

CUDA Libraries

Outline



CUDA includes 2 widely used libraries:

- CUBLAS: BLAS implementation
- CUFFT: FFT implementation

CUBLAS



CUBLAS is an implementation of BLAS (Basic Linear Algebra Subprograms) on top of the CUDA driver. It allows access to the computational resources of NVIDIA GPUs.

The library is self-contained at the API level, that is, no direct interaction with the CUDA driver is necessary.

The basic model by which applications use the CUBLAS library is to:
•create matrix and vector objects in GPU memory space,
•fill them with data,
•call a sequence of CUBLAS functions,
•upload the results from GPU memory space back to the host.

CUBLAS provides helper functions for creating and destroying objects in GPU space, and for writing data to and retrieving data from these objects.

Supported features



BLAS functions implemented (single precision only):
 Real data: Level 1, 2 and 3
 Complex data: Level1 and CGEMM

(Level 1=vector vector O(N), Level 2= matrix vector $O(N^2)$, Level 3=matrix matrix $O(N^3)$)

 For maximum compatibility with existing Fortran environments, CUBLAS uses column-major storage, and 1-based indexing:

Since C and C++ use row-major storage, this means applications cannot use the native C array semantics for two-dimensional arrays. Instead, macros or inline functions should be defined to implement matrices on top of onedimensional arrays.

Using CUBLAS



•The interface to the CUBLAS library is the header file cublas.h

•Function names: cublas(Original name). cublasSGEMM

•Because the CUBLAS core functions (as opposed to the helper functions) do not return error status directly, CUBLAS provides a separate function to retrieve the last error that was recorded, to aid in debugging

•CUBLAS is implemented using the C-based CUDA tool chain, and thus provides a C-style API. This makes interfacing to applications written in C or C++ trivial.

cublasInit, cublasShutdown



cublasStatus cublasInit()

initializes the CUBLAS library and must be called before any other CUBLAS API function is invoked. It allocates hardware resources necessary for accessing the GPU.

cublasStatus cublasShutdown()

releases CPU-side resources used by the CUBLAS library. The release of GPU-side resources may be deferred until the application shuts down.

CUBLAS performance



SGEMM performance



Ν

cublasGetError, cublasAlloc, cublasFree 💿

cublasStatus cublasGetError()

returns the last error that occurred on invocation of any of the CUBLAS core functions. While the CUBLAS helper functions return status directly, the CUBLAS core functions do not, improving compatibility with those existing environments that do not expect BLAS functions to return status. Reading the error status via cublasGetError() resets the internal error state to CUBLAS_STATUS_SUCCESS..

cublasStatus cublasAlloc (int n, int elemSize, void **devicePtr)

creates an object in GPU memory space capable of holding an array of n elements, where each element requires elemSize bytes of storage. Note that this is a device pointer that cannot be dereferenced in host code. cublasAlloc() is a wrapper around cudaMalloc(). Device pointers returned by cublasAlloc() can therefore be passed to any CUDA device kernels, not just CUBLAS functions.

cublasStatus cublasFree(const void *devicePtr)

destroys the object in GPU memory space referenced by devicePtr.

cublasSetVector, cublasGetVector



cublasStatus cublasSetVector(int n, int elemSize, const void *x, int incx, void *y, int incy)

copies n elements from a vector x in CPU memory space to a vector y in GPU memory space. Elements in both vectors are assumed to have a size of elemSize bytes. Storage spacing between consecutive elements is incx for the source vector x and incy for the destination vector y

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

copies n elements from a vector x in GPU memory space to a vector y in CPU memory space. Elements in both vectors are assumed to have a size of elemSize bytes. Storage spacing between consecutive elements is incx for the source vector x and incy for the destination vector y

cublasSetMatrix, cublasGetMatrix



cublasStatus cublasSetMatrix(int rows, int cols, int elemSize,

const void *A, int Ida, void *B, int Idb)

copies a tile of rows x cols elements from a matrix A in CPU memory space to a matrix B in GPU memory space. Each element requires storage of elemSize bytes. Both matrices are assumed to be stored in column-major format, with the leading dimension (that is, the number of rows) of source matrix A provided in Ida, and the leading dimension of destination matrix B provided in Idb.

cublasStatus cublasGetMatrix(int rows, int cols, int elemSize, const void *A, int Ida, void *B, int Idb)

copies a tile of rows x cols elements from a matrix A in GPU memory space to a matrix B in CPU memory space. Each element requires storage of elemSize bytes. Both matrices are assumed to be stored in column-major format, with the leading dimension (that is, the number of rows) of source matrix A provided in Ida, and the leading dimension of destination matrix B provided in Idb.

Calling CUBLAS from FORTRAN



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 type)

•CUBLAS provides wrapper functions (in the file fortran.c) that need to be compiled with the user preferred toolchain. Providing source code allows users to make any changes necessary for a particular platform and toolchain.

Calling CUBLAS from FORTRAN



Two different interfaces:

•**Thunking** (define CUBLAS_USE_THUNKING when compiling fortran.c): 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 GPGPU memory. As this process causes significant call overhead, these wrappers are intended for light testing, not for production code.

