The CUBLAS and
CULA libraries

Will Landau

CUBLAS overview

Using CUBLAS

CULA

# The CUBLAS and CULA libraries

Will Landau

Iowa State University

October 28, 2013

# Outline

CUBLAS overview

Using CUBLAS

CULA

# Outline

CUBLAS overview

Using CUBLAS

CULA

The CUBLAS and
CULA libraries

Will Landau

CUBLAS overview

Using CUBLAS

CULA

# CUBLAS

- ▶ **CUBLAS**: CUda Basic Linear Algebra Subroutines, the CUDA C implementation of BLAS.
- ▶ Consider scalars $\alpha, \beta$, vectors $x$, $y$, and matrices $A$, $B$, $C$.
- ▶ 3 "levels of functionality":
    - ▶ Level 1: $y \mapsto \alpha x + y$ and other vector-vector routines.
    - ▶ Level 2: $y \mapsto \alpha A x + \beta y$ and other vector-matrix routines.
    - ▶ Level 3: $C \mapsto \alpha A B + \beta C$ and other matrix-matrix routines.

# Level 1 functions

▶ Let $\alpha$ be a scalar, $x$, $y$, and $m$ be vectors, $G = \begin{bmatrix} c & s \\ -s & c \end{bmatrix}$ be some $2 \times 2$ rotation matrix, and $H$ be an arbitrary $2 \times 2$ matrix,

| in R | float | double |
|------|-------|--------|
| which.max($x$) | cublasIsamax() | cutlasIdamax() |
| which.min($x$) | cublasIsamin() | cublasIdamin() |
| sum(abs($x$)) | cublasSasum() | cublasDasum() |
| $\alpha$ * $x$ + $y$-> $y$ | cublasSaxpy() | cublasDaxpy() |
| $x$ -> $y$ | cublasScopy() | cublasDcopy() |
| sum($x$ * $y$) | cublasSdot() | cublasDdot() |
| sqrt(sum($x^2$)) | cublasSnrm2() | cublasDnrm2() |
| $G$ %*% $x$ | cublasSrot() | cublasDrot() |
| $H$ %*% $x$ | cublasSrotm() | cublasDrotm() |
| $\alpha$ * $x$ -> $x$ | cublasSscal() | cublasDscal() |
| $x$ -> $m$; $y$ -> $x$; $m$ -> $y$ | cublasSswap() | cublasDswap() |

▶ Like everything in CUBLAS, there are also analogous functions for cuComplex and cuDoubleComplex types.

# Example level 2 functions

$$\alpha \mathrm{op}(A) \cdot x + \beta y \mapsto y$$

where

$$\mathrm{op}(A) = \begin{cases} A & \texttt{transa == CUBLAS\_OP\_N} \\ A^T & \texttt{transa == CUBLAS\_OP\_T} \\ A^H & \texttt{transa == CUBLAS\_OP\_C} \end{cases}$$

| type of matrix, $A$ | float | double |
|---|---|---|
| any $m \times n$ | cublasSgemv() | cublasDgemv() |
| banded $m \times n$ | cublasSgbmv() | cublasDgbmv() |
| symmetric, banded | cublasSbmv() | cublasDbmv() |
| symmetric, packed format | cublasSspmv() | cublasDspmv() |
| symmetric, triangular | cublasSsymv() | cublasDsymv() |

# Example level 3 functions

- cublasSgemm() and cublasDgemm(): for any compatible matrices $A$, $B$, and $C$,

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

- cublasSgemmBatched() and cublasDgemmBatched(): for arrays of compatible matrices $A[]$, $B[]$, and $C[]$,

$$\alpha \cdot \text{op}(A[i])\text{op}(B[i]) + \beta C[i] \mapsto C[i]$$

- cublasStrsm() and cublasDtrsm() solve for $X$ when $A$ is triangular:

$$\begin{cases} \text{op}(A)X = \alpha B & \text{trans == CUBLAS\_SIDE\_LEFT} \\ X\text{op}(A) = \alpha B & \text{trans == CUBLAS\_SIDE\_RIGHT} \end{cases}$$

# Outline

CUBLAS overview

## Using CUBLAS

CULA

# Implementation of matrices

The CUBLAS and CULA libraries

Will Landau

CUBLAS overview

Using CUBLAS

CULA

▶ Matrices stored in column major order in linear arrays of memory. Array $A$,

| 1 | 1 | 2 | 3 | 5 | 8 | 13 | 21 | 34 | 55 | 89 | 144 |

encodes matrix $B$,

$$\begin{bmatrix} 1 & 5 & 32 \\ 1 & 8 & 55 \\ 2 & 13 & 89 \\ 3 & 21 & 144 \end{bmatrix}$$

▶ Index by

$$B[\text{row } i, \text{ col } j] = A[j \cdot ld + i]$$

where $ld$ is the lead dimension of the matrix (column length for column major order matrices).

