

CUDA SDK — LIBRARIES, NUMERICAL ACCURACY

Siegfried Höfinger

VSC Research Center, TU Wien

October 28, 2024

→ <https://tinyurl.com/cudafordummies/ii/14/notes-14.pdf>

CUDA 4 DUMMIES — OCT 29-30, 2024

OUTLINE

TO THE FORTRAN-MINDED

CUBLAS

CUSOLVER

MORE LIBRARIES

NUMERICAL ACCURACY & PERFORMANCE

5_DOMAIN_SPECIFIC/MARCHINGCUBES

PROFILING CUDA CODE

HIP

TAKE HOME MESSAGES

TO THE FORTRAN-MINDED

CUDA SDK

- Fortran is the language of choice in many scientific code development projects

→ <https://developer.nvidia.com/hpc-sdk>

TO THE FORTRAN-MINDED

CUDA SDK

- Fortran is the language of choice in many scientific code development projects
- Professional coverage from the NVIDIA hpc SDK — formerly portland group compilers — now an entire suite of compilers, libraries and developmental tools

→ <https://developer.nvidia.com/hpc-sdk>

TO THE FORTRAN-MINDED

CUDA SDK

- Fortran is the language of choice in many scientific code development projects
- Professional coverage from the NVIDIA hpc SDK — formerly portland group compilers — now an entire suite of compilers, libraries and developmental tools
- NVFORTRAN, NVC++ and NVC compilers (NVCC excluded)

→ <https://developer.nvidia.com/hpc-sdk>

TO THE FORTRAN-MINDED

CUDA SDK

- Fortran is the language of choice in many scientific code development projects
- Professional coverage from the NVIDIA hpc SDK — formerly portland group compilers — now an entire suite of compilers, libraries and developmental tools
- NVFORTRAN, NVC++ and NVC compilers (NVCC excluded)
- Simply download from <https://developer.nvidia.com/nvidia-hpc-sdk-249-downloads>

→ <https://developer.nvidia.com/hpc-sdk>

TO THE FORTRAN-MINDED

CUDA SDK

- Fortran is the language of choice in many scientific code development projects
- Professional coverage from the NVIDIA hpc SDK — formerly portland group compilers — now an entire suite of compilers, libraries and developmental tools
- NVFORTRAN, NVC++ and NVC compilers (NVCC excluded)
- Simply download from <https://developer.nvidia.com/nvidia-hpc-sdk-249-downloads>
- Follow guidelines in <https://docs.nvidia.com/hpc-sdk/archive/24.9/hpc-sdk-install-guide/index.html>

→ <https://developer.nvidia.com/hpc-sdk>

TO THE FORTRAN-MINDED

CUDA SDK

- Fortran is the language of choice in many scientific code development projects
- Professional coverage from the NVIDIA hpc SDK — formerly portland group compilers — now an entire suite of compilers, libraries and developmental tools
- NVFORTRAN, NVC++ and NVC compilers (NVCC excluded)
- Simply download from <https://developer.nvidia.com/nvidia-hpc-sdk-249-downloads>
- Follow guidelines in <https://docs.nvidia.com/hpc-sdk/archive/24.9/hpc-sdk-install-guide/index.html>
- Install and load the module to get ready

→ <https://developer.nvidia.com/hpc-sdk>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

vecadd.cuf

```
! host
program t1
use cudafor
use myvecadd
  integer, parameter :: n = 100
  integer :: i
  real, allocatable, device :: da(:), db(:), dc(:)
  real :: ha(n), hb(n), hc(n)
  istat = cudaSetDevice(0)
  allocate(da(n))
  ...
  da = ha
  ...
  call vecadd<<<< 1, n >>>>(da, db, dc)
  ha = da
  ...
  deallocate(da)
  ...
end program t1
```

```
! kernel
module myvecadd
contains
attributes(global) &
subroutine vecadd( a, b, c )
  integer :: i
  real :: a(*), b(*), c(*)

  i = (blockidx%x-1) * blockdim%x + threadidx%x
  c(i) = a(i) + b(i)
end subroutine vecadd
end module myvecadd
```

→ <https://tinyurl.com/cudafordummies/ii/14/vecadd.cuf>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

vecadd.cuf

```
! host
program t1
use cudafor
use myvecadd
  integer, parameter :: n = 100
  integer :: i
  real, allocatable, device :: da(:), db(:), dc(:)
  real :: ha(n), hb(n), hc(n)
  istat = cudaSetDevice(0)
  allocate(da(n))
  ...
  da = ha
  ...
  call vecadd<<<< 1, n >>>>(da, db, dc)
  ha = da
  ...
  deallocate(da)
  ...
end program t1
```