•Non-Thunking (default):

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 GPGPU 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).

FORTRAN 77 Code example:



```
program matrixmod
implicit none
integer M, N
parameter (M=6, N=5)
real*4 a(M,N)
integer i, j
do i = 1, N
 do i = 1, M
    a(i,j) = (i-1) * M + j
 enddo
enddo
call modify (a, M, N, 2, 3, 16.0, 12.0)
do_i = 1, N
 <u>do</u> i = 1, M
   write(*,"(F7.0$)") a(i,j)
 enddo
 write (*,*) ""
enddo
stop
end
```

subroutine modify (m, ldm, n, p, q, alpha, beta) implicit none integer ldm, n, p, q real*4 m(ldm,*), alpha, beta

external sscal

call sscal (n-p+1, alpha, m(p,q), ldm)

call sscal (ldm-p+1, beta, m(p,q), 1)

return end

FORTRAN 77 Code example: Non-thunking interface

```
program matrixmod
implicit none
integer M, N, sizeof_real, devPtrA
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
```

```
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"
    stop
endif
```

call cublas_set_matrix (M, N, sizeof_real, a, M, devPtrA, M) call modify (devPtrA, M, N, 2, 3, 16.0, 12.0) call cublas_get_matrix (M, N, sizeof_real, devPtrA, M, a, M) 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

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

If using fixed format check that the line length is below the 72 column limit !!!



CUFFT



The Fast Fourier Transform (FFT) is a divide-andconquer algorithm for efficiently computing discrete Fourier transform of complex or real-valued data sets.

The FFT is one of the most important and widely used numerical algorithms.

CUFFT, the "CUDA" FFT library, provides a simple interface for computing parallel FFT on an NVIDIA GPU. This allows users to leverage the floating-point power and parallelism of the GPU without having to develop a custom, GPU-based FFT implementation.

Supported features



- 1D, 2D and 3D transforms of complex and real-valued data
- Batched execution for doing multiple 1D transforms in parallel
- 1D transform size up to 8M elements
- 2D and 3D transform sizes in the range [2,16384]
- In-place and out-of-place transforms for real and complex data.

CUFFT Types and Definitions



type cufftHandle:

is a handle type used to store and access CUFFT plans

type cufftResults:

is an enumeration of values used as API function values return values.

CUFFT_SUCCESS CUFFT_INVALID_PLAN CUFFT_ALLOC_FAILED CUFFT_INVALID_TYPE CUFFT_INVALID_VALUE CUFFT_INTERNAL_ERROR CUFFT_EXEC_FAILED CUFFT_SETUP_FAILED CUFFT_SHUTDOWN_FAILED CUFFT_INVALID_SIZE

Any CUFFT operation is successful. CUFFT is passed an invalid plan handle. CUFFT failed to allocate GPU memory. The user requests an unsupported type. The user specifies a bad memory pointer. Used for all internal driver errors. CUFFT failed to execute an FFT on the GPU. The CUFFT library failed to initialize. The CUFFT library failed to shut down. The user specifies an unsupported FFT size.

Transform types



The library supports complex and real data transforms: CUFFT_C2C, CUFFT_C2R, CUFFT_R2C with directions: CUFFT_FORWARD (-1) and CUFFT_BACKWARD (1) according to the sign of the complex exponential term

For complex FFTs, the input and output arrays must interleave the real and imaginary part (cufftComplex type is defined for this purpose)

For real-to-complex FFTs, the output array holds only the nonredundant complex coefficients:

N -> N/2+1

N0 x N1 x x Nn -> N0 x N1 x X (Nn/2+1)

To perform in-place transform the input/output needs to be padded

More on transforms



- For 2D and 3D transforms, CUFFT performs transforms in rowmajor (C-order).
- If calling from FORTRAN or MATLAB, remember to change the order of size parameters during plan creation.
- CUFFT performs un-normalized transforms:

IFFT(FFT(A))= length(A)*A

CUFFT API is modeled after FFTW. Based on plans, that completely specify the optimal configuration to execute a particular size of FFT.

Once a plan is created, the library stores whatever state is needed to execute the plan multiple times without recomputing the configuration: it works very well for CUFFT, because different kinds of FFTs require different thread configurations and GPU resources.

cufftPlan1d()



cufftResult cufftPlan1d(cufftHandle *plan, int nx, cufftType type, int batch);

creates a 1D FFT plan configuration for a specified signal size and data type. The batch input parameter tells CUFFT how many 1D transforms to configure.

Input:

plan Pointer to a cufftHandle object
nx The transform size (e.g., 256 for a 256-point FFT)
type The transform data type (e.g., CUFFT_C2C for complex-to-complex)
batch Number of transforms of size nx

Output:

plan Contains a CUFFT 1D plan handle value

cufftPlan2d()



cufftResult cufftPlan2d(cufftHandle *plan, int nx, int ny, cufftType type);

creates a 2D FFT plan configuration for a specified signal size and data type.