▶ Use a macro for indexing:

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

# CUBLAS context

▶ For CUBLAS version $\geq$ 4.0, you must create a CUBLAS context:

```
1  cublasHandle_t handle;
2  cublasCreate(&handle);
3
4  // your code
5
6  cublasDestroy(handle);
```

▶ Pass handle to every CUBLAS function in your code.
▶ This approach allows the user to use multiple host threads and multiple GPUs.

# CUBLAS helper functions

▶ You don't actually need them, but you might see them:
  ▶ cublasSetVector()
  ▶ cublasGetVector()
  ▶ cublasSetMatrix()
  ▶ cublasGetMatrix()

# Choosing the right header file

► 2 choices of API

  ► cublas_v2.h: API for CUBLAS version 4.0 and above.
  ► cublas.h: older API for programs written with CUBLAS version < 4.0.

► Additions to newer API:

  ► cublasCreate() initializes the handle to the CUBLAS library context.
  ► Scalars can be passed by reference or by value to device functions.
  ► All CUBLAS functions return an error status, cublasStatus_t.
  ► cublasAlloc() and cublasFree() are deprecated. Use cudaMalloc() and cudaFree() instead.
  ► cublasSetKernelStream() was renamed cublasSetStream().

# Compiling with CUBLAS

1. Include either cublas_v2.h or cublas.h in your source
2. Link the CUBLAS library with the -lcublas flag.

▶ Example2.cu:

```
1 > nvcc −lcublas Example2.cu −o Example2.
2 > ./Example2
3        1       7      13      19      25      31
4        2       8      14      20      26      32
5        3    1728     180     252     324     396
6        4     160      16      22      28      34
7        5     176      17      23      29      35
```

# Example2.cu

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <math.h>
4  #include <cuda_runtime.h>
5  #include <cublas_v2.h>
6  #define M 6
7  #define N 5
8  #define IDX2F(i,j,ld) (((j−1)*ld)+(i−1))
9
10 static __inline__ void modify (cublasHandle_t handle, float *m, int
        ldm, int n, int p,
11   int q, float alpha, float beta){
12   cublasSscal (handle, n − p+1, &alpha, &m[IDX2F(p,q,ldm)], ldm);
13   cublasSscal (handle, ldm − p+1, &beta, &m[IDX2F(p,q,ldm)], 1);
14 }
15
16 int main (void){
17   cudaError_t cudaStat;
18   cublasStatus_t stat;
19   cublasHandle_t handle;
20   int i, j;
21   float* devPtrA;
22   float* a = 0;
23   a = (float *)malloc (M * N * sizeof (*a));
24   if (!a) {
25     printf ("host memory allocation failed");
26     return EXIT_FAILURE;
27   }
```

# Example2.cu

```
28    for (j = 1; j <= N; j++) {
29      for (i = 1; i <= M; i++) {
30        a[IDX2F(i,j,M)] = (float)((i-1) * M + j);
31      }
32    }
33
34    cudaStat = cudaMalloc ((void**)&devPtrA, M*N*sizeof(*a));
35    if ( cudaStat != cudaSuccess ) {
36      printf ("device memory allocation failed");
37      return EXIT_FAILURE;
38    }
39
40    stat = cublasCreate(&handle);
41    if ( stat != CUBLAS_STATUS_SUCCESS ) {
42      printf ("CUBLAS initialization failed\n");
43      return EXIT_FAILURE;
44    }
45
46    stat = cublasSetMatrix (M, N, sizeof(*a), a, M, devPtrA, M);
47
48     if(stat != CUBLAS_STATUS_SUCCESS) {
49       printf("data download failed");
50       cudaFree(devPtrA);
51       cublasDestroy(handle);
52       return EXIT_FAILURE;
53    }
```

# Example2.cu

