicates the number of rows of the allocated matrix, even if only a submatrix of it is being used. In general, `B` is a device pointer that points to an object, or part of an object, that was allocated in GPU memory space via `cublasAlloc()`.

Return Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	the operation completed successfully
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	the library was not initialized
<code>CUBLAS_STATUS_INVALID_VALUE</code>	the parameters <code>rows, cols < 0</code> or <code>elemSize, lda, ldb <= 0</code>
<code>CUBLAS_STATUS_MAPPING_ERROR</code>	there was an error accessing GPU memory

4.11 cublasGetMatrix()

```
cublasStatus_t
cublasGetMatrix(int rows, int cols, int elemSize,
                const void *A, int lda, void *B, int ldb)
```

This function copies a tile of `rows × cols` elements from a matrix `A` in GPU memory space to a matrix `B` in host memory space. It is assumed that each element requires storage of

`elemSize` bytes and that both matrices are stored in column-major format, with the leading dimension of the source matrix `A` and destination matrix `B` given in `lda` and `ldb`, respectively. The leading dimension indicates the number of rows of the allocated matrix, even if only a submatrix of it is being used. In general, `A` is a device pointer that points to an object, or part of an object, that was allocated in GPU memory space via `cublasAlloc()`.

Return Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	the operation completed successfully
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	the library was not initialized
<code>CUBLAS_STATUS_INVALID_VALUE</code>	the parameters <code>rows,cols<0</code> or <code>elemSize,lda,ldb<=0</code>
<code>CUBLAS_STATUS_MAPPING_ERROR</code>	there was an error accessing GPU memory

4.12 cublasSetVectorAsync()

```
cublasStatus_t
cublasSetVectorAsync(int n, int elemSize, const void *hostPtr, int incx,
                    void *devicePtr, int incy, cudaStream_t stream)
```

This function has the same functionality as `cublasSetVector()`, with the exception that the data transfer is done asynchronously (with respect to the host) using the given CUDA™ stream parameter.

Return Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	the operation completed successfully
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	the library was not initialized
<code>CUBLAS_STATUS_INVALID_VALUE</code>	the parameters <code>incx,incy,elemSize<=0</code>
<code>CUBLAS_STATUS_MAPPING_ERROR</code>	there was an error accessing GPU memory

4.13 cublasGetVectorAsync()

```
cublasStatus_t
cublasGetVectorAsync(int n, int elemSize, const void *devicePtr, int incx,
                    void *hostPtr, int incy, cudaStream_t stream)
```

This function has the same functionality as `cublasGetVector()`, with the exception that the data transfer is done asynchronously (with respect to the host) using the given CUDA™ stream parameter.

Return Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	the operation completed successfully
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	the library was not initialized
<code>CUBLAS_STATUS_INVALID_VALUE</code>	the parameters <code>incx,incy,elemSize<=0</code>
<code>CUBLAS_STATUS_MAPPING_ERROR</code>	there was an error accessing GPU memory

4.14 cublasSetMatrixAsync()

```
cublasStatus_t
cublasSetMatrixAsync(int rows, int cols, int elemSize, const void *A,
                    int lda, void *B, int ldb, cudaStream_t stream)
```

This function has the same functionality as `cublasSetMatrix()`, with the exception that the data transfer is done asynchronously (with respect to the host) using the given CUDA™ stream parameter.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>rows,cols<0</code> or <code>elemSize,lda,ldb<=0</code>
CUBLAS_STATUS_MAPPING_ERROR	there was an error accessing GPU memory

4.15 cublasGetMatrixAsync()

```
cublasStatus_t
cublasGetMatrixAsync(int rows, int cols, int elemSize, const void *A,
                    int lda, void *B, int ldb, cudaStream_t stream)
```

This function has the same functionality as `cublasGetMatrix()`, with the exception that the data transfer is done asynchronously (with respect to the host) using the given CUDA™ stream parameter.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>rows,cols<0</code> or <code>elemSize,lda,ldb<=0</code>
CUBLAS_STATUS_MAPPING_ERROR	there was an error accessing GPU memory

4.16 cublasGetMatrixAsync()

```
cublasStatus_t
cublasGetMatrixAsync(int rows, int cols, int elemSize, const void *A,
                    int lda, void *B, int ldb, cudaStream_t stream)
```

This function has the same functionality as `cublasGetMatrix()`, with the exception that the data transfer is done asynchronously (with respect to the host) using the given CUDA™ stream parameter.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>rows,cols<0</code> or <code>elemSize,lda,ldb<=0</code>
CUBLAS_STATUS_MAPPING_ERROR	there was an error accessing GPU memory

4.17 cublasSetAtomicMode()

```
cublasStatus_t
cublasSetAtomicMode(cublasHandle_t handle, cublasAtomicMode_t mode)
```

Some routines like `cublas<t>symv` and `cublas<t>hemv` have an alternate implementation that use atomics to cumulate results. This implementation is generally significantly faster but can generate results that are not strictly identical from one run to the others. Mathematically, those different results are not significant but when debugging those differences can be prejudicial.

This function allows or disallows the usage of atomics in the CUBLAS library for all routines which have an alternate implementation. When not explicitly specified in the documentation of any CUBLAS routine, it means that this routine does not have an alternate implementation that use atomics. When atomics mode is disabled, each CUBLAS routine should produce the same results from one run to the other when called with identical parameters on the same Hardware.

The *default* value of the atomics mode is `CUBLAS_ATOMICS_NOT_ALLOWED`. Please see the section on the `cublasAtomicMode_t` type for more details.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the atomics mode was set successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized

4.18 cublasGetAtomicMode()

```
cublasStatus_t
cublasGetAtomicMode(cublasHandle_t handle, cublasAtomicMode_t *mode)
```

This function queries the atomic mode of a specific CUBLAS context.

The *default* value of the atomics mode is `CUBLAS_ATOMICS_NOT_ALLOWED`. Please see the section on the `cublasAtomicMode_t` type for more details.

Return Value	Meaning
CUBLAS_STATUS_SUCCESS	the atomics mode was queried successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized

Chapter 5

CUBLAS Level-1 Function Reference

In this chapter we describe the Level-1 Basic Linear Algebra Subprograms (BLAS1) functions that perform scalar and vector based operations. We will use abbreviations $\langle type \rangle$ for type and $\langle t \rangle$ for the corresponding short type to make a more concise and clear presentation of the implemented functions. Unless otherwise specified $\langle type \rangle$ and $\langle t \rangle$ have the following meanings:

$\langle type \rangle$	$\langle t \rangle$	Meaning
float	's' or 'S'	real single-precision
double	'd' or 'D'	real double-precision
cuComplex	'c' or 'C'	complex single-precision
cuDoubleComplex	'z' or 'Z'	complex double-precision

When the parameters and returned values of the function differ, which sometimes happens for complex input, the $\langle t \rangle$ can also have the following meanings 'Sc', 'Cs', 'Dz' and 'Zd'.

The abbreviation **Re**(.) and **Im**(.) will stand for the real and imaginary part of a number, respectively. Since imaginary part of a real number does not exist, we will consider it to be zero and can usually simply discard it from the equation where it is being used. Also, the $\bar{\alpha}$ will denote the complex conjugate of α .

In general throughout the documentation, the lower case Greek symbols α and β will denote scalars, lower case English letters in bold type **x** and **y** will denote vectors and capital English letters *A*, *B* and *C* will denote matrices.

5.1 cublas $\langle t \rangle$ amax()

```
cublasStatus_t cublasIsamax(cublasHandle_t handle, int n,
                           const float      *x, int incx, int *result)
cublasStatus_t cublasIdamax(cublasHandle_t handle, int n,
                           const double     *x, int incx, int *result)
cublasStatus_t cublasIcamax(cublasHandle_t handle, int n,
                           const cuComplex  *x, int incx, int *result)
cublasStatus_t cublasIzamax(cublasHandle_t handle, int n,
                           const cuDoubleComplex *x, int incx, int *result)
```

This function finds the (smallest) index of the element of the maximum magnitude. Hence, the result is the first i such that $|\mathbf{Im}(\mathbf{x}[j])| + |\mathbf{Re}(\mathbf{x}[j])|$ is maximum for $i = 1, \dots, n$ and $j = 1 + (i - 1) * \text{incx}$. Notice that the last equation reflects 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle	device	input	handle to the CUBLAS library context.
n		input	number of elements in the vector \mathbf{x} .
x		input	<type> vector with n elements.
incx		input	stride between consecutive elements of \mathbf{x} .
result	host or device	output	the resulting index, which is 0 if $n, \text{incx} \leq 0$.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ALLOC_FAILED	the reduction buffer could not be allocated
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[isamax](#), [idamax](#), [icamax](#), [izamax](#)

5.2 cublasI<t>amin()

```

cublasStatus_t cublasIsamin(cublasHandle_t handle, int n,
                           const float      *x, int incx, int *result)
cublasStatus_t cublasIdamin(cublasHandle_t handle, int n,
                           const double     *x, int incx, int *result)
cublasStatus_t cublasIcamin(cublasHandle_t handle, int n,
                           const cuComplex  *x, int incx, int *result)
cublasStatus_t cublasIzamin(cublasHandle_t handle, int n,
                           const cuDoubleComplex *x, int incx, int *result)

```

This function finds the (smallest) index of the element of the minimum magnitude. Hence, the result is the first i such that $|\mathbf{Im}(\mathbf{x}[j])| + |\mathbf{Re}(\mathbf{x}[j])|$ is minimum for $i = 1, \dots, n$ and $j = 1 + (i - 1) * \text{incx}$. Notice that the last equation reflects 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle	device	input	handle to the CUBLAS library context.
n		input	number of elements in the vector \mathbf{x} .
x		input	<type> vector with n elements.
incx		input	stride between consecutive elements of \mathbf{x} .
result	host or device	output	the resulting index, which is 0 if $n, \text{incx} \leq 0$.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ALLOC_FAILED	the reduction buffer could not be allocated
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[isamin](#), [idamin](#)

5.3 cublas<t>asum()

```
cublasStatus_t cublasSasum(cublasHandle_t handle, int n,
                           const float          *x, int incx, float  *result)
```

```
cublasStatus_t cublasDasum(cublasHandle_t handle, int n,
                           const double         *x, int incx, double *result)
```

```
cublasStatus_t cublasScasum(cublasHandle_t handle, int n,
                           const cuComplex      *x, int incx, float  *result)
```

```
cublasStatus_t cublasDzasum(cublasHandle_t handle, int n,
                           const cuDoubleComplex *x, int incx, double *result)
```

This function computes the sum of the absolute values of the elements of vector **x**. Hence, the result is $\sum_{i=1}^n (|\mathbf{Im}(\mathbf{x}[j])| + |\mathbf{Re}(\mathbf{x}[j])|)$ where $j = 1 + (i - 1) * \text{incx}$. Notice that the last equation reflects 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle	device	input	handle to the CUBLAS library context.
n		input	number of elements in the vector x .
x		input	<type> vector with n elements.
incx		input	stride between consecutive elements of x .
result	host or device	output	the resulting sum, which is 0.0 if n, incx ≤ 0.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ALLOC_FAILED	the reduction buffer could not be allocated
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[sasum](#), [dasum](#), [scasum](#), [dzasum](#)

5.4 cublas<t>axpy()

```
cublasStatus_t cublasSaxpy(cublasHandle_t handle, int n,
                           const float          *alpha,
                           const float          *x, int incx,
                           float                *y, int incy)
cublasStatus_t cublasDaxpy(cublasHandle_t handle, int n,
                           const double         *alpha,
                           const double         *x, int incx,
                           double               *y, int incy)
cublasStatus_t cublasCaxpy(cublasHandle_t handle, int n,
                           const cuComplex      *alpha,
                           const cuComplex      *x, int incx,
                           cuComplex            *y, int incy)
cublasStatus_t cublasZaxpy(cublasHandle_t handle, int n,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *x, int incx,
                           cuDoubleComplex      *y, int incy)
```

This function multiplies the vector \mathbf{x} by the scalar α and adds it to the vector \mathbf{y} overwriting the latest vector with the result. Hence, the performed operation is $\mathbf{y}[j] = \alpha \times \mathbf{x}[k] + \mathbf{y}[j]$ for $i = 1, \dots, n$, $k = 1 + (i - 1) * \text{incx}$ and $j = 1 + (i - 1) * \text{incy}$. Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle	host or device	input	handle to the CUBLAS library context.
alpha		input	<type> scalar used for multiplication.
n		input	number of elements in the vector \mathbf{x} and \mathbf{y} .
x	device	input	<type> vector with n elements.
incx	device	input	stride between consecutive elements of \mathbf{x} .
y		in/out	<type> vector with n elements.
incy		input	stride between consecutive elements of \mathbf{y} .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[saxpy](#), [daxpy](#), [caxpy](#), [zaxpy](#)

5.5 cublas<t>copy()

```

cublasStatus_t cublasScopy(cublasHandle_t handle, int n,
                           const float          *x, int incx,
                           float                *y, int incy)
cublasStatus_t cublasDcopy(cublasHandle_t handle, int n,
                           const double         *x, int incx,
                           double               *y, int incy)
cublasStatus_t cublasCcopy(cublasHandle_t handle, int n,
                           const cuComplex      *x, int incx,
                           cuComplex            *y, int incy)
cublasStatus_t cublasZcopy(cublasHandle_t handle, int n,
                           const cuDoubleComplex *x, int incx,
                           cuDoubleComplex      *y, int incy)

```

This function copies the vector \mathbf{x} into the vector \mathbf{y} . Hence, the performed operation is $\mathbf{y}[j] = \mathbf{x}[k]$ for $i = 1, \dots, n$, $k = 1 + (i - 1) * \text{incx}$ and $j = 1 + (i - 1) * \text{incy}$. Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
n		input	number of elements in the vector \mathbf{x} and \mathbf{y} .
x	device	input	<type> vector with n elements.
incx		input	stride between consecutive elements of \mathbf{x} .
y	device	output	<type> vector with n elements.
incy		input	stride between consecutive elements of \mathbf{y} .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[scopy](#), [dcopy](#), [ccopy](#), [zcopy](#)

5.6 cublas<t>dot()

```

cublasStatus_t cublasSdot (cublasHandle_t handle, int n,
                           const float          *x, int incx,
                           const float          *y, int incy,
                           float                *result)
cublasStatus_t cublasDdot (cublasHandle_t handle, int n,

```

```

        const double      *x, int incx,
        const double      *y, int incy,
        double             *result)
cublasStatus_t cublasCdotu(cublasHandle_t handle, int n,
        const cuComplex   *x, int incx,
        const cuComplex   *y, int incy,
        cuComplex          *result)
cublasStatus_t cublasCdotc(cublasHandle_t handle, int n,
        const cuComplex   *x, int incx,
        const cuComplex   *y, int incy,
        cuComplex          *result)
cublasStatus_t cublasZdotu(cublasHandle_t handle, int n,
        const cuDoubleComplex *x, int incx,
        const cuDoubleComplex *y, int incy,
        cuDoubleComplex *result)
cublasStatus_t cublasZdotc(cublasHandle_t handle, int n,
        const cuDoubleComplex *x, int incx,
        const cuDoubleComplex *y, int incy,
        cuDoubleComplex *result)

```

This function computes the dot product of vectors \mathbf{x} and \mathbf{y} . Hence, the result is $\sum_{i=1}^n (\mathbf{x}[k] \times \mathbf{y}[j])$ where $k = 1 + (i - 1) * \text{incx}$ and $j = 1 + (i - 1) * \text{incy}$. Notice that in the first equation the conjugate of the element of vector \mathbf{x} should be used if the function name ends in character 'c' and that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
n		input	number of elements in the vectors \mathbf{x} and \mathbf{y} .
x	device	input	<type> vector with \mathbf{n} elements.
incx		input	stride between consecutive elements of \mathbf{x} .
y	device	input	<type> vector with \mathbf{n} elements.
incy		input	stride between consecutive elements of \mathbf{y} .
result	host or device	output	the resulting dot product, which is 0.0 if $\mathbf{n} \leq 0$.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ALLOC_FAILED	the reduction buffer could not be allocated
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[sdot](#), [ddot](#), [cdotu](#), [cdotc](#), [zdotu](#), [zdotc](#)