Input:

planPointer to a cufftHandle objectnxThe transform size in X dimensionnyThe transform size in Y dimensiontypeThe transform data type (e.g., CUFFT_C2C for complex-to-complex)

Output:

plan Contains a CUFFT 2D plan handle value

cufftPlan3d()



cufftResult cufftPlan3d(cufftHandle *plan, int nx, int ny, int nz, cufftType type);

creates a 3D FFT plan configuration for a specified signal size and data type.

Input:

- plan Pointer to a cufftHandle object
- nx The transform size in X dimension
- ny The transform size in Y dimension
- nz The transform size in Z dimension
- type The transform data type (e.g., CUFFT_C2C for complex-to-complex)

Output:

plan Contains a CUFFT 3D plan handle value





cufftResult cufftDestroy(cufftHandle plan);

frees all GPU resources associated with a CUFFT plan and destroys the internal plan data structure. This function should be called once a plan is no longer needed to avoid wasting GPU memory.

Input:

plan cufftHandle object





cufftResult cufftExecC2C(cufftHandle plan, cufftComplex *idata, cufftComplex *odata, int direction);

executes a CUFFT complex to complex transform plan.CUFFT uses as input data the GPU memory pointed to by the idata parameter. This function stores the Fourier coefficients in the odata array. If idata and odata are the same, this method does an in-place transform.

Input:

plancufftHandle object for the plane to updateidataPointer to the input data (in GPU memory) to transformodataPointer to the output data (in GPU memory)directionThe transform direction (CUFFT_FORWARD or CUFFT_BACKWARD)

Output:

odata

Contains the complex Fourier coefficients)

cufftExecR2C()



cufftResult cufftExecR2C(cufftHandle plan, cufftReal *idata, cufftComplex *odata);

executes a CUFFT real to complex transform plan.CUFFT uses as input data the GPU memory pointed to by the idata parameter. This function stores the Fourier coefficients in the odata array. If idata and odata are the same, this method does an in-place transform.

The output hold only the non-redundant complex Fourier coefficients.

Input:

plan	Pointer to a cufftHandle object
idata	Pointer to the input data (in GPU memory) to transform
odata	Pointer to the output data (in GPU memory)

Output: odata

Contains the complex Fourier coefficients

cufftExecC2R()



cufftResult cufftExecC2R(cufftHandle plan, cufftComplex *idata, cufftReal *odata);

executes a CUFFT complex to real transform plan. CUFFT uses as input data the GPU memory pointed to by the idata parameter. This function stores the Fourier coefficients in the odata array. If idata and odata are the same, this method does an in-place transform. The input hold only the non-redundant complex Fourier coefficients.

Input:

olan	Pointer to a cufftHandle object
data	Pointer to the complex input data (in GPU memory) to transform
odata	Pointer to the real output data (in GPU memory)

Output:

odata Contains the real-valued Fourier coefficients

Accuracy and performance



The CUFFT library implements several FFT algorithms, each with different performances and accuracy.

The best performance paths correspond to transform sizes that:

- 1. Fit in CUDA'a shared memory
- 2. Are powers of a single factor (e.g. power-of-two)

If only condition 1 is satisfied, CUFFT uses a more general mixed-radix factor algorithm that is slower and less accurate numerically.

If none of the above conditions is satisfied, CUFFT uses an out-of-place, mixedradix algorithm that stores all intermediate results in global GPU memory.

One notable exception is for long 1D transforms, where CUFFT uses a distributed algorithm that perform 1D FFT using 2D FFT.

CUFFT does not implement any specialized algorithms for real data, and so there is no direct performance benefit to using real to complex (or complex to real) plans instead of complex to complex. For this release, the real data API exists primarily for convenience

Code example: 1D complex to complex transforms



#define NX 256 #define BATCH 10

cufftHandle plan; cufftComplex *data; cudaMalloc((void**)&data, sizeof(cufftComplex)*NX*BATCH);

/* Create a 1D FFT plan. */ cufftPlan1d(&plan, NX, CUFFT_C2C, BATCH);

/* Use the CUFFT plan to transform the signal in place. cufftExecC2C(plan, data, data, CUFFT_FORWARD);

/* Inverse transform the signal in place. */ cufftExecC2C(plan, data, data, CUFFT_INVERSE);

/* Note:

- (1) Divide by number of elements in data-set to get back original data
- (2) Identical pointers to input and output arrays implies in-place transformation

/* Destroy the CUFFT plan. */ cufftDestroy(plan);

cudaFree(data);

Code example: 2D complex to complex transform



#define NX 256 #define NY 128

cufftHandle plan; cufftComplex *idata, *odata; cudaMalloc((void**)&idata, sizeof(cufftComplex)*NX*NY); cudaMalloc((void**)&odata, sizeof(cufftComplex)*NX*NY);

/* Create a 1D FFT plan. */ cufftPlan2d(&plan, NX,NY, CUFFT_C2C);

/* Use the CUFFT plan to transform the signal out of place. */
cufftExecC2C(plan, idata, odata, CUFFT_FORWARD);

/* Inverse transform the signal in place. */
cufftExecC2C(plan, odata, odata, CUFFT_INVERSE);

/* Note:

Different pointers to input and output arrays implies out of place transformation

/* Destroy the CUFFT plan. */ cufftDestroy(plan);

cudaFree(idata), cudaFree(odata);