```
55
56
57    modify ( handle , devPtrA , M, N, 2, 3, 16.0f , 12.0f );
58
59    stat = cublasGetMatrix (M, N, sizeof(*a), devPtrA , M, a, M);
60    if ( stat != CUBLAS_STATUS_SUCCESS ) {
61      printf ("data upload failed");
62      cudaFree (devPtrA );
63      cublasDestroy ( handle );
64      return EXIT_FAILURE;
65    }
66
67    cudaFree ( devPtrA );
68    cublasDestroy ( handle );
69
70    for (j = 1; j <= N; j++) {
71      for (i = 1; i <= M; i++) {
72        printf ("%7.0f", a[IDX2F(i,j,M)]);
73      }
74      printf ( "\n" );
75    }
76    return EXIT_SUCCESS;
77  }
```

# Example: `ols.cu`

▶ I will attempt to solve the least squares problem,

$$y = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

by computing the solution,

$$\widehat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T y$$

# Example: `ols.cu`

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <string.h>
4  #include <cuda_runtime.h>
5  #include <cublas_v2.h>
6  #include <cula.h>
7  #include <math.h>
8
9  #define I(i, j, ld) j * ld + i
10
11 #define CUDA_CALL(x) {if((x) != cudaSuccess){ \
12   printf("CUDA error at %s:%d\n",__FILE__,__LINE__); \
13   printf("  %s\n", cudaGetErrorString(cudaGetLastError())); \
14   exit(EXIT_FAILURE);}}
15
16 float rnorm(){
17   float r1 = ((float) rand()) / ((float) RAND_MAX);
18   float r2 = ((float) rand()) / ((float) RAND_MAX);
19   return sqrt( -2 * log(r1) ) * cos(2 * 3.1415 * r2);
20 }
21
22 int main(){
23   int i, j;
24   int n = 10;
25   int p = 3;
26   int* ipiv;
27   float k;
28   float *X, *XtX, *XtY, *beta, *Y, *dX, *dXtX, *dXtY, *dbeta, *dY;
```

# Example: `ols.cu`

```
29    float *a, *b;
30    a = (float*) malloc(sizeof(*X));
31    b = (float*) malloc(sizeof(*X));
32    *a = 1.0;
33    *b = 0.0;
34
35    cublasHandle_t handle;
36    cublasCreate(&handle);
37
38    X = (float*) malloc(n * p * sizeof(*X));
39    XtX = (float*) malloc(p * p * sizeof(*X));
40    XtY = (float*) malloc(p * sizeof(*X));
41    beta = (float*) malloc(p * sizeof(*X));
42    Y = (float*) malloc(n * sizeof(*X));
43
44    CUDA_CALL(cudaMalloc((void**) &ipiv, p * p * sizeof(*ipiv)));
45    CUDA_CALL(cudaMalloc((void**) &dX, n * p * sizeof(*X)));
46    CUDA_CALL(cudaMalloc((void**) &dXtX, p * p * sizeof(*X)));
47    CUDA_CALL(cudaMalloc((void**) &dXtY, p * sizeof(*X)));
48    CUDA_CALL(cudaMalloc((void**) &dbeta, p * sizeof(*X)));
49    CUDA_CALL(cudaMalloc((void**) &dY, n * sizeof(*X)));
```

# Example: `ols.cu`

```
51    printf("Y\t\tX\n");
52    for(i = 0; i < n; i++){
53      k = (float) i;
54      X[I(i, 0, n)] = 1.0;
55      X[I(i, 1, n)] = k / 10.0;
56      X[I(i, 2, n)] = k * k / 10.0;
57      Y[i] = (k - 5.0) * (k - 2.3) / 3.0 + rnorm();
58
59      printf("%0.2f\t\t", Y[i]);
60      for(j = 0; j < p; j++){
61        printf("%0.2f\t", X[I(i, j, n)]);
62      }
63      printf("\n");
64    }
65    printf("\n");
66
67    CUDA_CALL(cudaMemcpy(dX, X, n * p * sizeof(float),
          cudaMemcpyHostToDevice));
68    CUDA_CALL(cudaMemcpy(dY, Y, n * sizeof(float),
          cudaMemcpyHostToDevice));
69
70    cublasSgemm(handle, CUBLAS_OP_T, CUBLAS_OP_N, p, p, n,
71      a, dX, n, dX, n, b, dXtX, p);
72
73    CUDA_CALL(cudaMemcpy(XtX, dXtX, p * p * sizeof(float),
          cudaMemcpyDeviceToHost));
```