Select GPU 0

```
! kernel
module myvecadd
contains
attributes(global) &
subroutine vecadd( a, b, c )
  integer :: i
  real :: a(*), b(*), c(*)

  i = (blockidx%x-1) * blockdim%x + threadidx%x
  c(i) = a(i) + b(i)
end subroutine vecadd
end module myvecadd
```

→ <https://tinyurl.com/cudafordummies/ii/14/vecadd.cuf>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

vecadd.cuf

```
! host
program t1
use cudafor
use myvecadd
integer, parameter :: n = 100
integer :: i
real, allocatable, device :: da(:), db(:), dc(:)
real :: ha(n), hb(n), hc(n)
istat = cudaSetDevice(0)
allocate(da(n))
...
da = ha
...
call vecadd<<< 1, n >>>(da, db, dc)
ha = da
...
deallocate(da)
...
end program t1
```

Select GPU 0

Simple
htod/dtoh
copies

```
! kernel
module myvecadd
contains
attributes(global) &
subroutine vecadd( a, b, c )
integer :: i
real :: a(*), b(*), c(*)

i = (blockidx%x-1) * blockdim%x + threadidx%x
c(i) = a(i) + b(i)
end subroutine vecadd
end module myvecadd
```

→ <https://tinyurl.com/cudafordummies/ii/14/vecadd.cuf>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

vecadd.cuf

```
! host
program t1
use cudafor
use myvecadd
integer, parameter :: n = 100
integer :: i
real, allocatable, device :: da(:), db(:), dc(:)
real :: ha(n), hb(n), hc(n)
istat = cudaSetDevice(0)
allocate(da(n))
...
da = ha
...
call vecadd<<< 1, n >>>(da, db, dc)
ha = da
...
deallocate(da)
...
end program t1
```

Select GPU 0

Simple
htod/dtoh
copies

Slightly different syntax

```
! kernel
module myvecadd
contains
attributes(global) &
subroutine vecadd( a, b, c )
integer :: i
real :: a(*), b(*), c(*)

i = (blockidx%x-1) * blockdim%x + threadidx%x
c(i) = a(i) + b(i)
end subroutine vecadd
end module myvecadd
```

→ <https://tinyurl.com/cudafordummies/ii/14/vecadd.cuf>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

vecadd.cuf

```
! host
program t1
use cudafor
use myvecadd
integer, parameter :: n = 100
integer :: i
real, allocatable, device :: da(:), db(:), dc(:)
real :: ha(n), hb(n), hc(n)
istat = cudaSetDevice(0)
allocate(da(n))
...
da = ha
...
call vecadd<<< 1, n >>>(da, db, dc)
ha = da
...
deallocate(da)
...
end program t1
```

Select GPU 0

Simple
htod/dtoh
copies

```
! kernel
module myvecadd
contains
attributes(global) &
subroutine vecadd( a, b, c )
integer :: i
real :: a(*), b(*), c(*)

i = (blockidx%x-1) * blockdim%x + threadidx%x
c(i) = a(i) + b(i)
end subroutine vecadd
end module myvecadd
```

Slightly different syntax

Fortran counting
1,2,3...

→ <https://tinyurl.com/cudafordummies/ii/14/vecadd.cuf>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

```
cuda-zen sh@n3073-009:~$ nvfortran vecadd.cuf
cuda-zen sh@n3073-009:~$ ./a.out

  1  1.000000    99.00000    100.0000
  2  2.000000    98.00000    100.0000
  3  3.000000    97.00000    100.0000
  4  4.000000    96.00000    100.0000
  5  5.000000    95.00000    100.0000
.....
 98 98.00000     2.00000    100.0000
 99 99.00000     1.00000    100.0000
100 100.0000     0.00000    100.0000
```

→ <https://tinyurl.com/cudafordummies/ii/14/vecadd.cuf>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

- Frequently — only a single Fortran subroutine/function needs to be ported

→ <https://www.pgroup.com/resources/docs/18.4/pdf/pgi18cudaforg.pdf>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

- Frequently — only a single Fortran subroutine/function needs to be ported
- Basic feasibility studies can already be done with the help of short C routines called from Fortran

→ <https://www.pgroup.com/resources/docs/18.4/pdf/pgi18cudaforg.pdf>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

- Frequently — only a single Fortran subroutine/function needs to be ported
- Basic feasibility studies can already be done with the help of short C routines called from Fortran
- The usual Fortran \leftrightarrow C interfacing rules apply