5.7 cublas<t>nrm2()

```

cublasStatus_t cublasSnrm2(cublasHandle_t handle, int n,
                           const float          *x, int incx, float  *result)
cublasStatus_t cublasDnrm2(cublasHandle_t handle, int n,
                           const double         *x, int incx, double *result)
cublasStatus_t cublasScnrm2(cublasHandle_t handle, int n,
                            const cuComplex      *x, int incx, float  *result)
cublasStatus_t cublasDznrm2(cublasHandle_t handle, int n,
                            const cuDoubleComplex *x, int incx, double *result)

```

This function computes the Euclidean norm of the vector **x**. The code uses a multiphase model of accumulation to avoid intermediate underflow and overflow, with the result being equivalent to $\sqrt{\sum_{i=1}^n (\bar{x}[j] \times x[j])}$ where $j = 1 + (i - 1) * \text{incx}$ in exact arithmetic. Notice that the last equation reflects 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle	device	input	handle to the CUBLAS library context.
n		input	number of elements in the vector x .
x		input	<type> vector with n elements.
incx		input	stride between consecutive elements of x .
result	host or device	output	the resulting norm, which is 0.0 if n, incx ≤ 0.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ALLOC_FAILED	the reduction buffer could not be allocated
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[snrm2](#), [snrm2](#), [dnrm2](#), [dnrm2](#), [scnrm2](#), [scnrm2](#), [dznrm2](#)

5.8 cublas<t>rot()

```

cublasStatus_t cublasSrot (cublasHandle_t handle, int n,
                           float *x, int incx,
                           float *y, int incy,
                           const float  *c, const float  *s)
cublasStatus_t cublasDrot (cublasHandle_t handle, int n,
                           double      *x, int incx,
                           double      *y, int incy,

```

```

                                const double *c, const double      *s)

cublasStatus_t cublasCrot (cublasHandle_t handle, int n,
                           cuComplex      *x, int incx,
                           cuComplex      *y, int incy,
                           const float    *c, const cuComplex      *s)
cublasStatus_t cublasCsrot(cublasHandle_t handle, int n,
                           cuComplex      *x, int incx,
                           cuComplex      *y, int incy,
                           const float    *c, const float          *s)
cublasStatus_t cublasZrot (cublasHandle_t handle, int n,
                           cuDoubleComplex *x, int incx,
                           cuDoubleComplex *y, int incy,
                           const double    *c, const cuDoubleComplex *s)
cublasStatus_t cublasZdrot(cublasHandle_t handle, int n,
                           cuDoubleComplex *x, int incx,
                           cuDoubleComplex *y, int incy,
                           const double    *c, const double          *s)

```

This function applies Givens rotation matrix

$$G = \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$$

to vectors \mathbf{x} and \mathbf{y} .

Hence, the result is $\mathbf{x}[k] = c \times \mathbf{x}[k] + s \times \mathbf{y}[j]$ and $\mathbf{y}[j] = -s \times \mathbf{x}[k] + c \times \mathbf{y}[j]$ where $k = 1 + (i - 1) * \text{incx}$ and $j = 1 + (i - 1) * \text{incy}$. Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
n		input	number of elements in the vectors \mathbf{x} and \mathbf{y} .
x	device	in/out	<type> vector with n elements.
incx		input	stride between consecutive elements of \mathbf{x} .
y	device	in/out	<type> vector with n elements.
incy		input	stride between consecutive elements of \mathbf{y} .
c	host or device	input	cosine element of the rotation matrix.
s	host or device	input	sine element of the rotation matrix.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[srot](#), [drot](#), [crot](#), [csrot](#), [zrot](#), [zdrot](#)

5.9 cublas<t>rotg()

```
cublasStatus_t cublasSrotg(cublasHandle_t handle,
                           float          *a, float          *b,
                           float *c, float          *s)
cublasStatus_t cublasDrotg(cublasHandle_t handle,
                           double         *a, double         *b,
                           double *c, double         *s)
cublasStatus_t cublasCrotg(cublasHandle_t handle,
                           cuComplex      *a, cuComplex      *b,
                           float *c,      cuComplex *s)
cublasStatus_t cublasZrotg(cublasHandle_t handle,
                           cuDoubleComplex *a, cuDoubleComplex *b,
                           double *c, cuDoubleComplex *s)
```

This function constructs the Givens rotation matrix

$$G = \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$$

that zeros out the second entry of a 2×1 vector $(a, b)^T$.

Then, for real numbers we can write

$$\begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} r \\ 0 \end{pmatrix}$$

where $c^2 + s^2 = 1$ and $r = \sqrt{a^2 + b^2}$. The parameters a and b are overwritten with r and z , respectively. The value of z is such that c and s may be recovered using the following rules:

$$(c, s) = \begin{cases} (\sqrt{1 - z^2}, z) & \text{if } |z| < 1 \\ (0.0, 1.0) & \text{if } |z| = 1 \\ (1/z, \sqrt{1 - z^2}) & \text{if } |z| > 1 \end{cases}$$

For complex numbers we can write

$$\begin{pmatrix} c & s \\ -\bar{s} & c \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} r \\ 0 \end{pmatrix}$$

where $c^2 + (\bar{s} \times s) = 1$ and $r = \frac{a}{|a|} \times \|(a, b)^T\|_2$ with $\|(a, b)^T\|_2 = \sqrt{|a|^2 + |b|^2}$ for $a \neq 0$ and $r = b$ for $a = 0$. Finally, the parameter a is overwritten with r on exit.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
a	host or device	in/out	<type> scalar that is overwritten with r .
b	host or device	in/out	<type> scalar that is overwritten with z .
c	host or device	output	cosine element of the rotation matrix.
s	host or device	output	sine element of the rotation matrix.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[srotg](#), [drotg](#), [crotg](#), [zrotg](#)

5.10 cublas<t>rotm()

```
cublasStatus_t cublasSrotm(cublasHandle_t handle, int n, float *x, int incx,
                           float *y, int incy, const float* param)
cublasStatus_t cublasDrotm(cublasHandle_t handle, int n, double *x, int incx,
                           double *y, int incy, const double* param)
```

This function applies the modified Givens transformation

$$H = \begin{pmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{pmatrix}$$

to vectors \mathbf{x} and \mathbf{y} .

Hence, the result is $\mathbf{x}[k] = h_{11} \times \mathbf{x}[k] + h_{12} \times \mathbf{y}[j]$ and $\mathbf{y}[j] = h_{21} \times \mathbf{x}[k] + h_{22} \times \mathbf{y}[j]$ where $k = 1 + (i - 1) * \text{incx}$ and $j = 1 + (i - 1) * \text{incy}$. Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

The elements h_{11} , h_{21} , h_{12} and h_{22} of 2×2 matrix H are stored in `param[1]`, `param[2]`, `param[3]` and `param[4]`, respectively. The `flag = param[0]` defines the following predefined values for the matrix H entries

flag=-1.0	flag= 0.0	flag= 1.0	flag=-2.0
$\begin{pmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{pmatrix}$	$\begin{pmatrix} 1.0 & h_{12} \\ h_{21} & 1.0 \end{pmatrix}$	$\begin{pmatrix} h_{11} & 1.0 \\ -1.0 & h_{22} \end{pmatrix}$	$\begin{pmatrix} 1.0 & 0.0 \\ 0.0 & 1.0 \end{pmatrix}$

Notice that the values -1.0, 0.0 and 1.0 implied by the flag are not stored in `param`.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
n		input	number of elements in the vectors x and y .
x	device	in/out	<type> vector with n elements.
incx		input	stride between consecutive elements of x .
y	device	in/out	<type> vector with n elements.
incy		input	stride between consecutive elements of y .
param	host or device	input	<type> vector of 5 elements, where param [0] and param [1-4] contain the flag and matrix H .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[srotm](#), [drotm](#)

5.11 cublas<t>rotmg()

```
cublasStatus_t cublasSrotmg(cublasHandle_t handle, float *d1, float *d2,
                           float *x1, const float *y1, float *param)
cublasStatus_t cublasDrotmg(cublasHandle_t handle, double *d1, double *d2,
                           double *x1, const double *y1, double *param)
```

This function constructs the modified Givens transformation

$$H = \begin{pmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{pmatrix}$$

that zeros out the second entry of a 2×1 vector $(\sqrt{d1} * x1, \sqrt{d2} * y1)^T$.

The **flag** = **param**[0] defines the following predefined values for the matrix **H** entries

flag=-1.0	flag= 0.0	flag= 1.0	flag=-2.0
$\begin{pmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{pmatrix}$	$\begin{pmatrix} 1.0 & h_{12} \\ h_{21} & 1.0 \end{pmatrix}$	$\begin{pmatrix} h_{11} & 1.0 \\ -1.0 & h_{22} \end{pmatrix}$	$\begin{pmatrix} 1.0 & 0.0 \\ 0.0 & 1.0 \end{pmatrix}$

Notice that the values -1.0, 0.0 and 1.0 implied by the flag are not stored in **param**.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
d1	host or device	in/out	<type> scalar that is overwritten on exit.
d2	host or device	in/out	<type> scalar that is overwritten on exit.
x1	host or device	in/out	<type> scalar that is overwritten on exit.
y1	host or device	input	<type> scalar.
param	host or device	output	<type> vector of 5 elements, where <code>param[0]</code> and <code>param[1-4]</code> contain the flag and matrix H .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[srotmg](#), [drotmg](#)

5.12 cublas<t>scal()

```

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)

```

This function scales the vector \mathbf{x} by the scalar α and overwrites it with the result. Hence, the performed operation is $\mathbf{x}[j] = \alpha \times \mathbf{x}[j]$ for $i = 1, \dots, n$ and $j = 1 + (i - 1) * \text{incx}$. Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle	host or device	input	handle to the CUBLAS library context.
alpha		input	<type> scalar used for multiplication.
n		input	number of elements in the vector x .
x	device	in/out	<type> vector with n elements.
incx		input	stride between consecutive elements of x .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[sscal](#), [dscal](#), [csscal](#), [cscal](#), [zdscal](#), [zscal](#)

5.13 cublas<t>swap()

```

cublasStatus_t cublasSswap(cublasHandle_t handle, int n, float      *x,
                           int incx, float      *y, int incy)
cublasStatus_t cublasDswap(cublasHandle_t handle, int n, double    *x,
                           int incx, double    *y, int incy)
cublasStatus_t cublasCswap(cublasHandle_t handle, int n, cuComplex *x,
                           int incx, cuComplex *y, int incy)
cublasStatus_t cublasZswap(cublasHandle_t handle, int n, cuDoubleComplex *x,
                           int incx, cuDoubleComplex *y, int incy)

```

This function interchanges the elements of vector **x** and **y**. Hence, the performed operation is $\mathbf{y}[j] \Leftrightarrow \mathbf{x}[k]$ for $i = 1, \dots, n$, $k = 1 + (i - 1) * \text{incx}$ and $j = 1 + (i - 1) * \text{incy}$. Notice that the last two equations reflect 1-based indexing used for compatibility with Fortran.

Param.	Memory	In/out	Meaning
handle	device	input	handle to the CUBLAS library context.
n		input	number of elements in the vector x and y .
x		in/out	<type> vector with n elements.
incx		input	stride between consecutive elements of x .
y	device	in/out	<type> vector with n elements.
incy		input	stride between consecutive elements of y .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[sswap](#), [dswap](#), [cswap](#), [zswap](#)

Chapter 6

CUBLAS Level-2 Function Reference

In this chapter we describe the Level-2 Basic Linear Algebra Subprograms (BLAS2) functions that perform matrix-vector operations.

6.1 cublas<t>gbmv()

```
cublasStatus_t cublasSgbmv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n, int kl, int ku,
                           const float      *alpha,
                           const float      *A, int lda,
                           const float      *x, int incx,
                           const float      *beta,
                           float            *y, int incy)
cublasStatus_t cublasDgbmv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n, int kl, int ku,
                           const double     *alpha,
                           const double     *A, int lda,
                           const double     *x, int incx,
                           const double     *beta,
                           double           *y, int incy)
cublasStatus_t cublasCgbmv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n, int kl, int ku,
                           const cuComplex  *alpha,
                           const cuComplex  *A, int lda,
                           const cuComplex  *x, int incx,
                           const cuComplex  *beta,
                           cuComplex        *y, int incy)
cublasStatus_t cublasZgbmv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n, int kl, int ku,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *x, int incx,
```

```
const cuDoubleComplex *beta,
cuDoubleComplex *y, int incy)
```

This function performs the banded matrix-vector multiplication

$$\mathbf{y} = \alpha \text{op}(A)\mathbf{x} + \beta \mathbf{y}$$

where A is a $m \times n$ banded matrix with kl subdiagonals and ku superdiagonals, \mathbf{x} and \mathbf{y} are vectors, and α and β are scalars. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS_OP_N} \\ A^T & \text{if transa} == \text{CUBLAS_OP_T} \\ A^H & \text{if transa} == \text{CUBLAS_OP_C} \end{cases}$$

The banded matrix A is stored column by column, with the main diagonal stored in row $ku + 1$ (starting in first position), the first superdiagonal stored in row ku (starting in second position), the first subdiagonal stored in row $ku + 2$ (starting in first position), etc. So that in general, the element $A(i, j)$ is stored in the memory location $A(ku+1+i-j, j)$ for $j = 1, \dots, n$ and $i \in [\max(1, j - ku), \min(m, j + kl)]$. Also, the elements in the array A that do not conceptually correspond to the elements in the banded matrix (the top left $ku \times ku$ and bottom right $kl \times kl$ triangles) are not referenced.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
trans		input	operation $\text{op}(A)$ that is non- or (conj.) transpose.
m		input	number of rows of matrix A .
n		input	number of columns of matrix A .
kl		input	number of subdiagonals of matrix A .
ku		input	number of superdiagonals of matrix A .
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $\text{lda} \times n$ with $\text{lda} \geq \text{kl} + \text{ku} + 1$.
lda		input	leading dimension of two-dimensional array used to store matrix A .
x	device	input	<type> vector with n elements if $\text{transa} == \text{CUBLAS_OP_N}$ and m elements otherwise.
incx		input	stride between consecutive elements of \mathbf{x} .
beta	host or device	input	<type> scalar used for multiplication, if $\text{beta} == 0$ then \mathbf{y} does not have to be a valid input.
y	device	in/out	<type> vector with m elements if $\text{transa} == \text{CUBLAS_OP_N}$ and n elements otherwise.
incy		input	stride between consecutive elements of \mathbf{y} .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters m,n,kl,ku<0 or incx,incy=0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[sgbmv](#), [dgbmv](#), [cgbmv](#), [zgbmv](#)

6.2 cublas<t>gemv()

```

cublasStatus_t cublasSgemv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n,
                           const float      *alpha,
                           const float      *A, int lda,
                           const float      *x, int incx,
                           const float      *beta,
                           float            *y, int incy)
cublasStatus_t cublasDgemv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n,
                           const double     *alpha,
                           const double     *A, int lda,
                           const double     *x, int incx,
                           const double     *beta,
                           double           *y, int incy)
cublasStatus_t cublasCgemv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n,
                           const cuComplex  *alpha,
                           const cuComplex  *A, int lda,
                           const cuComplex  *x, int incx,
                           const cuComplex  *beta,
                           cuComplex        *y, int incy)
cublasStatus_t cublasZgemv(cublasHandle_t handle, cublasOperation_t trans,
                           int m, int n,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *x, int incx,
                           const cuDoubleComplex *beta,
                           cuDoubleComplex *y, int incy)