# Example: `ols.cu`

```
74     printf("XtX\n");
75     for(i = 0; i < p; i++){
76       for(j = 0; j < p; j++){
77         printf("%0.2f\t", XtX[I(i, j, p)]);
78       }
79       printf("\n");
80     }
81     printf("\n");
```

# Output of code so far

The CUBLAS and
CULA libraries

Will Landau

CUBLAS overview

Using CUBLAS

CULA

```
1 > nvcc -I /usr/local/cula/include -L /usr/local/cula/lib64 -
        lcula_core -lcula_lapack -lcublas -lcudart ols.cu -o ols
2 > ./ols
3 Y X
4 3.37 1.00 0.00 0.00
5 1.94 1.00 0.10 0.10
6 0.44 1.00 0.20 0.40
7 -0.30 1.00 0.30 0.90
8 -2.08 1.00 0.40 1.60
9 -0.84 1.00 0.50 2.50
10 -0.18 1.00 0.60 3.60
11 3.40 1.00 0.70 4.90
12 5.51 1.00 0.80 6.40
13 7.39 1.00 0.90 8.10
14
15 XtX
16 10.00 4.50 28.50
17 4.50 2.85 20.25
18 28.50 20.25 153.33
```

# Example: `ols.cu`

▶ We have $X^T X$, but which we need to invert in order to compute our solution,

$$\widehat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T y$$

▶ But CUBLAS can only invert triangular matrices!

# Enter CULA: CUDA LAPACK

```
82    culaInitialize();
83
84    culaDeviceSgetrf(p, p, dXtX, p, ipiv);
85    culaDeviceSgetri(p, dXtX, p, ipiv);
86
87    CUDA_CALL(cudaMemcpy(XtX, dXtX, p * p * sizeof(float),
              cudaMemcpyDeviceToHost));
88
89    printf("XtX^(-1)\n");
90    for(i = 0; i < p; i++){
91      for(j = 0; j < p; j++){
92        printf("%0.2f\t", XtX[I(i, j, p)]);
93      }
94      printf("\n");
95    }
96    printf("\n");
97
98    cublasSgemm(handle, CUBLAS_OP_T, CUBLAS_OP_N, p, 1, n,
99      a, dX, n, dY, n, b, dXtY, p);
100
101    cublasSgemv(handle, CUBLAS_OP_N, p, p,
102      a, dXtX, p, dXtY, 1, b, dbeta, 1);
103
104    CUDA_CALL(cudaMemcpy(beta, dbeta, p * sizeof(float),
              cudaMemcpyDeviceToHost));
105
106    printf("CUBLAS/CULA matrix algebra parameter estimates:\n");
107    for(i = 0; i < p; i++){
108      printf("beta_%i = %0.2f\n", i, beta[i]);
109    }
110    printf("\n");
```

# CULA's `culaSgels()` does least squares for you

```
111    culaSgels('N', n, p, 1, X, n, Y, n);
112
113    printf("culaSgels Parameter estimates:\n");
114    for(i = 0; i < p; i++){
115      printf("beta_%i = %0.2f\n", i, Y[i]);
116    }
117    printf("\n");
118
119    culaShutdown();
120    cublasDestroy(handle);
121
122    free(a);
123    free(b);
124    free(X);
125    free(XtX);
126    free(XtY);
127    free(beta);
128    free(Y);
129
130    CUDA_CALL(cudaFree(dX));
131    CUDA_CALL(cudaFree(dXtX));
132    CUDA_CALL(cudaFree(dXtY));
133    CUDA_CALL(cudaFree(dbeta));
134    CUDA_CALL(cudaFree(dY));
135  }
```

# Rest of the output

```
19  XtX^(−1)
20  0.62  −2.59  0.23
21  −2.59  16.55  −1.70
22  0.23  −1.70  0.19
23
24  CUBLAS/CULA matrix algebra parameter estimates:
25  beta_0 = 3.78
26  beta_1 = −25.53
27  beta_2 = 3.36
28
29  culaSgels Parameter estimates:
30  beta_0 = 3.78
31  beta_1 = −25.53
32  beta_2 = 3.36
```