→ <https://www.pgroup.com/resources/docs/18.4/pdf/pgi18cudaforg.pdf>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

- Frequently — only a single Fortran subroutine/function needs to be ported
- Basic feasibility studies can already be done with the help of short C routines called from Fortran
- The usual Fortran \leftrightarrow C interfacing rules apply
- Memory management is more complicated, would prefer `cudaMallocManaged()`

→ <https://www.pgroup.com/resources/docs/18.4/pdf/pgi18cudaforug.pdf>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

fvcadd.f

```
PROGRAM FVCADD
IMPLICIT NONE
INTEGER N
PARAMETER (N = 100)
INTEGER I
REAL A(N), B(N), C(N)

DO I=1, N
  A(I) = REAL(I)
  B(I) = REAL(N) - REAL(I)
  C(I) = REAL(0)
ENDDO

CALL NTMDTR(A, B, C, N)

DO I=1, N
  WRITE(6, '(I6F12.6)') I, C(I)
ENDDO

END
```

ntmdtr.cu

```
__global__ void VecAdd(float *A, float *B, float *C)
{
  int i = threadIdx.x;
  C[i] = A[i] + B[i];
}

extern "C" void ntmdtr_(float *A, float *B, float *C, int *N)
{
  dim3 numBlocks, threadsPerBlock;
  float *AD, *BD, *CD;
  threadsPerBlock.x = *N;
  numBlocks.x = 1;
  cudaMalloc((void **) &AD, (*N) * sizeof(float));
  cudaMalloc((void **) &BD, (*N) * sizeof(float));
  cudaMalloc((void **) &CD, (*N) * sizeof(float));
  cudaMemcpy(AD, A, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(BD, B, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(CD, C, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  VecAdd <<<< numBlocks, threadsPerBlock >>>> (AD, BD, CD);
  cudaDeviceSynchronize();
  cudaMemcpy(C, CD, (*N) * sizeof(float), cudaMemcpyDeviceToHost);
  cudaFree(AD);
  cudaFree(BD);
  cudaFree(CD);
  return;
}
```

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

fvcadd.f

```
PROGRAM FVCADD
IMPLICIT NONE
INTEGER N
PARAMETER (N = 100)
INTEGER I
REAL A(N), B(N), C(N)

DO I=1, N
  A(I) = REAL(I)
  B(I) = REAL(N) - REAL(I)
  C(I) = REAL(0)
ENDDO

CALL NTMDTR(A, B, C, N)

DO I=1, N
  WRITE(*, '(I6F12.6)') I, C(I)
ENDDO

END
```

ntmdtr.cu

```
__global__ void VecAdd(float *A, float *B, float *C)
{
  int i = threadIdx.x;
  C[i] = A[i] + B[i];
}

extern "C" void ntmdtr_(float *A, float *B, float *C, int *N)
{
  dim3 numBlocks, threadsPerBlock;
  float *AD, *BD, *CD;
  threadsPerBlock.x = *N;
  numBlocks.x = 1;
  cudaMalloc((void **) &AD, (*N) * sizeof(float));
  cudaMalloc((void **) &BD, (*N) * sizeof(float));
  cudaMalloc((void **) &CD, (*N) * sizeof(float));
  cudaMemcpy(AD, A, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(BD, B, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(CD, C, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  VecAdd <<<< numBlocks, threadsPerBlock >>>> (AD, BD, CD);
  cudaDeviceSynchronize();
  cudaMemcpy(C, CD, (*N) * sizeof(float), cudaMemcpyDeviceToHost);
  cudaFree(AD);
  cudaFree(BD);
  cudaFree(CD);
  return;
}
```

Intermediator in C

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

fvcadd.f

```
PROGRAM FVCADD
IMPLICIT NONE
INTEGER N
PARAMETER (N = 100)
INTEGER I
REAL A(N), B(N), C(N)

DO I=1, N
  A(I) = REAL(I)
  B(I) = REAL(N) - REAL(I)
  C(I) = REAL(0)
ENDDO

CALL NTMDTR(A, B, C, N)

DO I=1, N
  WRITE(*, '(I6F12.6)') I, C(I)
ENDDO

END
```

ntmdtr.cu

```
__global__ void VecAdd(float *A, float *B, float *C)
{
  int i = threadIdx.x;
  C[i] = A[i] + B[i];
}

extern "C" void ntmdtr_(float *A, float *B, float *C, int *N)
{
  dim3 numBlocks, threadsPerBlock;
  float *AD, *BD, *CD;
  threadsPerBlock.x = *N;
  numBlocks.x = 1;
  cudaMalloc((void **) &AD, (*N) * sizeof(float));
  cudaMalloc((void **) &BD, (*N) * sizeof(float));
  cudaMalloc((void **) &CD, (*N) * sizeof(float));
  cudaMemcpy(AD, A, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(BD, B, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(CD, C, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  VecAdd <<<< numBlocks, threadsPerBlock >>>> (AD, BD, CD);
  cudaDeviceSynchronize();
  cudaMemcpy(C, CD, (*N) * sizeof(float), cudaMemcpyDeviceToHost);
  cudaFree(AD);
  cudaFree(BD);
  cudaFree(CD);
  return;
}
```

Trailing underscore !