```

This function performs the matrix-vector multiplication

$$\mathbf{y} = \alpha \mathbf{op}(A)\mathbf{x} + \beta \mathbf{y}$$

where A is a $m \times n$ matrix stored in column-major format, \mathbf{x} and \mathbf{y} are vectors, and α and β are scalars. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS_OP_N} \\ A^T & \text{if transa} == \text{CUBLAS_OP_T} \\ A^H & \text{if transa} == \text{CUBLAS_OP_C} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
trans		input	operation op(A) that is non- or (conj.) transpose.
m		input	number of rows of matrix A.
n		input	number of columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $\text{lda} \times \text{n}$ with $\text{lda} \geq \max(1, \text{m})$ if $\text{transa} == \text{CUBLAS_OP_N}$ and $\text{lda} \times \text{m}$ with $\text{lda} \geq \max(1, \text{n})$ otherwise.
lda		input	leading dimension of two-dimensional array used to store matrix A.
x	device	input	<type> vector with n elements if $\text{transa} == \text{CUBLAS_OP_N}$ and m elements otherwise.
incx		input	stride between consecutive elements of \mathbf{x} .
beta	host or device	input	<type> scalar used for multiplication, if $\text{beta} == 0$ then \mathbf{y} does not have to be a valid input.
y	device	input	<type> vector with m elements if $\text{transa} == \text{CUBLAS_OP_N}$ and n elements otherwise.
incy		in/out	stride between consecutive elements of \mathbf{y} .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters $\text{m}, \text{n} < 0$ or $\text{incx}, \text{incy} = 0$
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[sgemv](#), [dgemv](#), [cgemv](#), [zgemv](#)

6.3 cublas<t>ger()

```
cublasStatus_t cublasSger (cublasHandle_t handle, int m, int n,
                          const float      *alpha,
```

```

        const float          *x, int incx,
        const float          *y, int incy,
        float                *A, int lda)
cublasStatus_t cublasDger (cublasHandle_t handle, int m, int n,
        const double         *alpha,
        const double         *x, int incx,
        const double         *y, int incy,
        double               *A, int lda)
cublasStatus_t cublasCgeru(cublasHandle_t handle, int m, int n,
        const cuComplex      *alpha,
        const cuComplex      *x, int incx,
        const cuComplex      *y, int incy,
        cuComplex            *A, int lda)
cublasStatus_t cublasCgerc(cublasHandle_t handle, int m, int n,
        const cuComplex      *alpha,
        const cuComplex      *x, int incx,
        const cuComplex      *y, int incy,
        cuComplex            *A, int lda)
cublasStatus_t cublasZgeru(cublasHandle_t handle, int m, int n,
        const cuDoubleComplex *alpha,
        const cuDoubleComplex *x, int incx,
        const cuDoubleComplex *y, int incy,
        cuDoubleComplex *A, int lda)
cublasStatus_t cublasZgerc(cublasHandle_t handle, int m, int n,
        const cuDoubleComplex *alpha,
        const cuDoubleComplex *x, int incx,
        const cuDoubleComplex *y, int incy,
        cuDoubleComplex *A, int lda)

```

This function performs the rank-1 update

$$A = \begin{cases} \alpha \mathbf{xy}^T + A & \text{if } \text{ger}(), \text{geru}() \text{ is called} \\ \alpha \mathbf{xy}^H + A & \text{if } \text{gerc}() \text{ is called} \end{cases}$$

where A is a $m \times n$ matrix stored in column-major format, \mathbf{x} and \mathbf{y} are vectors, and α is a scalar.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
m		input	number of rows of matrix A.
n		input	number of columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with m elements.
incx		input	stride between consecutive elements of x.
y	device	input	<type> vector with n elements.
incy		input	stride between consecutive elements of y.
A	device	in/out	<type> array of dimension lda × n with lda ≥ max(1,m).
lda		input	leading dimension of two-dimensional array used to store matrix A.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters m,n<0 or incx,incy=0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[sger](#), [dger](#), [cgeru](#), [cgerc](#), [zgeru](#), [zgerc](#)

6.4 cublas<t>sbmv()

```

cublasStatus_t cublasSsbmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, int k, const float *alpha,
                           const float *A, int lda,
                           const float *x, int incx,
                           const float *beta, float *y, int incy)
cublasStatus_t cublasDsbmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, int k, const double *alpha,
                           const double *A, int lda,
                           const double *x, int incx,
                           const double *beta, double *y, int incy)

```

This function performs the symmetric banded matrix-vector multiplication

$$\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$$

where A is a $n \times n$ symmetric banded matrix with k subdiagonals and superdiagonals, \mathbf{x} and \mathbf{y} are vectors, and α and β are scalars.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the symmetric banded matrix A is stored column by column, with the main diagonal of the matrix stored in row 1, the first subdiagonal in row 2 (starting at first position), the second subdiagonal in row 3 (starting at first position), etc. So that in general, the element $A(i, j)$ is stored in the memory location $A(1+i-j, j)$ for $j = 1, \dots, n$ and $i \in [j, \min(m, j+k)]$. Also, the elements in the array A that do not conceptually correspond to the elements in the banded matrix (the bottom right $k \times k$ triangle) are not referenced.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the symmetric banded matrix A is stored column by column, with the main diagonal of the matrix stored in row $k+1$, the first superdiagonal in row k (starting at second position), the second superdiagonal in row $k-1$ (starting at third position), etc. So that in general, the element $A(i, j)$ is stored in the memory location $A(1+k+i-j, j)$ for $j = 1, \dots, n$ and $i \in [\max(1, j-k), j]$. Also, the elements in the array A that do not conceptually correspond to the elements in the banded matrix (the top left $k \times k$ triangle) are not referenced.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A .
k		input	number of sub- and super-diagonals of matrix A .
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $lda \times n$ with $lda \geq k+1$.
lda		input	leading dimension of two-dimensional array used to store matrix A .
x	device	input	<type> vector with n elements.
incx		input	stride between consecutive elements of x .
beta	host or device	input	<type> scalar used for multiplication, if <code>beta==0</code> then y does not have to be a valid input.
y	device	in/out	<type> vector with n elements.
incy		input	stride between consecutive elements of y .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters $n, k < 0$ or <code>incx, incy=0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[ssbmv](#), [dsbmv](#)

6.5 cublas<t>spmv()

```

cublasStatus_t cublasSspmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const float *alpha, const float *AP,
                           const float *x, int incx, const float *beta,
                           float *y, int incy)
cublasStatus_t cublasDspmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const double *alpha, const double *AP,
                           const double *x, int incx, const double *beta,
                           double *y, int incy)

```

This function performs the symmetric packed matrix-vector multiplication

$$\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$$

where A is a $n \times n$ symmetric matrix stored in packed format, \mathbf{x} and \mathbf{y} are vectors, and α and β are scalars.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the symmetric matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+((2*n-j+1)*j)/2]` for $j = 1, \dots, n$ and $i \geq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the symmetric matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+(j*(j+1))/2]` for $j = 1, \dots, n$ and $i \leq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A .
k		input	number of sub- and super-diagonals of matrix A .
alpha	host or device	input	<type> scalar used for multiplication.
AP	device	input	<type> array with A stored in packed format.
x	device	input	<type> vector with <code>n</code> elements.
incx		input	stride between consecutive elements of \mathbf{x} .
beta	host or device	input	<type> scalar used for multiplication, if <code>beta==0</code> then \mathbf{y} does not have to be a valid input.
y	device	input	<type> vector with <code>n</code> elements.
incy		input	stride between consecutive elements of \mathbf{y} .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>n<0</code> or <code>incx,incy=0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[sspmv](#), [dspmv](#)

6.6 cublas<t>spr()

```
cublasStatus_t cublasSspr(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const float *alpha,
                          const float *x, int incx, float *AP)
cublasStatus_t cublasDspr(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const double *alpha,
                          const double *x, int incx, double *AP)
```

This function performs the packed symmetric rank-1 update

$$A = \alpha \mathbf{x} \mathbf{x}^T + A$$

where A is a $n \times n$ symmetric matrix stored in packed format, \mathbf{x} is a vector, and α is a scalar.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the symmetric matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+((2*n-j+1)*j)/2]` for $j = 1, \dots, n$ and $i \geq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the symmetric matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+(j*(j+1))/2]` for $j = 1, \dots, n$ and $i \leq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A .
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with <code>n</code> elements.
incx		input	stride between consecutive elements of \mathbf{x} .
AP	device	in/out	<type> array with A stored in packed format.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>n<0</code> or <code>incx,incy=0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[sspr](#), [dspr](#)

6.7 cublas<t>spr2()

```
cublasStatus_t cublasSspr2(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const float *alpha,
                          const float *x, int incx,
                          const float *y, int incy, float *AP)
cublasStatus_t cublasDspr2(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const double *alpha,
                          const double *x, int incx,
                          const double *y, int incy, double *AP)
```

This function performs the packed symmetric rank-2 update

$$A = \alpha (\mathbf{x}\mathbf{y}^T + \mathbf{y}\mathbf{x}^T) + A$$

where A is a $n \times n$ symmetric matrix stored in packed format, \mathbf{x} and \mathbf{y} are vectors, and α is a scalar.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the symmetric matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+((2*n-j+1)*j)/2]` for $j = 1, \dots, n$ and $i \geq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the symmetric matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+(j*(j+1))/2]` for $j = 1, \dots, n$ and $i \leq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A .
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with n elements.
incx		input	stride between consecutive elements of x .
y	device	input	<type> vector with n elements.
incy		input	stride between consecutive elements of y .
AP	device	in/out	<type> array with A stored in packed format.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters $n < 0$ or $incx, incy = 0$
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[sspr2](#), [dspr2](#)

6.8 cublas<t>symv()

```
cublasStatus_t cublasSsymv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const float *alpha,
                           const float *A, int lda,
                           const float *x, int incx, const float *beta,
                           float *y, int incy)
```

```
cublasStatus_t cublasDsymv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const double *alpha,
                           const double *A, int lda,
                           const double *x, int incx, const double *beta,
                           double *y, int incy)
```

```
cublasStatus_t cublasCsymv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n,
                           const cuComplex *alpha, /* host or device pointer */
                           const cuComplex *A, int lda,
                           const cuComplex *x, int incx, const cuComplex *beta,
                           cuComplex *y, int incy)
```

```

cublasStatus_t cublasZsymv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *x, int incx,
                           const cuDoubleComplex *beta,
                           cuDoubleComplex *y, int incy)

```

This function performs the symmetric matrix-vector multiplication.

$$\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$$

where A is a $n \times n$ symmetric matrix stored in lower or upper mode, \mathbf{x} and \mathbf{y} are vectors, and α and β are scalars.

This function has an alternate faster implementation using atomics that can be enabled with `cublasSetAtomicMode()`

Please see the section on the `cublasSetAtomicMode()` function for more details about the usage of atomics

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $\text{lda} \times n$ with $\text{lda} \geq \max(1, n)$.
lda		input	leading dimension of two-dimensional array used to store matrix A.
x	device	input	<type> vector with n elements.
incx		input	stride between consecutive elements of \mathbf{x} .
beta	host or device	input	<type> scalar used for multiplication, if <code>beta==0</code> then \mathbf{y} does not have to be a valid input.
y	device	in/out	<type> vector with n elements.
incy		input	stride between consecutive elements of \mathbf{y} .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters $n < 0$ or $\text{incx}, \text{incy} = 0$
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[ssymv](#), [dsymv](#)

6.9 cublas<t>syr()

```
cublasStatus_t cublasSsyr(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const float *alpha,
                          const float *x, int incx, float *A, int lda)
cublasStatus_t cublasDsyr(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const double *alpha,
                          const double *x, int incx, double *A, int lda)
cublasStatus_t cublasCsyr(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const cuComplex *alpha,
                          const cuComplex *x, int incx, cuComplex *A, int lda)
cublasStatus_t cublasZsyr(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const cuDoubleComplex *alpha,
                          const cuDoubleComplex *x, int incx, cuDoubleComplex *A, int lda)
```

This function performs the symmetric rank-1 update

$$A = \alpha \mathbf{x} \mathbf{x}^T + A$$

where A is a $n \times n$ symmetric matrix stored in column-major format, \mathbf{x} is a vector, and α is a scalar.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part, is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A .
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with n elements.
incx		input	stride between consecutive elements of \mathbf{x} .
A	device	in/out	<type> array of dimensions $lda \times n$, with $lda \geq \max(1, n)$.
lda		input	leading dimension of two-dimensional array used to store matrix A .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>n<0</code> or <code>incx=0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[ssyr](#), [dsyr](#)

6.10 cublas<t>syr2()

```

cublasStatus_t cublasSsyr2(cublasHandle_t handle, cublasFillMode_t uplo, int n,
                          const float *alpha, const float *x, int incx,
                          const float *y, int incy, float *A, int lda)
cublasStatus_t cublasDsyr2(cublasHandle_t handle, cublasFillMode_t uplo, int n,
                          const double *alpha, const double *x, int incx,
                          const double *y, int incy, double *A, int lda)
cublasStatus_t cublasCsyr2(cublasHandle_t handle, cublasFillMode_t uplo, int n,
                          const cuComplex *alpha, const cuComplex *x, int incx,
                          const cuComplex *y, int incy, cuComplex *A, int lda)
cublasStatus_t cublasZsyr2(cublasHandle_t handle, cublasFillMode_t uplo, int n,
                          const cuDoubleComplex *alpha, const cuDoubleComplex *x, int incx,
                          const cuDoubleComplex *y, int incy, cuDoubleComplex *A, int lda)

```

This function performs the symmetric rank-2 update

$$A = \alpha (\mathbf{x}\mathbf{y}^T + \mathbf{y}\mathbf{x}^T) + A$$

where A is a $n \times n$ symmetric matrix stored in column-major format, \mathbf{x} and \mathbf{y} are vectors, and α is a scalar.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with n elements.
incx		input	stride between consecutive elements of x.
y	device	input	<type> vector with n elements.
incy		input	stride between consecutive elements of y.
A	device	in/out	<type> array of dimensions lda × n, with $lda \geq \max(1,n)$.
lda		input	leading dimension of two-dimensional array used to store matrix A.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters n<0 or incx,incy=0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[ssyr2](#), [dsyr2](#)

6.11 cublas<t>tbmv()

```

cublasStatus_t cublasStbmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const float      *A, int lda,
                           float      *x, int incx)
cublasStatus_t cublasDtbmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const double     *A, int lda,
                           double     *x, int incx)
cublasStatus_t cublasCtbmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const cuComplex  *A, int lda,
                           cuComplex  *x, int incx)
cublasStatus_t cublasZtbmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,

```

```
int n, int k, const cuDoubleComplex *A, int lda,
cuDoubleComplex *x, int incx)
```

This function performs the triangular banded matrix-vector multiplication

$$\mathbf{x} = \text{op}(A)\mathbf{x}$$

where A is a triangular banded matrix, and \mathbf{x} is a vector. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if trans == CUBLAS_OP_N} \\ A^T & \text{if trans == CUBLAS_OP_T} \\ A^H & \text{if trans == CUBLAS_OP_C} \end{cases}$$

If `uplo == CUBLAS_FILL_MODE_LOWER` then the triangular banded matrix A is stored column by column, with the main diagonal of the matrix stored in row 1, the first subdiagonal in row 2 (starting at first position), the second subdiagonal in row 3 (starting at first position), etc. So that in general, the element $A(i, j)$ is stored in the memory location $A(1+i-j, j)$ for $j = 1, \dots, n$ and $i \in [j, \min(m, j+k)]$. Also, the elements in the array A that do not conceptually correspond to the elements in the banded matrix (the bottom right $k \times k$ triangle) are not referenced.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the triangular banded matrix A is stored column by column, with the main diagonal of the matrix stored in row $k+1$, the first superdiagonal in row k (starting at second position), the second superdiagonal in row $k-1$ (starting at third position), etc. So that in general, the element $A(i, j)$ is stored in the memory location $A(1+k+i-j, j)$ for $j = 1, \dots, n$ and $i \in [\max(1, j-k), j]$. Also, the elements in the array A that do not conceptually correspond to the elements in the banded matrix (the top left $k \times k$ triangle) are not referenced.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	operation $\text{op}(A)$ that is non- or (conj.) transpose.
diag		input	indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
n			number of rows and columns of matrix A .
k		input	number of sub- and super-diagonals of matrix A .
A	device	input	<type> array of dimension $\text{lda} \times n$, with $\text{lda} \geq k+1$.
lda		input	leading dimension of two-dimensional array used to store matrix A .
x	device	in/out	<type> vector with n elements.
incx		input	stride between consecutive elements of \mathbf{x} .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>n</code> , <code>k</code> < 0 or <code>incx</code> = 0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_ALLOC_FAILED	the allocation of internal scratch memory failed
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[stbmv](#), [dtbmv](#), [ctbmv](#), [ztbmv](#)

6.12 cublas<t>tbsv()

```

cublasStatus_t cublasStbsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const float      *A, int lda,
                           float      *x, int incx)
cublasStatus_t cublasDtbsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const double     *A, int lda,
                           double     *x, int incx)
cublasStatus_t cublasCtbsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const cuComplex  *A, int lda,
                           cuComplex  *x, int incx)
cublasStatus_t cublasZtbsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, int k, const cuDoubleComplex *A, int lda,
                           cuDoubleComplex *x, int incx)

```

This function solves the triangular banded linear system with a single right-hand-side

$$\text{op}(A)\mathbf{x} = \mathbf{b}$$

where A is a triangular banded matrix, and \mathbf{x} and \mathbf{b} are vectors. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if trans == CUBLAS_OP_N} \\ A^T & \text{if trans == CUBLAS_OP_T} \\ A^H & \text{if trans == CUBLAS_OP_C} \end{cases}$$

The solution \mathbf{x} overwrites the right-hand-sides \mathbf{b} on exit.