# Outline

CUBLAS overview

Using CUBLAS

CULA

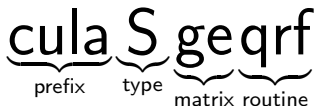The CUBLAS and
CULA libraries

Will Landau

CUBLAS overview
Using CUBLAS
CULA

# More on CULA

- ▶ CULA: the CUDA C implementation of LAPACK
- ▶ Features:
  - ▶ More matrix algebra routines
  - ▶ Factorizations: LU, QR, RQ, QL, SVD, and Cholesky
  - ▶ Solving systems of linear equations (matrix inversion)
  - ▶ Least squares
  - ▶ Eigenvalue solvers
- ▶ Interfaces (collections of functions)
  - ▶ **Standard**: users need not micromanage GPU memory or copy data to or from the GPU.
  - ▶ **Device**: users need explicitly to allocate GPU memory and copy to and from the GPU.
- ▶ Be careful of the standard interface functions: they're convenient, but they copy to and from the GPU with every call.

# CULA naming conventions

$$\underbrace{\text{cula}}_{\text{prefix}}\underbrace{\text{S}}_{\text{type}}\underbrace{\text{ge}}_{\text{matrix}}\underbrace{\text{qrf}}_{\text{routine}}$$

▶ Prefix: `cula` for standard interface, `culaDevice` for device interface:

▶ Type: single precision (S), single precision complex (C), double precision real (D), or double precision complex (Z).

▶ Matrix:

| | |
|----|------------------------------------------|
| bd | Bidiagonal |
| ge | General |
| gg | General matrices, generalized problem |
| he | Hermitian symmetric |
| or | (Real) orthogonal |
| sb | Symmetric, banded |
| sy | Symmetric |
| tr | Triangular |
| un | (Complex) unitary |

# CULA naming conventions

$$\underbrace{\text{cula}}_{\text{prefix}}\underbrace{\text{S}}_{\text{type}}\underbrace{\text{ge}}_{\text{matrix}}\underbrace{\text{qrf}}_{\text{routine}}$$

▶ Routine:

| trf | Triangular factorization |
| sv | Factor a matrix and solve system of linear equations |
| qrf | QR factorization without pivoting |
| svd | Singular value decomposition |
| ls | Solve over- or under-determined linear system |

▶ Consult the CULA manual for other routines.

# Compiling with CULA

The CUBLAS and
CULA libraries

Will Landau

CUBLAS overview

Using CUBLAS

CULA

▶ Include cula_lapack.h for the standard interface,
  cula_lapack_device.h for the device interface, or
  cula.h for both.

▶ Compile on impact1 with:

```
1  nvcc −I / usr / local / cula / include −L / usr /
       local / cula / lib64 −lcula_core −
       lcula_lapack −lcublas −lcudart
       your_source . cu −o your_binary
```

  ▶ -I /usr/local/cula/include tells the compiler,
    nvcc, where to find the header files.
  ▶ -L /usr/local/cula/lib64 tells nvcc where the
    CULA library is (the 64-bit version in this case).
  ▶ -lcula_core -lcula_lapack -lcublas -lcudart
    links the required libraries to your binary.

# Minimal working example: `mwe.cu`

```
1  #include <cula.h>
2  #include <stdlib.h>
3  #include <stdio.h>
4
5  int main(){
6
7    culaStatus s;
8    s = culaInitialize();
9
10   if(s != culaNoError)
11     printf("%s\n", culaGetErrorInfo());
12
13   /* ... Your code ... */
14
15   culaShutdown();
16 }
```

# Minimal working example: `mwe.cu`

```
17 > nvcc −I /usr/local/cula/include −L /usr/local/
        cula/lib64 −lcula_core −lcula_lapack −
        lcublas −lcudart mwe.cu −o mwe
18 > ./mwe
19 >
```

# Outline

CUBLAS overview

Using CUBLAS

CULA

# Resources

The CUBLAS and
CULA libraries

Will Landau

CUBLAS overview

Using CUBLAS

CULA

- ▶ Guides:

  1. CUDA Toolkit 4.2 CUBLAS Library
  2. "CULA Programmers Guide". CULA Tools.
  3. "CULA Reference Manual". CULA Tools.

- ▶ Code from today:

  - ▶ Example2.cu
  - ▶ ols.cu
  - ▶ mwe.cu

- ▶ Other example code:

  - ▶ simpleCUBLAS.cpp
  - ▶ ae.cu
  - ▶ de.cu
  - ▶ deviceInterface.c
  - ▶ ll.cu
  - ▶ systemSolve.c

# That's all for today.

▶ Series materials are available at
  http://will-landau.com/gpu.