Intermediator in C

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

fvcadd.f

```
PROGRAM FVCADD
IMPLICIT NONE
INTEGER N
PARAMETER (N = 100)
INTEGER I
REAL A(N), B(N), C(N)

DO I=1, N
  A(I) = REAL(I)
  B(I) = REAL(N) - REAL(I)
  C(I) = REAL(0)
ENDDO

CALL NTMDTR(A, B, C, N)

DO I=1, N
  WRITE(*, '(I6F12.6)') I, C(I)
ENDDO

END
```

ntmdtr.cu

```
__global__ void VecAdd(float *A, float *B, float *C)
{
  int i = threadIdx.x;
  C[i] = A[i] + B[i];
}

extern "C" void ntmdtr_(float *A, float *B, float *C, int *N)
{
  dim3 numBlocks, threadsPerBlock;
  float *AD, *BD, *CD;
  threadsPerBlock.x = *N;
  numBlocks.x = 1;
  cudaMalloc((void **) &AD, (*N) * sizeof(float));
  cudaMalloc((void **) &BD, (*N) * sizeof(float));
  cudaMalloc((void **) &CD, (*N) * sizeof(float));
  cudaMemcpy(AD, A, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(BD, B, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(CD, C, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  VecAdd <<<< numBlocks, threadsPerBlock >>>> (AD, BD, CD);
  cudaDeviceSynchronize();
  cudaMemcpy(C, CD, (*N) * sizeof(float), cudaMemcpyDeviceToHost);
  cudaFree(AD);
  cudaFree(BD);
  cudaFree(CD);
  return;
}
```

Trailing underscore !

Call by Reference

Intermediator in C

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

fvcadd.f

```
PROGRAM FVCADD
IMPLICIT NONE
INTEGER N
PARAMETER (N = 100)
INTEGER I
REAL A(N), B(N), C(N)

DO I=1, N
  A(I) = REAL(I)
  B(I) = REAL(N) - REAL(I)
  C(I) = REAL(0)
ENDDO

CALL NTMDTR(A, B, C, N)

DO I=1, N
  WRITE(*, '(I6F12.6)') I, C(I)
ENDDO

END
```

ntmdtr.cu

```
__global__ void VecAdd(float *A, float *B, float *C)
{
  int i = threadIdx.x;
  C[i] = A[i] + B[i];
}

extern "C" void ntmdtr_(float *A, float *B, float *C, int *N)
{
  dim3 numBlocks, threadsPerBlock;
  float *AD, *BD, *CD;
  threadsPerBlock.x = *N;
  numBlocks.x = 1;
  cudaMalloc((void **) &AD, (*N) * sizeof(float));
  cudaMalloc((void **) &BD, (*N) * sizeof(float));
  cudaMalloc((void **) &CD, (*N) * sizeof(float));
  cudaMemcpy(AD, A, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(BD, B, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  cudaMemcpy(CD, C, (*N) * sizeof(float), cudaMemcpyHostToDevice);
  VecAdd <<<< numBlocks, threadsPerBlock >>>> (AD, BD, CD);
  cudaDeviceSynchronize();
  cudaMemcpy(C, CD, (*N) * sizeof(float), cudaMemcpyDeviceToHost);
  cudaFree(AD);
  cudaFree(BD);
  cudaFree(CD);
  return;
}
```

Trailing underscore !