No test for singularity or near-singularity is included in this function.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the triangular banded matrix A is stored column by column, with the main diagonal of the matrix stored in row 1, the first subdiagonal in row 2 (starting at first position), the second subdiagonal in row 3 (starting at first position), etc. So that in general, the element $A(i, j)$ is stored in the memory location $A(1+i-j, j)$ for $j = 1, \dots, n$ and $i \in [j, \min(m, j+k)]$. Also, the elements in the array **A** that do not conceptually correspond to the elements in the banded matrix (the bottom right $k \times k$ triangle) are not referenced.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the triangular banded matrix A is stored column by column, with the main diagonal of the matrix stored in row $k+1$, the first superdiagonal in row k (starting at second position), the second superdiagonal in row $k-1$ (starting at third position), etc. So that in general, the element $A(i, j)$ is stored in the memory location $A(1+k+i-j, j)$ for $j = 1, \dots, n$ and $i \in [\max(1, j-k), j]$. Also, the elements in the array **A** that do not conceptually correspond to the elements in the banded matrix (the top left $k \times k$ triangle) are not referenced.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	operation <code>op(A)</code> that is non- or (conj.) transpose.
diag		input	indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
n			number of rows and columns of matrix A .
k		input	number of sub- and super-diagonals of matrix A .
A	device	input	<type> array of dimension $lda \times n$, with $lda \geq k+1$.
lda			leading dimension of two-dimensional array used to store matrix A .
x	device	in/out	<type> vector with n elements.
incx		input	stride between consecutive elements of x .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>n<0</code> or <code>incx=0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[stbsv](#), [dtbsv](#), [ctbsv](#), [ztbsv](#)

6.13 cublas<t>tpmv()

```

cublasStatus_t cublasStpmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const float      *AP,
                           float      *x, int incx)
cublasStatus_t cublasDtpmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const double     *AP,
                           double     *x, int incx)
cublasStatus_t cublasCtpmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const cuComplex  *AP,
                           cuComplex  *x, int incx)
cublasStatus_t cublasZtpmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const cuDoubleComplex *AP,
                           cuDoubleComplex *x, int incx)

```

This function performs the triangular packed matrix-vector multiplication

$$\mathbf{x} = \text{op}(A)\mathbf{x}$$

where A is a triangular matrix stored in packed format, and \mathbf{x} is a vector. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if trans == CUBLAS_OP_N} \\ A^T & \text{if trans == CUBLAS_OP_T} \\ A^H & \text{if trans == CUBLAS_OP_C} \end{cases}$$

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the triangular matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+((2*n-j+1)*j)/2]` for $j = 1, \dots, n$ and $i \geq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the triangular matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+(j*(j+1))/2]` for $j = 1, \dots, n$ and $i \leq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	operation op(A) that is non- or (conj.) transpose.
diag		input	indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
n		input	number of rows and columns of matrix A .
AP	device	input	<type> array with A stored in packed format.
x	device	in/out	<type> vector with n elements.
incx		input	stride between consecutive elements of x .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters n <0 or incx =0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_ALLOC_FAILED	the allocation of internal scratch memory failed
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[stpmv](#), [dtpmv](#), [ctpmv](#), [ztpmv](#)

6.14 cublas<t>tpsv()

```

cublasStatus_t cublasStpsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const float      *AP,
                           float      *x, int incx)
cublasStatus_t cublasDtpsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const double     *AP,
                           double     *x, int incx)
cublasStatus_t cublasCtpsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const cuComplex  *AP,
                           cuComplex  *x, int incx)
cublasStatus_t cublasZtpsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const cuDoubleComplex *AP,
                           cuDoubleComplex *x, int incx)

```

This function solves the packed triangular linear system with a single right-hand-side

$$\text{op}(A)\mathbf{x} = \mathbf{b}$$

where A is a triangular matrix stored in packed format, and \mathbf{x} and \mathbf{b} are vectors. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if trans == CUBLAS_OP_N} \\ A^T & \text{if trans == CUBLAS_OP_T} \\ A^H & \text{if trans == CUBLAS_OP_C} \end{cases}$$

The solution \mathbf{x} overwrites the right-hand-sides \mathbf{b} on exit.

No test for singularity or near-singularity is included in this function.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the triangular matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+((2*n-j+1)*j)/2]` for $j = 1, \dots, n$ and $i \geq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the triangular matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+(j*(j+1))/2]` for $j = 1, \dots, n$ and $i \leq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	operation $\text{op}(A)$ that is non- or (conj.) transpose.
diag		input	indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
n		input	number of rows and columns of matrix A .
AP	device	input	<type> array with A stored in packed format.
x	device	in/out	<type> vector with n elements.
incx		input	stride between consecutive elements of \mathbf{x} .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>n<0</code> or <code>incx=0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[stpsv](#), [dtpsv](#), [ctpsv](#), [ztpsv](#)

6.15 cublas<t>trmv()

```

cublasStatus_t cublasStrmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const float      *A, int lda,
                           float      *x, int incx)
cublasStatus_t cublasDtrmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const double     *A, int lda,
                           double     *x, int incx)
cublasStatus_t cublasCtrmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const cuComplex  *A, int lda,
                           cuComplex  *x, int incx)
cublasStatus_t cublasZtrmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const cuDoubleComplex *A, int lda,
                           cuDoubleComplex *x, int incx)

```

This function performs the triangular matrix-vector multiplication

$$\mathbf{x} = \text{op}(\mathbf{A})\mathbf{x}$$

where \mathbf{A} is a triangular matrix stored in lower or upper mode with or without the main diagonal, and \mathbf{x} is a vector. Also, for matrix \mathbf{A}

$$\text{op}(\mathbf{A}) = \begin{cases} \mathbf{A} & \text{if trans == CUBLAS_OP_N} \\ \mathbf{A}^T & \text{if trans == CUBLAS_OP_T} \\ \mathbf{A}^H & \text{if trans == CUBLAS_OP_C} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix \mathbf{A} lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	operation $\text{op}(\mathbf{A})$ that is non- or (conj.) transpose.
diag		input	indicates if the elements on the main diagonal of matrix \mathbf{A} are unity and should not be accessed.
n		input	number of rows and columns of matrix \mathbf{A} .
A	device	input	<type> array of dimensions $\text{lda} \times \text{n}$, with $\text{lda} \geq \max(1, \text{n})$.
lda		input	leading dimension of two-dimensional array used to store matrix \mathbf{A} .
x	device	in/out	<type> vector with n elements.
incx		input	stride between consecutive elements of \mathbf{x} .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>n<0</code> or <code>incx=0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_ALLOC_FAILED	the allocation of internal scratch memory failed
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[strmv](#), [dtrmv](#), [ctrmv](#), [ztrmv](#)

6.16 cublas<t>trsv()

```

cublasStatus_t cublasStrsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const float      *A, int lda,
                           float      *x, int incx)
cublasStatus_t cublasDtrsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const double     *A, int lda,
                           double     *x, int incx)
cublasStatus_t cublasCtrsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const cuComplex  *A, int lda,
                           cuComplex  *x, int incx)
cublasStatus_t cublasZtrsv(cublasHandle_t handle, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int n, const cuDoubleComplex *A, int lda,
                           cuDoubleComplex *x, int incx)

```

This function solves the triangular linear system with a single right-hand-side

$$\text{op}(A)\mathbf{x} = \mathbf{b}$$

where A is a triangular matrix stored in lower or upper mode with or without the main diagonal, and \mathbf{x} and \mathbf{b} are vectors. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if trans == CUBLAS_OP_N} \\ A^T & \text{if trans == CUBLAS_OP_T} \\ A^H & \text{if trans == CUBLAS_OP_C} \end{cases}$$

The solution \mathbf{x} overwrites the right-hand-sides \mathbf{b} on exit.

No test for singularity or near-singularity is included in this function.

Param.	Memory	In/out	Meaning
handle	device	input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	operation $\text{op}(\mathbf{A})$ that is non- or (conj.) transpose.
diag		input	indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
n		input	number of rows and columns of matrix A .
A	device	input	<type> array of dimension $\text{lda} \times \text{n}$, with $\text{lda} \geq \max(1, \text{n})$.
lda		input	leading dimension of two-dimensional array used to store matrix A .
x	device	in/out	<type> vector with n elements.
incx		input	stride between consecutive elements of x .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters $\text{n} < 0$ or $\text{incx} = 0$
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[strsv](#), [dtrsv](#), [ctrsv](#), [ztrsv](#)

6.17 cublas<t>hemv()

```

cublasStatus_t cublasChemv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuComplex      *alpha,
                           const cuComplex      *A, int lda,
                           const cuComplex      *x, int incx,
                           const cuComplex      *beta,
                           cuComplex      *y, int incy)
cublasStatus_t cublasZhemv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *x, int incx,
                           const cuDoubleComplex *beta,
                           cuDoubleComplex *y, int incy)

```


This function performs the Hermitian matrix-vector multiplication

$$\mathbf{y} = \alpha A \mathbf{x} + \beta \mathbf{y}$$

where A is a $n \times n$ Hermitian matrix stored in lower or upper mode, \mathbf{x} and \mathbf{y} are vectors, and α and β are scalars.

This function has an alternate faster implementation using atomics that can be enabled with `cublasSetAtomicsMode()`

Please see the section on the `cublasSetAtomicsMode()` for more details about the usage of atomics

Param.	Memory	In/out	Meaning
handle	host or device device	input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A .
alpha		input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $\text{lda} \times \text{n}$, with $\text{lda} \geq \max(1, \text{n})$ The imaginary parts of the diagonal elements are assumed to be zero.
lda		input	leading dimension of two-dimensional array used to store matrix A .
x	device	input	<type> vector with n elements.
incx	host or device	input	stride between consecutive elements of \mathbf{x} .
beta		input	<type> scalar used for multiplication, if <code>beta==0</code> then \mathbf{y} does not have to be a valid input.
y	device	in/out	<type> vector with n elements.
incy		input	stride between consecutive elements of \mathbf{y} .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>n<0</code> or <code>incx,incy=0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[chemv](#), [zhemv](#)

6.18 cublas<t>hbmV()

```
cublasStatus_t cublasChbmV(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, int k, const cuComplex      *alpha,
```

```

        const cuComplex      *A, int lda,
        const cuComplex      *x, int incx,
        const cuComplex      *beta,
        cuComplex             *y, int incy)
cublasStatus_t cublasZhbmvc(cublasHandle_t handle, cublasFillMode_t uplo,
        int n, int k, const cuDoubleComplex *alpha,
        const cuDoubleComplex *A, int lda,
        const cuDoubleComplex *x, int incx,
        const cuDoubleComplex *beta,
        cuDoubleComplex *y, int incy)

```

This function performs the Hermitian banded matrix-vector multiplication

$$\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$$

where A is a $n \times n$ Hermitian banded matrix with k subdiagonals and superdiagonals, \mathbf{x} and \mathbf{y} are vectors, and α and β are scalars.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the Hermitian banded matrix A is stored column by column, with the main diagonal of the matrix stored in row 1, the first subdiagonal in row 2 (starting at first position), the second subdiagonal in row 3 (starting at first position), etc. So that in general, the element $A(i, j)$ is stored in the memory location $A(1+i-j, j)$ for $j = 1, \dots, n$ and $i \in [j, \min(m, j+k)]$. Also, the elements in the array A that do not conceptually correspond to the elements in the banded matrix (the bottom right $k \times k$ triangle) are not referenced.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the Hermitian banded matrix A is stored column by column, with the main diagonal of the matrix stored in row $k+1$, the first superdiagonal in row k (starting at second position), the second superdiagonal in row $k-1$ (starting at third position), etc. So that in general, the element $A(i, j)$ is stored in the memory location $A(1+k+i-j, j)$ for $j = 1, \dots, n$ and $i \in [\max(1, j-k), j]$. Also, the elements in the array A that do not conceptually correspond to the elements in the banded matrix (the top left $k \times k$ triangle) are not referenced.