Call by Reference

Intermediator in C

Extra data transfer ↔

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

```
cuda-zen sh@n3073-009:~$ ls
fvcadd.f ntmctr.cu
cuda-zen sh@n3073-009:~$ gfortran -c fvcadd.f
cuda-zen sh@n3073-009:~$ nvcc -c ntmctr.cu
cuda-zen sh@n3073-009:~$ nvcc fvcadd.o ntmctr.o -lcudart -lgfortran
cuda-zen sh@n3073-009:~$ ./a.out
 1 100.000000
 2 100.000000
 3 100.000000
 4 100.000000
 5 100.000000
 6 100.000000
 7 100.000000
...
98 100.000000
99 100.000000
100 100.000000
```

→ <https://www.olcf.ornl.gov/tutorials/compiling-mixed-gpu-and-cpu-code>

→ <https://tinyurl.com/cudafordummies/ii/14/fvcadd.f>

→ <https://tinyurl.com/cudafordummies/ii/14/ntmctr.cu>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

```
cuda-zen sh@n3073-009:~$ ls
fvcadd.f ntmctr.cu
cuda-zen sh@n3073-009:~$ gfortran -c fvcadd.f
cuda-zen sh@n3073-009:~$ nvcc -c ntmctr.cu
cuda-zen sh@n3073-009:~$ nvcc fvcadd.o ntmctr.o -lcudart -lgfortran
cuda-zen sh@n3073-009:~$ ./a.out
 1 100.000000
 2 100.000000
 3 100.000000
 4 100.000000
 5 100.000000
 6 100.000000
 7 100.000000
...
98 100.000000
99 100.000000
100 100.000000
```

Name mangling 2b considered !

- <https://www.olcf.ornl.gov/tutorials/compiling-mixed-gpu-and-cpu-code>
- <https://tinyurl.com/cudafordummies/ii/14/fvcadd.f>
- <https://tinyurl.com/cudafordummies/ii/14/ntmctr.cu>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

```
cuda-zen sh@n3073-009:~$ ls
fvcadd.f ntmctr.cu
cuda-zen sh@n3073-009:~$ gfortran -c fvcadd.f
cuda-zen sh@n3073-009:~$ nvcc -c ntmctr.cu
cuda-zen sh@n3073-009:~$ nvcc fvcadd.o ntmctr.o -lcudart -lgfortran
cuda-zen sh@n3073-009:~$ ./a.out
```

Name mangling 2b considered !

C++ like compiler !

```
1 100.000000
2 100.000000
3 100.000000
4 100.000000
5 100.000000
6 100.000000
7 100.000000
...
98 100.000000
99 100.000000
100 100.000000
```

- <https://www.olcf.ornl.gov/tutorials/compiling-mixed-gpu-and-cpu-code>
- <https://tinyurl.com/cudafordummies/ii/14/fvcadd.f>
- <https://tinyurl.com/cudafordummies/ii/14/ntmctr.cu>

TO THE FORTRAN-MINDED CONT.

CUDA SDK CONT.

- Type analogues,
int ↔ INTEGER
float ↔ REAL
double ↔ DOUBLE PRECISION
- Name mangling: compiling Fortran code will alter names of called subroutines/function by adding a trailing underscore, _ Therefore C-functions will need to be aware of this and add this underscore ahead of time
- Fortran uses column-major order for storing multidimensional arrays while C uses row-major order
- Linking may need the addition of -lcudart -lgfortran etc

→ http://www.computationalmathematics.org/topics/files/calling_cuda_from_fortran.html

CUBLAS

CUDA SDK CONT.

- BLAS/LAPACK are among the most widely used libraries in scientific computing

→ <http://www.netlib.org/blas/>

CUBLAS

CUDA SDK CONT.

- BLAS/LAPACK are among the most widely used libraries in scientific computing
- General tasks in linear algebra with systematic naming scheme and a focus on high performance computing

→ <http://www.netlib.org/blas/>

CUBLAS

CUDA SDK CONT.

- BLAS/LAPACK are among the most widely used libraries in scientific computing
- General tasks in linear algebra with systematic naming scheme and a focus on high performance computing
- Standard API for many prominent implementations (MKL, OpenBLAS, ATLAS...)

→ <http://www.netlib.org/blas/>

CUBLAS

CUDA SDK CONT.

- BLAS/LAPACK are among the most widely used libraries in scientific computing
- General tasks in linear algebra with systematic naming scheme and a focus on high performance computing
- Standard API for many prominent implementations (MKL, OpenBLAS, ATLAS...)
- Originally written in Fortran

→ <http://www.netlib.org/blas/>

CUBLAS CONT.

CUDA SDK CONT.

S real, single precision

D real, double precision

C complex, single precision

Z complex, double precision

DGEMM(...)

→ <http://www.netlib.org/blas/>

CUBLAS CONT.

CUDA SDK CONT.

S real, single precision

D real, double precision

C complex, single precision

Z complex, double precision



DGEMM(...)

→ <http://www.netlib.org/blas/>

CUBLAS CONT.

CUDA SDK CONT.

SY symmetric	HE Hermitian	HP Hermitian packed	TP triangular packed
SP symmetric packed	GB general band	HB Hermitian band	TB triangular band
SB symmetric band	GE general	TR triangular	

S real, single precision

D real, double precision

C complex, single precision

Z complex, double precision



DGEMM(...)

→ <http://www.netlib.org/blas/>

CUBLAS CONT.

CUDA SDK CONT.

SY symmetric **HE** Hermitian **HP** Hermitian packed **TP** triangular packed
SP symmetric packed **GB** general band **HB** Hermitian band **TB** triangular band
SB symmetric band **GE** general **TR** triangular

S real, single precision

D real, double precision

C complex, single precision

Z complex, double precision



MV matrix · vector

SV syst eq. 1 vec unk

R rank-1 matrix update

R2 rank-2 matrix update

MM matrix · matrix

SM syst eq. matrix unk

C conjugate vector

G Givens rotation

MG mod. Givens rot constr

RK rank-k matrix update

R2K rank-2k matrix update

U unconjugate vector

M modified Givens rot

→ <http://www.netlib.org/blas/>

CUBLAS CONT.

CUDA SDK CONT.

```
...
DGEMM(TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB, BETA, C, LDC)
...

...
Atype = Btype = 'T';
Adim = Bdim = Cdim = n;
alpha = (double) 1;
beta = (double) 0;
info = dgemm( &Atype, &Btype, &Adim, &Bdim, &Cdim, &alpha,
             &A[0], &Adim, &B[0], &Bdim, &beta, &C[0], &Cdim );
if ( info != 0 ) {
    printf("blas error in dgemm returning %d\n", info);
    exit(99);
}
...

...
stat = cublasDgemm ( handle, dev_Atype, dev_Btype, Adim, Bdim, Cdim, &alpha,
                   dev_ptr_A, Adim, dev_ptr_B, Bdim, &beta, dev_ptr_C, Cdim );
if ( stat != CUBLAS_STATUS_SUCCESS ) {
    printf("cublas error in dgemm \n");
    exit(99);
}
```

→ <http://www.netlib.org/blas/>

CUBLAS CONT.

CUDA SDK CONT.

```
...
DGEMM(TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB, BETA, C, LDC)
...

...
Atype = Btype = 'T';
Adim = Bdim = Cdim = n;
alpha = (double) 1;
beta = (double) 0;
info = dgemm( &Atype, &Btype, &Adim, &Bdim, &Cdim, &alpha,
             &A[0], &Adim, &B[0], &Bdim, &beta, &C[0], &Cdim );
if ( info != 0 ) {
    printf("blas error in dgemm returning %d\n", info);
    exit(99);
}
...

...
stat = cublasDgemm ( handle, dev_Atype, dev_Btype, Adim, Bdim, Cdim, &alpha,
                   dev_ptr_A, Adim, dev_ptr_B, Bdim, &beta, dev_ptr_C, Cdim );
if ( stat != CUBLAS_STATUS_SUCCESS ) {
    printf("cublas error in dgemm \n");
    exit(99);
}
}
```

gfortran ... -lblas

→ <http://www.netlib.org/blas/>

CUBLAS CONT.

CUDA SDK CONT.

```
...  
DGEMM(TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB, BETA, C, LDC)  
...
```

gfortran ... -lblas

```
...  
Atype = Btype = 'T';  
Adim = Bdim = Cdim = n;  
alpha = (double) 1;  
beta = (double) 0;  
info = dgemm( &Atype, &Btype, &Adim, &Bdim, &Cdim, &alpha,  
             &A[0], &Adim, &B[0], &Bdim, &beta, &C[0], &Cdim );
```

```
if ( info != 0 ) {  
    printf("blas error in dgemm returning %d\n", info);  
    exit(99);  
}
```

gcc -Ddgemm=dgemm_ ... -lblas

```
...  
stat = cublasDgemm ( handle, dev_Atype, dev_Btype, Adim, Bdim, Cdim, &alpha,  
                   dev_ptr_A, Adim, dev_ptr_B, Bdim, &beta, dev_ptr_C, Cdim );  
if ( stat != CUBLAS_STATUS_SUCCESS ) {  
    printf("cublas error in dgemm \n");  
    exit(99);  
}
```

→ <http://www.netlib.org/blas/>

CUBLAS CONT.