Param.	Memory	In/out	Meaning
handle	host or device device	input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A.
k		input	number of sub- and super-diagonals of matrix A.
alpha		input	<type> scalar used for multiplication.
A	device	input	<type> array of dimensions $lda \times n$, with $lda \geq k+1$. The imaginary parts of the diagonal elements are assumed to be zero.
lda		input	leading dimension of two-dimensional array used to store matrix A.
x		input	<type> vector with n elements.
incx	host or device	input	stride between consecutive elements of x.
beta		input	<type> scalar used for multiplication, if <code>beta==0</code> then y does not have to be a valid input.
y	device	in/out	<type> vector with n elements.
incy		input	stride between consecutive elements of y.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters n,k<0 or incx,incy=0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[chbm](#), [zhbm](#)

6.19 cublas<t>hpmv()

```

cublasStatus_t cublasChpmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuComplex      *alpha,
                           const cuComplex      *AP,
                           const cuComplex      *x, int incx,
                           const cuComplex      *beta,
                           cuComplex      *y, int incy)
cublasStatus_t cublasZhpmv(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuDoubleComplex *alpha,
                           const cuDoubleComplex *AP,
                           const cuDoubleComplex *x, int incx,

```

```
const cuDoubleComplex *beta,
cuDoubleComplex *y, int incy)
```

This function performs the Hermitian packed matrix-vector multiplication

$$\mathbf{y} = \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$$

where A is a $n \times n$ Hermitian matrix stored in packed format, \mathbf{x} and \mathbf{y} are vectors, and α and β are scalars.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the Hermitian matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+((2*n-j+1)*j)/2]` for $j = 1, \dots, n$ and $i \geq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the Hermitian matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+(j*(j+1))/2]` for $j = 1, \dots, n$ and $i \leq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A .
k		input	number of sub- and super-diagonals of matrix A .
alpha	host or device	input	<type> scalar used for multiplication.
AP	device	input	<type> array with A stored in packed format The imaginary parts of the diagonal elements are assumed to be zero.
x	device	input	<type> vector with <code>n</code> elements.
incx		input	stride between consecutive elements of \mathbf{x} .
beta	host or device	input	<type> scalar used for multiplication, if <code>beta==0</code> then \mathbf{y} does not have to be a valid input.
y	device	in/out	<type> vector with <code>n</code> elements.
incy		input	stride between consecutive elements of \mathbf{y} .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>n<0</code> or <code>incx, incy=0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[chpmv](#), [zhpmv](#)

6.20 cublas<t>her()

```

cublasStatus_t cublasCher(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const float  *alpha,
                          const cuComplex    *x, int incx,
                          cuComplex          *A, int lda)
cublasStatus_t cublasZher(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const double *alpha,
                          const cuDoubleComplex *x, int incx,
                          cuDoubleComplex *A, int lda)

```

This function performs the Hermitian rank-1 update

$$A = \alpha \mathbf{x} \mathbf{x}^H + A$$

where A is a $n \times n$ Hermitian matrix stored in column-major format, \mathbf{x} is a vector, and α is a scalar.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with n elements.
incx		input	stride between consecutive elements of x.
A	device	in/out	<type> array of dimensions $\text{lda} \times \text{n}$, with $\text{lda} \geq \max(1, \text{n})$. The imaginary parts of the diagonal elements are assumed and set to zero.
lda		input	leading dimension of two-dimensional array used to store matrix A.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters $\text{n} < 0$ or $\text{incx} = 0$
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[cher](#), [zher](#)

6.21 cublas<t>her2()

```

cublasStatus_t cublasCher2(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuComplex      *alpha,
                           const cuComplex      *x, int incx,
                           const cuComplex      *y, int incy,
                           cuComplex      *A, int lda)
cublasStatus_t cublasZher2(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuDoubleComplex *alpha,
                           const cuDoubleComplex *x, int incx,
                           const cuDoubleComplex *y, int incy,
                           cuDoubleComplex *A, int lda)

```

This function performs the Hermitian rank-2 update

$$A = \alpha \mathbf{x} \mathbf{y}^H + \bar{\alpha} \mathbf{y} \mathbf{x}^H + A$$

where A is a $n \times n$ Hermitian matrix stored in column-major format, \mathbf{x} and \mathbf{y} are vectors, and α is a scalar.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A.
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with n elements.
incx		input	stride between consecutive elements of x.
y	device	input	<type> vector with n elements.
incy		input	stride between consecutive elements of y.
A	device	in/out	<type> array of dimension $\text{lda} \times n$ with $\text{lda} \geq \max(1, n)$. The imaginary parts of the diagonal elements are assumed and set to zero.
lda		input	leading dimension of two-dimensional array used to store matrix A.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters $n < 0$ or $\text{incx}, \text{incy} = 0$
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[cher2](#), [zher2](#)

6.22 cublas<t>hpr()

```

cublasStatus_t cublasChpr(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const float *alpha,
                          const cuComplex *x, int incx,
                          cuComplex *AP)
cublasStatus_t cublasZhpr(cublasHandle_t handle, cublasFillMode_t uplo,
                          int n, const double *alpha,
                          const cuDoubleComplex *x, int incx,
                          cuDoubleComplex *AP)

```

This function performs the packed Hermitian rank-1 update

$$A = \alpha \mathbf{x} \mathbf{x}^H + A$$

where A is a $n \times n$ Hermitian matrix stored in packed format, \mathbf{x} is a vector, and α is a scalar.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the Hermitian matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+((2*n-j+1)*j)/2]` for $j = 1, \dots, n$ and $i \geq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the Hermitian matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+(j*(j+1))/2]` for $j = 1, \dots, n$ and $i \leq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A .
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with <code>n</code> elements.
incx		input	stride between consecutive elements of \mathbf{x} .
AP	device	in/out	<type> array with A stored in packed format The imaginary parts of the diagonal elements are assumed and set to zero.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>n<0</code> or <code>incx=0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[chpr](#), [zhpr](#)

6.23 cublas<t>hpr2()

```

cublasStatus_t cublasChpr2(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuComplex      *alpha,
                           const cuComplex      *x, int incx,
                           const cuComplex      *y, int incy,
                           cuComplex      *AP)
cublasStatus_t cublasZhpr2(cublasHandle_t handle, cublasFillMode_t uplo,
                           int n, const cuDoubleComplex *alpha,
                           const cuDoubleComplex *x, int incx,
                           const cuDoubleComplex *y, int incy,
                           cuDoubleComplex *AP)

```

This function performs the packed Hermitian rank-2 update

$$A = \alpha \mathbf{x} \mathbf{y}^H + \bar{\alpha} \mathbf{y} \mathbf{x}^H + A$$

where A is a $n \times n$ Hermitian matrix stored in packed format, \mathbf{x} and \mathbf{y} are vectors, and α is a scalar.

If `uplo == CUBLAS_FILL_MODE_LOWER` then the elements in the lower triangular part of the Hermitian matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+((2*n-j+1)*j)/2]` for $j = 1, \dots, n$ and $i \geq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

If `uplo == CUBLAS_FILL_MODE_UPPER` then the elements in the upper triangular part of the Hermitian matrix A are packed together column by column without gaps, so that the element $A(i, j)$ is stored in the memory location `AP[i+(j*(j+1))/2]` for $j = 1, \dots, n$ and $i \leq j$. Consequently, the packed format requires only $\frac{n(n+1)}{2}$ elements for storage.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
n		input	number of rows and columns of matrix A .
alpha	host or device	input	<type> scalar used for multiplication.
x	device	input	<type> vector with n elements.
incx		input	stride between consecutive elements of x .
y	device	input	<type> vector with n elements.
incy		input	stride between consecutive elements of y .
AP	device	in/out	<type> array with A stored in packed format The imaginary parts of the diagonal elements are assumed and set to zero.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters $n < 0$ or $incx, incy = 0$
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[chpr2](#), [zhpr2](#)

Chapter 7

CUBLAS Level-3 Function Reference

In this chapter we describe the Level-3 Basic Linear Algebra Subprograms (BLAS3) functions that perform matrix-matrix operations.

7.1 cublas<t>gemm()

```
cublasStatus_t cublasSgemm(cublasHandle_t handle,
                           cublasOperation_t transa, cublasOperation_t transb,
                           int m, int n, int k,
                           const float          *alpha,
                           const float          *A, int lda,
                           const float          *B, int ldb,
                           const float          *beta,
                           float                *C, int ldc)

cublasStatus_t cublasDgemm(cublasHandle_t handle,
                           cublasOperation_t transa, cublasOperation_t transb,
                           int m, int n, int k,
                           const double         *alpha,
                           const double         *A, int lda,
                           const double         *B, int ldb,
                           const double         *beta,
                           double               *C, int ldc)

cublasStatus_t cublasCgemm(cublasHandle_t handle,
                           cublasOperation_t transa, cublasOperation_t transb,
                           int m, int n, int k,
                           const cuComplex      *alpha,
                           const cuComplex      *A, int lda,
                           const cuComplex      *B, int ldb,
                           const cuComplex      *beta,
                           cuComplex            *C, int ldc)

cublasStatus_t cublasZgemm(cublasHandle_t handle,
                           cublasOperation_t transa, cublasOperation_t transb,
```

```

int m, int n, int k,
const cuDoubleComplex *alpha,
const cuDoubleComplex *A, int lda,
const cuDoubleComplex *B, int ldb,
const cuDoubleComplex *beta,
cuDoubleComplex *C, int ldc)

```

This function performs the matrix-matrix multiplication

$$C = \alpha \text{op}(A) \text{op}(B) + \beta C$$

where α and β are scalars, and A , B and C are matrices stored in column-major format with dimensions $\text{op}(A)$ $m \times k$, $\text{op}(B)$ $k \times n$ and C $m \times n$, respectively. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS_OP_N} \\ A^T & \text{if transa} == \text{CUBLAS_OP_T} \\ A^H & \text{if transa} == \text{CUBLAS_OP_C} \end{cases}$$

and $\text{op}(B)$ is defined similarly for matrix B .

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
transa		input	operation $\text{op}(A)$ that is non- or (conj.) transpose.
transb		input	operation $\text{op}(B)$ that is non- or (conj.) transpose.
m		input	number of rows of matrix $\text{op}(A)$ and C .
n		input	number of columns of matrix $\text{op}(B)$ and C .
k		input	number of columns of $\text{op}(A)$ and rows of $\text{op}(B)$.
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimensions $\text{lda} \times k$ with $\text{lda} \geq \max(1, m)$ if $\text{transa} == \text{CUBLAS_OP_N}$ and $\text{lda} \times m$ with $\text{lda} \geq \max(1, k)$ otherwise.
lda		input	leading dimension of two-dimensional array used to store the matrix A .
B	device	input	<type> array of dimension $\text{ldb} \times n$ with $\text{ldb} \geq \max(1, k)$ if $\text{transa} == \text{CUBLAS_OP_N}$ and $\text{ldb} \times k$ with $\text{ldb} \geq \max(1, n)$ otherwise.
ldb		input	leading dimension of two-dimensional array used to store matrix B .
beta	host or device	input	<type> scalar used for multiplication.
C	device	in/out	If $\text{beta} == 0$, C does not have to be a valid input. <type> array of dimensions $\text{ldc} \times n$ with $\text{ldc} \geq \max(1, m)$.
ldc		input	leading dimension of a two-dimensional array used to store the matrix C .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters m,n,k<0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[sgemm](#), [dgemm](#), [cgemm](#), [zgemm](#)

7.2 cublas<t>gemmBatched()

```

cublasStatus_t cublasSgemmBatched(cublasHandle_t handle,
                                cublasOperation_t transa, cublasOperation_t transb,
                                int m, int n, int k,
                                const float      *alpha,
                                const float      *Aarray[], int lda,
                                const float      *Barray[], int ldb,
                                const float      *beta,
                                float      *Carray[], int ldc, int batchCount)
cublasStatus_t cublasDgemmBatched(cublasHandle_t handle,
                                cublasOperation_t transa, cublasOperation_t transb,
                                int m, int n, int k,
                                const double     *alpha,
                                const double     *Aarray[], int lda,
                                const double     *Barray[], int ldb,
                                const double     *beta,
                                double     *Carray[], int ldc, int batchCount)
cublasStatus_t cublasCgemmBatched(cublasHandle_t handle,
                                cublasOperation_t transa, cublasOperation_t transb,
                                int m, int n, int k,
                                const cuComplex  *alpha,
                                const cuComplex  *Aarray[], int lda,
                                const cuComplex  *Barray[], int ldb,
                                const cuComplex  *beta,
                                cuComplex  *Carray[], int ldc, int batchCount)
cublasStatus_t cublasZgemmBatched(cublasHandle_t handle,
                                cublasOperation_t transa, cublasOperation_t transb,
                                int m, int n, int k,
                                const cuDoubleComplex *alpha,
                                const cuDoubleComplex *Aarray[], int lda,
                                const cuDoubleComplex *Barray[], int ldb,

```

```
const cuDoubleComplex *beta,
cuDoubleComplex *Carray[], int ldc, int batchCount)
```

This function performs the matrix-matrix multiplications of an array of matrices.

$$C[i] = \alpha \text{op}(A[i])\text{op}(B[i]) + \beta C[i], \text{ for } i \in [0, \text{batchCount} - 1]$$

where α and β are scalars, and A , B and C are arrays of pointers to matrices stored in column-major format with dimensions $\text{op}(A[i])$ $m \times k$, $\text{op}(B[i])$ $k \times n$ and $C[i]$ $m \times n$, respectively. Also, for matrix $A[i]$

$$\text{op}(A[i]) = \begin{cases} A[i] & \text{if transa} == \text{CUBLAS_OP_N} \\ A^T[i] & \text{if transa} == \text{CUBLAS_OP_T} \\ A^H[i] & \text{if transa} == \text{CUBLAS_OP_C} \end{cases}$$

and $\text{op}(B[i])$ is defined similarly for matrix $B[i]$.

This function is intended to be used for matrices of small sizes where the launch overhead is a significant factor. For small sizes, typically smaller than 100x100, this function improves significantly performance compared to making calls to its corresponding `cublas<t>gemm` routine. However, on GPU architectures that support concurrent kernels, it might be advantageous to make multiple calls to `cublas<t>gemm` into different streams as the matrix sizes increase.

Param.	Memory	In/out	Meaning
handle	host or device	input	handle to the CUBLAS library context.
transa		input	operation op(A[i]) that is non- or (conj.) transpose.
transb		input	operation op(B[i]) that is non- or (conj.) transpose.
m		input	number of rows of matrix op(A[i]) and C[i].
n		input	number of columns of op(B[i]) and C[i].
k		input	number of columns of op(A[i]) and rows of op(B[i]).
alpha		input	<type> scalar used for multiplication.
A	device	input	array of pointers to <type> array, with each array of dim. $\text{lda} \times k$ with $\text{lda} \geq \max(1, m)$ if <code>transa==CUBLAS_OP_N</code> and $\text{lda} \times m$ with $\text{lda} \geq \max(1, k)$ otherwise.
lda	device	input	leading dimension of two-dimensional array used to store each matrix A[i].
B		input	array of pointers to <type> array, with each array of dim. $\text{ldb} \times n$ with $\text{ldb} \geq \max(1, k)$ if <code>transa==CUBLAS_OP_N</code> and $\text{ldb} \times k$ with $\text{ldb} \geq \max(1, n)$ otherwise.
ldb		input	leading dimension of two-dimensional array used to store each matrix B[i].
beta	host or device	input	<type> scalar used for multiplication If <code>beta == 0</code> , C does not have to be a valid input.
C	device	in/out	array of pointers to <type> array It has dimensions $\text{ldc} \times n$ with $\text{ldc} \geq \max(1, m)$.
ldc		input	leading dimension of two-dimensional array used to store each matrix C[i].
batchCount		input	number of pointers contained in A, B and C.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters m,n,k, batchCount<0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

7.3 cublas<t>symm()

```

cublasStatus_t cublasSsymm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           int m, int n,
                           const float *alpha,
                           const float *A, int lda,
                           const float *B, int ldb,
                           const float *beta,
```

```

float          *C, int ldc)
cublasStatus_t cublasDsymm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           int m, int n,
                           const double      *alpha,
                           const double      *A, int lda,
                           const double      *B, int ldb,
                           const double      *beta,
                           double          *C, int ldc)
cublasStatus_t cublasCsymm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           int m, int n,
                           const cuComplex   *alpha,
                           const cuComplex   *A, int lda,
                           const cuComplex   *B, int ldb,
                           const cuComplex   *beta,
                           cuComplex        *C, int ldc)
cublasStatus_t cublasZsymm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           int m, int n,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *B, int ldb,
                           const cuDoubleComplex *beta,
                           cuDoubleComplex *C, int ldc)

```

This function performs the symmetric matrix-matrix multiplication

$$C = \begin{cases} \alpha AB + \beta C & \text{if side == CUBLAS_SIDE_LEFT} \\ \alpha BA + \beta C & \text{if side == CUBLAS_SIDE_RIGHT} \end{cases}$$

where A is a symmetric matrix stored in lower or upper mode, B and C are $m \times n$ matrices, and α and β are scalars.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
side		input	indicates if matrix A is on the left or right of B.
uplo		input	indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
m		input	number of rows of matrix C and B, with matrix A sized accordingly.
n		input	number of columns of matrix C and B, with matrix A sized accordingly.
alpha	host or device	input	<type> scalar used for multiplication.
A		input	<type> array of dimension $\text{lda} \times \text{m}$ with $\text{lda} \geq \max(1, \text{m})$ if <code>side == CUBLAS_SIDE_LEFT</code> and $\text{lda} \times \text{n}$ with $\text{lda} \geq \max(1, \text{n})$ otherwise.
lda	device		leading dimension of two-dimensional array used to store matrix A.
B			<type> array of dimension $\text{ldb} \times \text{n}$ with $\text{ldb} \geq \max(1, \text{m})$.
ldb			leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	<type> scalar used for multiplication, if <code>beta==0</code> then C does not have to be a valid input.
C	device	in/out	<type> array of dimension $\text{ldc} \times \text{n}$ with $\text{ldc} \geq \max(1, \text{m})$.
ldc		input	leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>m, n < 0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[ssymm](#), [dsymm](#), [csymm](#), [zsymm](#)