CUDA SDK CONT.

```
...  
DGEMM(TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB, BETA, C, LDC)  
...
```

gfortran ... -lblas

```
...  
Atype = Btype = 'T';  
Adim = Bdim = Cdim = n;  
alpha = (double) 1;  
beta = (double) 0;  
info = dgemm( &Atype, &Btype, &Adim, &Bdim, &Cdim, &alpha,  
             &A[0], &Adim, &B[0], &Bdim, &beta, &C[0], &Cdim );
```

```
if ( info != 0 ) {  
    printf("blas error in dgemm returning %d\n", info);  
    exit(99);  
}
```

gcc -Ddgemm=dgemm_ ... -lblas

```
...  
...  
stat = cublasDgemm ( handle, dev_Atype, dev_Btype, Adim, Bdim, Cdim, &alpha,  
                   dev_ptr_A, Adim, dev_ptr_B, Bdim, &beta, dev_ptr_C, Cdim );
```

```
if ( stat != CUBLAS_STATUS_SUCCESS ) {  
    printf("cublas error in dgemm \n");  
    exit(99);  
}
```

nvcc ... -lcublas

→ <http://www.netlib.org/blas/>

CUBLAS CONT.

CUDA SDK CONT.

col-wise
linearized
matrices;
∀ pointers

```
...  
DGEMM(TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB, BETA, C, LDC)  
...
```

gfortran ... -lblas

```
...  
Atype = Btype = 'T';  
Adim = Bdim = Cdim = n;  
alpha = (double) 1;  
beta = (double) 0;  
info = dgemm( &Atype, &Btype, &Adim, &Bdim, &Cdim, &alpha,  
             &A[0], &Adim, &B[0], &Bdim, &beta, &C[0], &Cdim );
```

```
if ( info != 0 ) {  
    printf("blas error in dgemm returning %d\n", info);  
    exit(99);  
}  
...
```

gcc -Ddgemm=dgemm_ ... -lblas

```
...  
stat = cublasDgemm ( handle, dev_Atype, dev_Btype, Adim, Bdim, Cdim, &alpha,  
                   dev_ptr_A, Adim, dev_ptr_B, Bdim, &beta, dev_ptr_C, Cdim );  
if ( stat != CUBLAS_STATUS_SUCCESS ) {  
    printf("cublas error in dgemm \n");  
    exit(99);  
}
```

nvcc ... -lcublas

→ <http://www.netlib.org/blas/>

CUBLAS CONT.

CUDA SDK CONT.

```
...  
DGEMM(TRANSA, TRANSB, M, N, K, ALPHA, A, LDA, B, LDB, BETA, C, LDC)  
...
```

gfortran ... -lblas

```
...  
Atype = Btype = 'T';  
Adim = Bdim = Cdim = n;  
alpha = (double) 1;  
beta = (double) 0;  
info = dgemm( &Atype, &Btype, &Adim, &Bdim, &Cdim, &alpha,  
             &A[0], &Adim, &B[0], &Bdim, &beta, &C[0], &Cdim );
```

```
if ( info != 0 ) {  
    printf("blas error in dgemm returning %d\n", info);  
    exit(99);  
}  
...
```

gcc -Ddgemm=dgemm_ ... -lblas

```
...  
stat = cublasDgemm ( handle, dev_Atype, dev_Btype, Adim, Bdim, Cdim, &alpha,  
                   dev_ptr_A, Adim, dev_ptr_B, Bdim, &beta, dev_ptr_C, Cdim );  
if ( stat != CUBLAS_STATUS_SUCCESS ) {  
    printf("cublas error in dgemm \n");  
    exit(99);  
}
```

nvcc ... -lcublas

col-wise
linearized
matrices;
∇pointers

separate
cublas
arg/types

→ <http://www.netlib.org/blas/>

Porting code to CUBLAS requires a more-or-less standard procedure of the following steps:

1. Initiate the CUBLAS context (`cublasCreate()`)
2. Allocate device memory using `cudaMalloc()` !
3. Transfer content of host arrays to the device (`cublasSetMatrix()`)
4. Call a specific CUBLAS routine, e.g. `cublasDgemm()`
5. Transfer back the result from device memory to host (`cublasGetMatrix()`)
6. Free device memory
7. Destroy CUBLAS context

→ <https://docs.nvidia.com/cuda/cublas/index.html>

→ <https://devtalk.nvidia.com/default/topic/1047981/b/t/post/5318441>

Example cublasDgemm()

```
#include <stdio.h>
#include <stdlib.h>
#include <cuda_runtime.h>
#include "cublas_v2.h"

int main(int argc, char **argv)
{
    int N, i, j;
    double alpha, beta, *A, *Adev, *B, *Bdev, *C, *Cdev;
    cublasStatus_t stat;
    cublasHandle_t handle;
    cublasOperation_t Atype, Btype;

    // memory allocation and parameter set up
    N = 5;
    alpha = (double) 1;
    beta = (double) 0;
    Atype = CUBLAS_OP_N;
    Btype = CUBLAS_OP_N;
    A = (double *) malloc(N * N * sizeof(double));
    B = (double *) malloc(N * N * sizeof(double));
    C = (double *) malloc(N * N * sizeof(double));

    // fill matrices A[] and B[] with data
```

→ <https://docs.nvidia.com/cuda/cublas/index.html>

CUBLAS CONT.