7.4 cublas<t>syrk()

```
cublasStatus_t cublasSsyrk(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
```



```

        const float          *alpha,
        const float          *A, int lda,
        const float          *beta,
        float                *C, int ldc)
cublasStatus_t cublasDsyrk(cublasHandle_t handle,
        cublasFillMode_t uplo, cublasOperation_t trans,
        int n, int k,
        const double         *alpha,
        const double         *A, int lda,
        const double         *beta,
        double               *C, int ldc)
cublasStatus_t cublasCsyrk(cublasHandle_t handle,
        cublasFillMode_t uplo, cublasOperation_t trans,
        int n, int k,
        const cuComplex      *alpha,
        const cuComplex      *A, int lda,
        const cuComplex      *beta,
        cuComplex            *C, int ldc)
cublasStatus_t cublasZsyrk(cublasHandle_t handle,
        cublasFillMode_t uplo, cublasOperation_t trans,
        int n, int k,
        const cuDoubleComplex *alpha,
        const cuDoubleComplex *A, int lda,
        const cuDoubleComplex *beta,
        cuDoubleComplex *C, int ldc)

```

This function performs the symmetric rank- k update

$$C = \alpha \text{op}(A) \text{op}(A)^T + \beta C$$

where α and β are scalars, C is a symmetric matrix stored in lower or upper mode, and A is a matrix with dimensions $\text{op}(A)$ $n \times k$. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if trans == CUBLAS_OP_N} \\ A^T & \text{if trans == CUBLAS_OP_T} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other symmetric part is not referenced and is inferred from the stored elements.
trans		input	operation op(A) that is non- or transpose.
n		input	number of rows of matrix op(A) and C.
k		input	number of columns of matrix op(A).
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $lda \times k$ with $lda \geq \max(1,n)$ if <code>transa == CUBLAS_OP_N</code> and $lda \times n$ with $lda \geq \max(1,k)$ otherwise.
lda		input	leading dimension of two-dimensional array used to store matrix A.
beta	host or device	input	<type> scalar used for multiplication, if <code>beta==0</code> then C does not have to be a valid input.
C	device	in/out	<type> array of dimension $ldc \times n$, with $ldc \geq \max(1,n)$.
ldc		input	leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>n,k<0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[ssyrk](#), [dsyrk](#), [csyrk](#), [zsyrk](#)

7.5 cublas<t>syr2k()

```

cublasStatus_t cublasSsyr2k(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const float          *alpha,
                           const float          *A, int lda,
                           const float          *B, int ldb,
                           const float          *beta,
                           float                *C, int ldc)
cublasStatus_t cublasDsyr2k(cublasHandle_t handle,

```

```

        cublasFillMode_t uplo, cublasOperation_t trans,
        int n, int k,
        const double      *alpha,
        const double      *A, int lda,
        const double      *B, int ldb,
        const double      *beta,
        double            *C, int ldc)
cublasStatus_t cublasCsyr2k(cublasHandle_t handle,
        cublasFillMode_t uplo, cublasOperation_t trans,
        int n, int k,
        const cuComplex   *alpha,
        const cuComplex   *A, int lda,
        const cuComplex   *B, int ldb,
        const cuComplex   *beta,
        cuComplex         *C, int ldc)
cublasStatus_t cublasZsyr2k(cublasHandle_t handle,
        cublasFillMode_t uplo, cublasOperation_t trans,
        int n, int k,
        const cuDoubleComplex *alpha,
        const cuDoubleComplex *A, int lda,
        const cuDoubleComplex *B, int ldb,
        const cuDoubleComplex *beta,
        cuDoubleComplex *C, int ldc)

```

This function performs the symmetric rank- $2k$ update

$$C = \alpha (\text{op}(A)\text{op}(B)^T + \text{op}(B)\text{op}(A)^T) + \beta C$$

where α and β are scalars, C is a symmetric matrix stored in lower or upper mode, and A and B are matrices with dimensions $\text{op}(A)$ $n \times k$ and $\text{op}(B)$ $n \times k$, respectively. Also, for matrix A and B

$$\text{op}(A) \text{ and } \text{op}(B) = \begin{cases} A \text{ and } B & \text{if } \text{trans} == \text{CUBLAS_OP_N} \\ A^T \text{ and } B^T & \text{if } \text{trans} == \text{CUBLAS_OP_T} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part, is stored, the other symmetric part is not referenced and is inferred from the stored elements.
trans		input	operation op(A) that is non- or transpose.
n		input	number of rows of matrix op(A), op(B) and C.
k		input	number of columns of matrix op(A) and op(B).
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $\text{lda} \times \text{k}$ with $\text{lda} \geq \max(1, \text{n})$ if <code>transa == CUBLAS_OP_N</code> and $\text{lda} \times \text{n}$ with $\text{lda} \geq \max(1, \text{k})$ otherwise.
lda		input	leading dimension of two-dimensional array used to store matrix A.
B	device	input	<type> array of dimensions $\text{ldb} \times \text{k}$ with $\text{ldb} \geq \max(1, \text{n})$ if <code>transa == CUBLAS_OP_N</code> and $\text{ldb} \times \text{n}$ with $\text{ldb} \geq \max(1, \text{k})$ otherwise.
ldb		input	leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	<type> scalar used for multiplication, if <code>beta==0</code> , then C does not have to be a valid input.
C	device	in/out	<type> array of dimensions $\text{ldc} \times \text{n}$ with $\text{ldc} \geq \max(1, \text{n})$.
ldc		input	leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>n, k < 0</code>
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[ssyr2k](#), [dsyr2k](#), [csyr2k](#), [zsyr2k](#)

7.6 cublas<t>trmm()

```
cublasStatus_t cublasStrmm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int m, int n,
```

```

        const float      *alpha,
        const float      *A, int lda,
        const float      *B, int ldb,
        float            *C, int ldc)
cublasStatus_t cublasDtrmm(cublasHandle_t handle,
        cublasSideMode_t side, cublasFillMode_t uplo,
        cublasOperation_t trans, cublasDiagType_t diag,
        int m, int n,
        const double     *alpha,
        const double     *A, int lda,
        const double     *B, int ldb,
        double           *C, int ldc)
cublasStatus_t cublasCtrmm(cublasHandle_t handle,
        cublasSideMode_t side, cublasFillMode_t uplo,
        cublasOperation_t trans, cublasDiagType_t diag,
        int m, int n,
        const cuComplex *alpha,
        const cuComplex *A, int lda,
        const cuComplex *B, int ldb,
        cuComplex       *C, int ldc)
cublasStatus_t cublasZtrmm(cublasHandle_t handle,
        cublasSideMode_t side, cublasFillMode_t uplo,
        cublasOperation_t trans, cublasDiagType_t diag,
        int m, int n, const
        cuDoubleComplex *alpha,
        const cuDoubleComplex *A, int lda,
        const cuDoubleComplex *B, int ldb,
        cuDoubleComplex     *C, int ldc)

```

This function performs the triangular matrix-matrix multiplication

$$C = \begin{cases} \alpha \text{op}(A)B & \text{if side} == \text{CUBLAS_SIDE_LEFT} \\ \alpha B \text{op}(A) & \text{if side} == \text{CUBLAS_SIDE_RIGHT} \end{cases}$$

where A is a triangular matrix stored in lower or upper mode with or without the main diagonal, B and C are $m \times n$ matrix, and α is a scalar. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if trans} == \text{CUBLAS_OP_N} \\ A^T & \text{if trans} == \text{CUBLAS_OP_T} \\ A^H & \text{if trans} == \text{CUBLAS_OP_C} \end{cases}$$

Notice that in order to achieve better parallelism CUBLAS differs from the BLAS API only for this routine. The BLAS API assumes an in-place implementation (with results written back to B), while the CUBLAS API assumes an out-of-place implementation (with

results written into C). The application can obtain the in-place functionality of BLAS in the CUBLAS API by passing the address of the matrix B in place of the matrix C. No other overlapping in the input parameters is supported.

Param.	Memory	In/out	Meaning
handle	host or device	input	handle to the CUBLAS library context.
side		input	indicates if matrix A is on the left or right of B.
uplo		input	indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	operation op(A) that is non- or (conj.) transpose.
diag		input	indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
m		input	number of rows of matrix B, with matrix A sized accordingly.
n	device	input	number of columns of matrix B, with matrix A sized accordingly.
alpha		input	<type> scalar used for multiplication, if <code>alpha==0</code> then A is not referenced and B does not have to be a valid input.
A		input	<type> array of dimension $lda \times m$ with $lda \geq \max(1,m)$ if <code>side == CUBLAS_SIDE_LEFT</code> and $lda \times n$ with $lda \geq \max(1,n)$ otherwise.
lda		input	leading dimension of two-dimensional array used to store matrix A.
B		input	<type> array of dimension $ldb \times n$ with $ldb \geq \max(1,m)$.
ldb		input	leading dimension of two-dimensional array used to store matrix B.
C	device	in/out	<type> array of dimension $ldc \times n$ with $ldc \geq \max(1,m)$.
ldc		input	leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>m,n</code> <0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[strmm](#), [dtrmm](#), [ctrmm](#), [ztrmm](#)

7.7 cublas<t>trsm()

```

cublasStatus_t cublasStrsm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int m, int n,
                           const float          *alpha,
                           const float          *A, int lda,
                           float                *B, int ldb)
cublasStatus_t cublasDtrsm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int m, int n,
                           const double         *alpha,
                           const double         *A, int lda,
                           double               *B, int ldb)
cublasStatus_t cublasCtrsm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int m, int n,
                           const cuComplex      *alpha,
                           const cuComplex      *A, int lda,
                           cuComplex            *B, int ldb)
cublasStatus_t cublasZtrsm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           cublasOperation_t trans, cublasDiagType_t diag,
                           int m, int n,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           cuDoubleComplex *B, int ldb)

```

This function solves the triangular linear system with multiple right-hand-sides

$$\begin{cases} \text{op}(A)X = \alpha B & \text{if side} == \text{CUBLAS_SIDE_LEFT} \\ X\text{op}(A) = \alpha B & \text{if side} == \text{CUBLAS_SIDE_RIGHT} \end{cases}$$

where A is a triangular matrix stored in lower or upper mode with or without the main diagonal, X and B are $m \times n$ matrices, and α is a scalar. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if trans} == \text{CUBLAS_OP_N} \\ A^T & \text{if trans} == \text{CUBLAS_OP_T} \\ A^H & \text{if trans} == \text{CUBLAS_OP_C} \end{cases}$$

The solution X overwrites the right-hand-sides B on exit.

No test for singularity or near-singularity is included in this function.

Param.	Memory	In/out	Meaning
handle	host or device	input	handle to the CUBLAS library context.
side		input	indicates if matrix A is on the left or right of X .
uplo		input	indicates if matrix A lower or upper part is stored, the other part is not referenced and is inferred from the stored elements.
trans		input	operation $\text{op}(\mathbf{A})$ that is non- or (conj.) transpose.
diag		input	indicates if the elements on the main diagonal of matrix A are unity and should not be accessed.
m		input	number of rows of matrix B, with matrix A sized accordingly.
n	device	input	number of columns of matrix B, with matrix A sized accordingly.
alpha		input	<type> scalar used for multiplication, if <code>alpha==0</code> then A is not referenced and B does not have to be a valid input.
A		input	<type> array of dimension $\text{lda} \times \text{m}$ with $\text{lda} \geq \max(1, \text{m})$ if <code>side == CUBLAS_SIDE_LEFT</code> and $\text{lda} \times \text{n}$ with $\text{lda} \geq \max(1, \text{n})$ otherwise.
lda		input	leading dimension of two-dimensional array used to store matrix A.
B		in/out	<type> array It has dimensions $\text{ldb} \times \text{n}$ with $\text{ldb} \geq \max(1, \text{m})$.
ldb		input	leading dimension of two-dimensional array used to store matrix B.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters <code>m, n</code> < 0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[strsm](#), [dtrsm](#), [ctrsm](#), [ztrsm](#)

7.8 cublas<t>hemm()

```
cublasStatus_t cublasChemm(cublasHandle_t handle,
                           cublasSideMode_t side, cublasFillMode_t uplo,
                           int m, int n,
```



```

        const cuComplex      *alpha,
        const cuComplex      *A, int lda,
        const cuComplex      *B, int ldb,
        const cuComplex      *beta,
        cuComplex            *C, int ldc)
cublasStatus_t cublasZhemm(cublasHandle_t handle,
        cublasSideMode_t side, cublasFillMode_t uplo,
        int m, int n,
        const cuDoubleComplex *alpha,
        const cuDoubleComplex *A, int lda,
        const cuDoubleComplex *B, int ldb,
        const cuDoubleComplex *beta,
        cuDoubleComplex *C, int ldc)

```

This function performs the Hermitian matrix-matrix multiplication

$$C = \begin{cases} \alpha AB + \beta C & \text{if side == CUBLAS_SIDE_LEFT} \\ \alpha BA + \beta C & \text{if side == CUBLAS_SIDE_RIGHT} \end{cases}$$

where A is a Hermitian matrix stored in lower or upper mode, B and C are $m \times n$ matrices, and α and β are scalars.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
side		input	indicates if matrix A is on the left or right of B.
uplo		input	indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
m		input	number of rows of matrix C and B, with matrix A sized accordingly.
n		input	number of columns of matrix C and B, with matrix A sized accordingly.
alpha	host or device	input	<type> scalar used for multiplication.
A		input	<type> array of dimension $\text{lda} \times m$ with $\text{lda} \geq \max(1,m)$ if <code>side==CUBLAS_SIDE_LEFT</code> and $\text{lda} \times n$ with $\text{lda} \geq \max(1,n)$ otherwise. The imaginary parts of the diagonal elements are assumed to be zero.
lda	device	input	leading dimension of two-dimensional array used to store matrix A.
B		input	<type> array of dimension $\text{ldb} \times n$ with $\text{ldb} \geq \max(1,m)$.
ldb		input	leading dimension of two-dimensional array used to store matrix B.
beta		input	<type> scalar used for multiplication, if <code>beta==0</code> then C does not have to be a valid input.
C	device	in/out	<type> array of dimensions $\text{ldc} \times n$ with $\text{ldc} \geq \max(1,m)$.
ldc		input	leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters m,n<0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[chemm](#), [zhemm](#)

7.9 cublas<t>herk()

```

cublasStatus_t cublasCherk(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const float *alpha,
                           const cuComplex *A, int lda,
                           const float *beta,
                           cuComplex *C, int ldc)
cublasStatus_t cublasZherk(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const double *alpha,
                           const cuDoubleComplex *A, int lda,
                           const double *beta,
                           cuDoubleComplex *C, int ldc)

```

This function performs the Hermitian rank- k update

$$C = \alpha \text{op}(A) \text{op}(A)^H + \beta C$$

where α and β are scalars, C is a Hermitian matrix stored in lower or upper mode, and A is a matrix with dimensions $\text{op}(A) \ n \times k$. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if trans} == \text{CUBLAS_OP_N} \\ A^H & \text{if trans} == \text{CUBLAS_OP_C} \end{cases}$$

Param.	Memory	In/out	Meaning
handle	host or device device	input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
trans		input	operation op(A) that is non- or (conj.) transpose.
n		input	number of rows of matrix op(A) and C.
k		input	number of columns of matrix op(A).
alpha		input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $lda \times k$ with $lda \geq \max(1,n)$ if <code>transa == CUBLAS_OP_N</code> and $lda \times n$ with $lda \geq \max(1,k)$ otherwise.
lda		input	leading dimension of two-dimensional array used to store matrix A.
beta		input	<type> scalar used for multiplication, if <code>beta==0</code> then C does not have to be a valid input.
C		in/out	<type> array of dimension $ldc \times n$, with $ldc \geq \max(1,n)$. The imaginary parts of the diagonal elements are assumed and set to zero.
ldc		input	leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters n,k<0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[cherk](#), [zherk](#)

7.10 cublas<t>her2k()

```

cublasStatus_t cublasCher2k(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const cuComplex      *alpha,
                           const cuComplex      *A, int lda,
                           const cuComplex      *B, int ldb,
                           const float          *beta,
                           cuComplex            *C, int ldc)

```

```

cublasStatus_t cublasZher2k(cublasHandle_t handle,
                           cublasFillMode_t uplo, cublasOperation_t trans,
                           int n, int k,
                           const cuDoubleComplex *alpha,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *B, int ldb,
                           const double *beta,
                           cuDoubleComplex *C, int ldc)

```

This function performs the Hermitian rank- $2k$ update

$$C = \alpha \text{op}(A) \text{op}(B)^H + \bar{\alpha} \text{op}(B) \text{op}(A)^H + \beta C$$

where α and β are scalars, C is a Hermitian matrix stored in lower or upper mode, and A and B are matrices with dimensions $\text{op}(A)$ $n \times k$ and $\text{op}(B)$ $n \times k$, respectively. Also, for matrix A and B

$$\text{op}(A) \text{ and } \text{op}(B) = \begin{cases} A \text{ and } B & \text{if } \text{trans} == \text{CUBLAS_OP_N} \\ A^H \text{ and } B^H & \text{if } \text{trans} == \text{CUBLAS_OP_C} \end{cases}$$

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
uplo		input	indicates if matrix A lower or upper part is stored, the other Hermitian part is not referenced and is inferred from the stored elements.
trans		input	operation op(A) that is non- or (conj.) transpose.
n		input	number of rows of matrix op(A), op(B) and C.
k		input	number of columns of matrix op(A) and op(B).
alpha	host or device	input	<type> scalar used for multiplication.
A	device	input	<type> array of dimension $\text{lda} \times k$ with $\text{lda} \geq \max(1, n)$ if $\text{transa} == \text{CUBLAS_OP_N}$ and $\text{lda} \times n$ with $\text{lda} \geq \max(1, k)$ otherwise.
lda		input	leading dimension of two-dimensional array used to store matrix A.
B	device	input	<type> array of dimension $\text{ldb} \times k$ with $\text{ldb} \geq \max(1, n)$ if $\text{transa} == \text{CUBLAS_OP_N}$ and $\text{ldb} \times n$ with $\text{ldb} \geq \max(1, k)$ otherwise.
ldb		input	leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	<type> scalar used for multiplication, if $\text{beta} == 0$ then C does not have to be a valid input.
C	device	in/out	<type> array of dimension $\text{ldc} \times n$, with $\text{ldc} \geq \max(1, n)$. The imaginary parts of the diagonal elements are assumed and set to zero.
ldc		input	leading dimension of two-dimensional array used to store matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters n , k < 0
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

For references please refer to:

[cher2k](#), [zher2k](#)

Chapter 8

BLAS-like Extension

In this chapter we describe the BLAS-extension functions that perform matrix-matrix operations.

8.1 cublas<t>geam()

```
cublasStatus_t cublasSgeam(cublasHandle_t handle,
                           cublasOperation_t transa, cublasOperation_t transb,
                           int m, int n,
                           const float *alpha, const float *A, int lda,
                           const float *beta , const float *B, int ldb,
                           float *C, int ldc)
```

```
cublasStatus_t cublasDgeam(cublasHandle_t handle,
                           cublasOperation_t transa, cublasOperation_t transb,
                           int m, int n,
                           const double *alpha, const double *A, int lda,
                           const double *beta , const double *B, int ldb,
                           double *C, int ldc)
```

```
cublasStatus_t cublasCgeam(cublasHandle_t handle,
                           cublasOperation_t transa, cublasOperation_t transb,
                           int m, int n,
                           const cuComplex *alpha, const cuComplex *A, int lda,
                           const cuComplex *beta , const cuComplex *B, int ldb,
                           cuComplex *C, int ldc)
```

```
cublasStatus_t cublasZgeam(cublasHandle_t handle,
                           cublasOperation_t transa, cublasOperation_t transb,
                           int m, int n,
                           const cuDoubleComplex *alpha, const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *beta , const cuDoubleComplex *B, int ldb,
```

```
cuDoubleComplex *C, int ldc)
```

This function performs the matrix-matrix addition/transposition

$$C = \alpha \text{op}(A) + \beta \text{op}(B)$$

where α and β are scalars, and A , B and C are matrices stored in column-major format with dimensions $\text{op}(A)$ $m \times n$, $\text{op}(B)$ $m \times n$ and C $m \times n$, respectively. Also, for matrix A

$$\text{op}(A) = \begin{cases} A & \text{if transa} == \text{CUBLAS_OP_N} \\ A^T & \text{if transa} == \text{CUBLAS_OP_T} \\ A^H & \text{if transa} == \text{CUBLAS_OP_C} \end{cases}$$

and $\text{op}(B)$ is defined similarly for matrix B .

The operation is out-of-place and CUBLAS would not check range of pointer A, B and C. If C overlaps A or B, then behaviour is undefined. The operation includes the following special cases:

1. the user can reset matrix C to zero by setting `*alpha=*beta=0`.
2. the user can transpose matrix A by setting `*alpha=1` and `*beta=0`.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
transa		input	operation op(A) that is non- or (conj.) transpose.
transb		input	operation op(B) that is non- or (conj.) transpose.
m		input	number of rows of matrix op(A) and C.
n		input	number of columns of matrix op(B) and C.
alpha	host or device	input	<type> scalar used for multiplication. If <code>*alpha == 0</code> , A does not have to be a valid input.
A	device	input	<type> array of dimensions <code>lda × n</code> with <code>lda ≥ max(1,m)</code> if <code>transa == CUBLAS_OP_N</code> and <code>lda × m</code> with <code>lda ≥ max(1,n)</code> otherwise.
lda		input	leading dimension of two-dimensional array used to store the matrix A.
B	device	input	<type> array of dimension <code>ldb × n</code> with <code>ldb ≥ max(1,m)</code> if <code>transa == CUBLAS_OP_N</code> and <code>ldb × m</code> with <code>ldb ≥ max(1,n)</code> otherwise.
ldb		input	leading dimension of two-dimensional array used to store matrix B.
beta	host or device	input	<type> scalar used for multiplication. If <code>*beta == 0</code> , B does not have to be a valid input.
C	device	in/out	<type> array of dimensions <code>ldc × n</code> with <code>ldc ≥ max(1,m)</code> .
ldc		input	leading dimension of a two-dimensional array used to store the matrix C.

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
CUBLAS_STATUS_SUCCESS	the operation completed successfully
CUBLAS_STATUS_NOT_INITIALIZED	the library was not initialized
CUBLAS_STATUS_INVALID_VALUE	the parameters $m, n < 0$ or $\alpha, \beta = \text{NULL}$
CUBLAS_STATUS_ARCH_MISMATCH	the device does not support double-precision
CUBLAS_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

8.2 cublas<t>dggmm()

```
cublasStatus_t cublasSdggmm(cublasHandle_t handle,
                           cublasSideMode_t mode, int m, int n,
                           const float *A, int lda,
                           const float *x, int incx,
                           float *C, int ldc)
```

```
cublasStatus_t cublasDdggmm(cublasHandle_t handle,
                           cublasSideMode_t mode, int m, int n,
                           const double *A, int lda,
                           const double *x, int incx,
                           double *C, int ldc)
```

```
cublasStatus_t cublasCdggmm(cublasHandle_t handle,
                           cublasSideMode_t mode, int m, int n,
                           const cuComplex *A, int lda,
                           const cuComplex *x, int incx,
                           cuComplex *C, int ldc)
```

```
cublasStatus_t cublasZdggmm(cublasHandle_t handle,
                           cublasSideMode_t mode, int m, int n,
                           const cuDoubleComplex *A, int lda,
                           const cuDoubleComplex *x, int incx,
                           cuDoubleComplex *C, int ldc)
```

This function performs the matrix-matrix multiplication

$$C = \begin{cases} A \times \text{diag}(X) & \text{if mode} == \text{CUBLAS_SIDE_RIGHT} \\ \text{diag}(X) \times A & \text{if mode} == \text{CUBLAS_SIDE_LEFT} \end{cases}$$

where A and C are matrices stored in column-major format with dimensions $m \times n$. X is a vector of size n if `mode == CUBLAS_SIDE_RIGHT` and of size m if `mode == CUBLAS_SIDE_LEFT`. X is gathered from one-dimensional array `x` with stride `incx`. The

absolute value of `incx` is the stride and the sign of `incx` is direction of the stride. If `incx` is positive, then we forward `x` from the first element. Otherwise, we backward `x` from the last element. The formula of `X` is

$$X[j] = \begin{cases} x[j \times incx] & \text{if } incx \geq 0 \\ x[(\chi - 1) \times |incx| - j \times |incx|] & \text{if } incx < 0 \end{cases}$$

where $\chi = m$ if `mode == CUBLAS_SIDE_LEFT` and $\chi = n$ if `mode == CUBLAS_SIDE_RIGHT`.

Example 1: if the user wants to perform $diag(diag(B)) \times A$, then `incx = ldb + 1` where `ldb` is leading dimension of matrix `B`, either row-major or column-major.

Example 2: if the user wants to perform $\alpha \times A$, then there are two choices, either `cublasgemv` with `*beta=0` and `transa == CUBLAS_OP_N` or `cublasdggmm` with `incx=0` and `x[0]= α` .

The operation is out-of-place. The in-place only works if `lda = ldc`.

Param.	Memory	In/out	Meaning
handle		input	handle to the CUBLAS library context.
mode		input	left multiply if <code>mode == CUBLAS_SIDE_LEFT</code> or right multiply if <code>mode == CUBLAS_SIDE_RIGHT</code>
m		input	number of rows of matrix <code>A</code> and <code>C</code> .
n		input	number of columns of matrix <code>A</code> and <code>C</code> .
A	device	input	<type> array of dimensions <code>lda</code> \times <code>n</code> with <code>lda</code> \geq <code>max(1,m)</code>
lda		input	leading dimension of two-dimensional array used to store the matrix <code>A</code> .
x	device	input	one-dimensional <type> array of size <code> incx </code> \times <code>m</code> if <code>mode == CUBLAS_SIDE_LEFT</code> and <code> incx </code> \times <code>n</code> if <code>mode == CUBLAS_SIDE_RIGHT</code>
incx		input	stride of one-dimensional array <code>x</code> .
C	device	in/out	<type> array of dimensions <code>ldc</code> \times <code>n</code> with <code>ldc</code> \geq <code>max(1,m)</code> .
ldc		input	leading dimension of a two-dimensional array used to store the matrix <code>C</code> .

The possible error values returned by this function and their meanings are listed below.

Error Values	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	the operation completed successfully
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	the library was not initialized
<code>CUBLAS_STATUS_INVALID_VALUE</code>	the parameters <code>m,n</code> <0 or <code>mode != CUBLAS_SIDE_LEFT, CUBLAS_SIDE_RIGHT</code>
<code>CUBLAS_STATUS_ARCH_MISMATCH</code>	the device does not support double-precision
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU

Chapter 9

Appendix A: Using the CUBLAS Legacy API

This appendix does not provide a full reference of each Legacy API datatype and entry point. Instead, it describes how to use the API, especially where this is different from the regular CUBLAS API.

Note that in this section, all references to the “CUBLAS Library” refer to the Legacy CUBLAS API only.

9.1 Error Status

The `cublasStatus` type is used for function status returns. The CUBLAS Library helper functions return status directly, while the status of core functions can be retrieved using `cublasGetError()`. Notice that reading the error status via `cublasGetError()`, resets the internal error state to `CUBLAS_STATUS_SUCCESS`. Currently, the following values for `cublasStatus` are defined:

Value	Meaning
<code>CUBLAS_STATUS_SUCCESS</code>	the operation completed successfully
<code>CUBLAS_STATUS_NOT_INITIALIZED</code>	the library was not initialized
<code>CUBLAS_STATUS_ALLOC_FAILED</code>	the resource allocation failed
<code>CUBLAS_STATUS_INVALID_VALUE</code>	an invalid numerical value was used as an argument
<code>CUBLAS_STATUS_ARCH_MISMATCH</code>	an absent device architectural feature is required
<code>CUBLAS_STATUS_MAPPING_ERROR</code>	an access to GPU memory space failed
<code>CUBLAS_STATUS_EXECUTION_FAILED</code>	the GPU program failed to execute
<code>CUBLAS_STATUS_INTERNAL_ERROR</code>	an internal operation failed

This legacy type corresponds to type `cublasStatus_t` in the CUBLAS library API.

9.2 Initialization and Shutdown

The functions `cublasInit()` and `cublasShutdown()` are used to initialize and shutdown the CUBLAS library. It is recommended for `cublasInit()` to be called before any other

function is invoked. It allocates hardware resources on the GPU device that is currently bound to the host thread from which it was invoked.

The legacy initialization and shutdown functions are similar to the CUBLAS library API routines `cublasCreate()` and `cublasDestroy()`.

9.3 Thread Safety

The legacy API is not thread safe when used with multiple host threads and devices. It is recommended to be used only when utmost compatibility with Fortran is required and when a single host thread is used to setup the library and make all the functions calls.

9.4 Memory Management

The memory used by the legacy CUBLAS library API is allocated and released using functions `cublasAlloc()` and `cublasFree()`, respectively. These functions create and destroy an object in the GPU memory space capable of holding an array of `n` elements, where each element requires `elemSize` bytes of storage. Please see the legacy CUBLAS API header file “cublas.h” for the prototypes of these functions.

The function `cublasAlloc()` is a wrapper around the function `cudaMalloc()`, therefore device pointers returned by `cublasAlloc()` can be passed to any CUDA™ device kernel functions. However, these device pointers can not be dereferenced in the host code. The function `cublasFree()` is a wrapper around the function `cudaFree()`.

9.5 Scalar Parameters

In the legacy CUBLAS API, scalar parameters are passed by value from the host. Also, the few functions that do return a scalar result, such as `dot()` and `nrm2()`, return the resulting value on the host, and hence these routines will wait for kernel execution on the device to complete before returning, which makes parallelism with streams impractical. However, the majority of functions do not return any value, in order to be more compatible with Fortran and the existing BLAS libraries.

9.6 Helper Functions

In this section we list the helper functions provided by the legacy CUBLAS API and their functionality. For the exact prototypes of these functions please refer to the legacy CUBLAS API header file “cublas.h”.

Helper function	Meaning
<code>cublasInit()</code>	initialize the library
<code>cublasShutdown()</code>	shuts down the library
<code>cublasGetError()</code>	retrieves the error status of the library
<code>cublasSetKernelStream()</code>	sets the stream to be used by the library
<code>cublasAlloc()</code>	allocates the device memory for the library
<code>cublasFree()</code>	releases the device memory allocated for the library
<code>cublasSetVector()</code>	copies a vector x on the host to a vector y on the GPU
<code>cublasGetVector()</code>	copies a vector x on the GPU to a vector y on the host
<code>cublasSetMatrix()</code>	copies a $m \times n$ tile from a matrix on the host to the GPU
<code>cublasGetMatrix()</code>	copies a $m \times n$ tile from a matrix on the GPU to the host
<code>cublasSetVectorAsync()</code>	similar to <code>cublasSetVector()</code> , but the copy is asynchronous
<code>cublasGetVectorAsync()</code>	similar to <code>cublasGetVector()</code> , but the copy is asynchronous
<code>cublasSetMatrixAsync()</code>	similar to <code>cublasSetMatrix()</code> , but the copy is asynchronous
<code>cublasGetMatrixAsync()</code>	similar to <code>cublasGetMatrix()</code> , but the copy is asynchronous

9.7 Level-1,2,3 Functions

The Level-1,2,3 CUBLAS functions (also called core functions) have the same name and behavior as the ones listed in the chapters 3, 4 and 5 in this document. Please refer to the legacy CUBLAS API header file “`cublas.h`” for their exact prototype. Also, the next section talks a bit more about the differences between the legacy and the CUBLAS API prototypes, more specifically how to convert the function calls from one API to another.

9.8 Converting Legacy to the CUBLAS API

There are a few general rules that can be used to convert from legacy to the CUBLAS API.

1. Exchange the header file “`cublas.h`” for “`cublas_v2.h`”.
2. Exchange the type `cublasStatus` for `cublasStatus_t`.
3. Exchange the function `cublasSetKernelStream()` for `cublasSetStream()`.
4. Exchange the function `cublasAlloc()` and `cublasFree()` for `cudaMalloc()` and `cudaFree()`, respectively. Notice that `cudaMalloc()` expects the size of the allocated memory to be provided in bytes (usually simply provide `n × elemSize` to allocate `n` elements, each of size `elemSize` bytes).
5. Declare the `cublasHandle_t` CUBLAS library handle.
6. Initialize the handle using `cublasCreate()`. Also, release the handle once finished using `cublasDestroy()`.
7. Add the handle as the first parameter to all the CUBLAS library function calls.

8. Change the scalar parameters to be passed by reference, instead of by value (usually simply adding “&” symbol in C/C++ is enough, because the parameters are passed by reference on the host by *default*). However, note that if the routine is running asynchronously, then the variable holding the scalar parameter cannot be changed until the kernels that the routine dispatches are completed. See the CUDA C Programming Guide for a detailed discussion of how to use streams.
9. Change the parameter characters ‘N’ or ‘n’ (non-transpose operation), ‘T’ or ‘t’ (transpose operation) and ‘C’ or ‘c’ (conjugate transpose operation) to CUBLAS_OP_N, CUBLAS_OP_T and CUBLAS_OP_C, respectively.
10. Change the parameter characters ‘L’ or ‘l’ (lower part filled) and ‘U’ or ‘u’ (upper part filled) to CUBLAS_FILL_MODE_LOWER and CUBLAS_FILL_MODE_UPPER, respectively.
11. Change the parameter characters ‘N’ or ‘n’ (non-unit diagonal) and ‘U’ or ‘u’ (unit diagonal) to CUBLAS_DIAG_NON_UNIT and CUBLAS_DIAG_UNIT, respectively.
12. Change the parameter characters ‘L’ or ‘l’ (left side) and ‘R’ or ‘r’ (right side) to CUBLAS_SIDE_LEFT and CUBLAS_SIDE_RIGHT, respectively.
13. If the legacy API function returns a scalar value, add an extra scalar parameter of the same type passed by reference, as the last parameter to the same function.
14. Instead of using `cublasGetError()`, use the return value of the function itself to check for errors.

Finally, please use the function prototypes in the header files “cublas.h” and “cublas_v2.h” to check the code for correctness.

9.9 Examples

For sample code references that use the legacy CUBLAS API please see the two examples below. They show an application written in C using the legacy CUBLAS library API with two indexing styles (Example A.1. "Application Using C and CUBLAS: 1-based indexing" and Example A.2. "Application Using C and CUBLAS: 0-based Indexing"). This application is analogous to the one using the CUBLAS library API that is shown in the Introduction chapter.

```
//Example A.1. Application Using C and CUBLAS: 1-based indexing
//-----
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#include "cublas.h"
#define M 6
#define N 5
#define IDX2F(i,j,ld) (((j)-1)*(ld))+((i)-1))
```

```

static __inline__ void modify (float *m, int ldm, int n, int p, int q, float alpha,
                                float beta){
    cublasSscal (n-p+1, alpha, &m[IDX2F(p,q,ldm)], ldm);
    cublasSscal (ldm-p+1, beta, &m[IDX2F(p,q,ldm)], 1);
}

int main (void){
    int i, j;
    cublasStatus stat;
    float* devPtrA;
    float* a = 0;
    a = (float *)malloc (M * N * sizeof (*a));
    if (!a) {
        printf ("host memory allocation failed");
        return EXIT_FAILURE;
    }
    for (j = 1; j <= N; j++) {
        for (i = 1; i <= M; i++) {
            a[IDX2F(i,j,M)] = (float)((i-1) * M + j);
        }
    }
    cublasInit();
    stat = cublasAlloc (M*N, sizeof(*a), (void**)&devPtrA);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("device memory allocation failed");
        cublasShutdown();
        return EXIT_FAILURE;
    }
    stat = cublasSetMatrix (M, N, sizeof(*a), a, M, devPtrA, M);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("data download failed");
        cublasFree (devPtrA);
        cublasShutdown();
        return EXIT_FAILURE;
    }
    modify (devPtrA, M, N, 2, 3, 16.0f, 12.0f);
    stat = cublasGetMatrix (M, N, sizeof(*a), devPtrA, M, a, M);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("data upload failed");
        cublasFree (devPtrA);
        cublasShutdown();
        return EXIT_FAILURE;
    }
    cublasFree (devPtrA);
    cublasShutdown();
    for (j = 1; j <= N; j++) {
        for (i = 1; i <= M; i++) {
            printf ("%7.0f", a[IDX2F(i,j,M)]);
        }
        printf ("\n");
    }
    return EXIT_SUCCESS;
}

```

```

//Example A.2. Application Using C and CUBLAS: 0-based indexing
//-----
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
#include "cublas.h"

```

```

#define M 6
#define N 5
#define IDX2C(i,j,ld) (((j)*(ld))+(i))

static __inline__ void modify (float *m, int ldm, int n, int p, int q, float alpha,
float beta){
    cublasSscal (n-p, alpha, &m[IDX2C(p,q,ldm)], ldm);
    cublasSscal (ldm-p, beta, &m[IDX2C(p,q,ldm)], 1);
}

int main (void){
    int i, j;
    cublasStatus stat;
    float* devPtrA;
    float* a = 0;
    a = (float *)malloc (M * N * sizeof (*a));
    if (!a) {
        printf ("host memory allocation failed");
        return EXIT_FAILURE;
    }
    for (j = 0; j < N; j++) {
        for (i = 0; i < M; i++) {
            a[IDX2C(i,j,M)] = (float)(i * M + j + 1);
        }
    }
    cublasInit();
    stat = cublasAlloc (M*N, sizeof(*a), (void**)&devPtrA);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("device memory allocation failed");
        cublasShutdown();
        return EXIT_FAILURE;
    }
    stat = cublasSetMatrix (M, N, sizeof(*a), a, M, devPtrA, M);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("data download failed");
        cublasFree (devPtrA);
        cublasShutdown();
        return EXIT_FAILURE;
    }
    modify (devPtrA, M, N, 1, 2, 16.0f, 12.0f);
    stat = cublasGetMatrix (M, N, sizeof(*a), devPtrA, M, a, M);
    if (stat != CUBLAS_STATUS_SUCCESS) {
        printf ("data upload failed");
        cublasFree (devPtrA);
        cublasShutdown();
        return EXIT_FAILURE;
    }
    cublasFree (devPtrA);
    cublasShutdown();
    for (j = 0; j < N; j++) {
        for (i = 0; i < M; i++) {
            printf ("%7.0f", a[IDX2C(i,j,M)]);
        }
        printf ("\n");
    }
    return EXIT_SUCCESS;
}

```


Chapter 10

Appendix B: CUBLAS Fortran Bindings

The CUBLAS library is implemented using the C-based CUDA toolchain, and thus provides a C-style API. This makes interfacing to applications written in C and C++ trivial, but the library can also be used by applications written in Fortran. In particular, the CUBLAS library uses 1-based indexing and Fortran-style column-major storage for multidimensional data to simplify interfacing to Fortran applications. Unfortunately, Fortran-to-C calling conventions are not standardized and differ by platform and toolchain. In particular, differences may exist in the following areas:

- ▶ symbol names (capitalization, name decoration)
- ▶ argument passing (by value or reference)
- ▶ passing of string arguments (length information)
- ▶ passing of pointer arguments (size of the pointer)
- ▶ returning floating-point or compound data types (for example single-precision or complex data types)

To provide maximum flexibility in addressing those differences, the CUBLAS Fortran interface is provided in the form of wrapper functions. These wrapper functions, written in C, are provided in two forms:

1. the thunking wrapper interface located in the file `fortran_thunking.c`
2. the direct wrapper interface located in the file `fortran.c`

The code of one of those 2 files needs to be compiled into an application for it to call the CUBLAS API functions. Providing source code allows users to make any changes necessary for a particular platform and toolchain.

The code in those two C files has been used to demonstrate interoperability with the compilers g77 3.2.3 and g95 0.91 on 32-bit Linux, g77 3.4.5 and g95 0.91 on 64-bit Linux, Intel Fortran 9.0 and Intel Fortran 10.0 on 32-bit and 64-bit Microsoft Windows XP, and g77 3.4.0 and g95 0.92 on Mac OS X.

Note that for g77, use of the compiler flag `-fno-second-underscore` is required to use these wrappers as provided. Also, the use of the default calling conventions with regard to

argument and return value passing is expected. Using the flag `-fno-f2c` changes the default calling convention with respect to these two items.

The thunking wrappers allow interfacing to existing Fortran applications without any changes to the application. During each call, the wrappers allocate GPU memory, copy source data from CPU memory space to GPU memory space, call CUBLAS, and finally copy back the results to CPU memory space and deallocate the GPU memory. As this process causes very significant call overhead, these wrappers are intended for light testing, not for production code. To use the thunking wrappers, the application needs to be compiled with the file `fortran_thunking.c`

The direct wrappers, intended for production code, substitute device pointers for vector and matrix arguments in all BLAS functions. To use these interfaces, existing applications need to be modified slightly to allocate and deallocate data structures in GPU memory space (using `CUBLAS_ALLOC` and `CUBLAS_FREE`) and to copy data between GPU and CPU memory spaces (using `CUBLAS_SET_VECTOR`, `CUBLAS_GET_VECTOR`, `CUBLAS_SET_MATRIX`, and `CUBLAS_GET_MATRIX`). The sample wrappers provided in `fortran.c` map device pointers to the OS-dependent type `size_t`, which is 32-bit wide on 32-bit platforms and 64-bit wide on a 64-bit platforms.

!Example B.1. Fortran 77 Application Executing on the Host

```
!
subroutine modify (m, ldm, n, p, q, alpha, beta)
implicit none
integer ldm, n, p, q
real*4 m(ldm,*), alpha, beta
external cublas_sscal
call cublas_sscal (n-p+1, alpha, m(p,q), ldm)
call cublas_sscal (ldm-p+1, beta, m(p,q), 1)
return
end

program matrixmod
implicit none
integer M, N
parameter (M=6, N=5)
real*4 a(M,N)
integer i, j
external cublas_init
external cublas_shutdown

do j = 1, N
  do i = 1, M
    a(i,j) = (i-1) * M + j
  enddo
enddo

call cublas_init
call modify (a, M, N, 2, 3, 16.0, 12.0)
call cublas_shutdown
do j = 1, N
  do i = 1, M
    write(*,"(F7.0$)") a(i,j)
  enddo
  write (*,*) " "
enddo
```

```
stop
end
```

One approach to deal with index arithmetic on device pointers in Fortran code is to use C-style macros, and use the C preprocessor to expand these, as shown in the example below. On Linux and Mac OS X, one way of pre-processing is to use the option '-E -x f77-cpp-input' when using g77 compiler, or simply the option '-cpp' when using g95 or gfortran. On Windows platforms with Microsoft Visual C/C++, using 'cl -EP' achieves similar results.

When traditional fixed-form Fortran 77 code is ported to use the CUBLAS library, line length often increases when the BLAS calls are exchanged for CUBLAS calls. Longer function names and possible macro expansion are contributing factors. Inadvertently exceeding the maximum line length can lead to run-time errors that are difficult to find, so care should be taken not to exceed the 72-column limit if fixed form is retained.

The examples in this chapter show a small application implemented in Fortran 77 on the host and the same application with the non-thunking wrappers after it has been ported to use the CUBLAS library.

The second example should be compiled with ARCH_64 defined as 1 on 64-bit OS system and as 0 on 32-bit OS system. For example for g95 or gfortran, this can be done directly on the command line by using the option '-cpp -DARCH_64=1'.

```
!Example B.2. Same Application Using Non-thunking CUBLAS Calls
!-----
#define IDX2F(i,j,ld) (((j)-1)*(ld))+((i)-1))

      subroutine modify (devPtrM, ldm, n, p, q, alpha, beta)
      implicit none
      integer sizeof_real
      parameter (sizeof_real=4)
      integer ldm, n, p, q
#if ARCH_64
      integer*8 devPtrM
#else
      integer*4 devPtrM
#endif
      real*4 alpha, beta
      call cublas_sscal (n-p+1, alpha,
1      devPtrM+IDX2F(p,q,ldm)*sizeof_real,
2      ldm)
      call cublas_sscal (ldm-p+1, beta,
1      devPtrM+IDX2F(p,q,ldm)*sizeof_real,
2      1)
      return
      end

      program matrixmod
      implicit none
      integer M, N, sizeof_real
#if ARCH_64
      integer*8 devPtrA
#else
```

```

integer*4 devPtrA
#endif
parameter (M=6, N=5, sizeof_real=4)
real*4 a(M,N)
integer i, j, stat
external cublas_init, cublas_set_matrix, cublas_get_matrix
external cublas_shutdown, cublas_alloc
integer cublas_alloc, cublas_set_matrix, cublas_get_matrix
do j = 1, N
  do i = 1, M
    a(i,j) = (i-1) * M + j
  enddo
enddo
call cublas_init
stat = cublas_alloc(M*N, sizeof_real, devPtrA)
if (stat .NE. 0) then
  write(*,*) "device memory allocation failed"
  call cublas_shutdown
  stop
endif
stat = cublas_set_matrix (M, N, sizeof_real, a, M, devPtrA, M)
if (stat .NE. 0) then
  call cublas_free (devPtrA)
  write(*,*) "data download failed"
  call cublas_shutdown
  stop
endif
call modify (devPtrA, M, N, 2, 3, 16.0, 12.0)
stat = cublas_get_matrix (M, N, sizeof_real, devPtrA, M, a, M)
if (stat .NE. 0) then
  call cublas_free (devPtrA)
  write(*,*) "data upload failed"
  call cublas_shutdown
  stop
endif
call cublas_free (devPtrA)
call cublas_shutdown
do j = 1, N
  do i = 1, M
    write(*,"(F7.0$)") a(i,j)
  enddo
  write (*,*) " "
enddo
stop
end

```

Notice

ALL NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, "MATERIALS") ARE BEING PROVIDED "AS IS." NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE.

Information furnished is believed to be accurate and reliable. However, NVIDIA Corporation assumes no responsibility for the consequences of use of such information or for any infringement of patents or other rights of third parties that may result from its use. No license is granted by implication of otherwise under any patent rights of NVIDIA Corporation. Specifications mentioned in this publication are subject to change without notice. This publication supersedes and replaces all other information previously supplied. NVIDIA Corporation products are not authorized as critical components in life support devices or systems without express written approval of NVIDIA Corporation.

Trademarks

NVIDIA and the NVIDIA logo are trademarks and/or registered trademarks of NVIDIA Corporation in the U.S. and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

Copyright

© 2012 NVIDIA Corporation. All rights reserved.