CUDA SDK CONT.

Example cublasDgemm()

```
#include <stdio.h>
#include <stdlib.h>
#include <cuda_runtime.h>
#include "cublas_v2.h"

int main(int argc, char **argv)
{
    int N, i, j;
    double alpha, beta, *A, *Adev, *B, *Bdev, *C, *Cdev;
    cublasStatus_t stat;
    cublasHandle_t handle;
    cublasOperation_t Atype, Btype;

    // memory allocation and parameter set up
    N = 5;
    alpha = (double) 1;
    beta = (double) 0;
    Atype = CUBLAS_OP_N;
    Btype = CUBLAS_OP_N;
    A = (double *) malloc(N * N * sizeof(double));
    B = (double *) malloc(N * N * sizeof(double));
    C = (double *) malloc(N * N * sizeof(double));

    // fill matrices A[] and B[] with data
```

go for the new
CUBLAS library API

→ <https://docs.nvidia.com/cuda/cublas/index.html>

CUBLAS CONT.

CUDA SDK CONT.

Example cublasDgemm()

```
#include <stdio.h>
#include <stdlib.h>
#include <cuda_runtime.h>
#include "cublas_v2.h"

int main(int argc, char **argv)
{
    int N, i, j;
    double alpha, beta, *A, *Adev, *B, *Bdev, *C, *Cdev;
    cublasStatus_t stat;
    cublasHandle_t handle;
    cublasOperation_t Atype, Btype;

    // memory allocation and parameter set up
    N = 5;
    alpha = (double) 1;
    beta = (double) 0;
    Atype = CUBLAS_OP_N;
    Btype = CUBLAS_OP_N;
    A = (double *) malloc(N * N * sizeof(double));
    B = (double *) malloc(N * N * sizeof(double));
    C = (double *) malloc(N * N * sizeof(double));

    // fill matrices A[] and B[] with data
```

go for the new
CUBLAS library API

will establish CUBLAS context

→ <https://docs.nvidia.com/cuda/cublas/index.html>

CUBLAS CONT.

CUDA SDK CONT.

Example cublasDgemm()

```
#include <stdio.h>
#include <stdlib.h>
#include <cuda_runtime.h>
#include "cublas_v2.h"

int main(int argc, char **argv)
{
    int N, i, j;
    double alpha, beta, *A, *Adev, *B, *Bdev, *C, *Cdev;
    cublasStatus_t stat;
    cublasHandle_t handle;
    cublasOperation_t Atype, Btype;

    // memory allocation and parameter set up
    N = 5;
    alpha = (double) 1;
    beta = (double) 0;
    Atype = CUBLAS_OP_N;
    Btype = CUBLAS_OP_N;
    A = (double *) malloc(N * N * sizeof(double));
    B = (double *) malloc(N * N * sizeof(double));
    C = (double *) malloc(N * N * sizeof(double));

    // fill matrices A[] and B[] with data
```

go for the new
CUBLAS library API

will establish CUBLAS context

CUBLAS way of
TRANSA/B='T'V'N'

→ <https://docs.nvidia.com/cuda/cublas/index.html>

CUBLAS CONT.

CUDA SDK CONT.

Example cublasDgemm()

```
#include <stdio.h>
#include <stdlib.h>
#include <cuda_runtime.h>
#include "cublas_v2.h"

int main(int argc, char **argv)
{
    int N, i, j;
    double alpha, beta, *A, *Adev, *B, *Bdev, *C, *Cdev;
    cublasStatus_t stat;
    cublasHandle_t handle;
    cublasOperation_t Atype, Btype;

    // memory allocation and parameter set up
    N = 5;
    alpha = (double) 1;
    beta = (double) 0;
    Atype = CUBLAS_OP_N;
    Btype = CUBLAS_OP_N;
    A = (double *) malloc(N * N * sizeof(double));
    B = (double *) malloc(N * N * sizeof(double));
    C = (double *) malloc(N * N * sizeof(double));

    // fill matrices A[] and B[] with data
```

go for the new
CUBLAS library API

will establish CUBLAS context

CUBLAS way of
TRANSA/B='T'∨'N'

normal memory
allocation

→ <https://docs.nvidia.com/cuda/cublas/index.html>

Example cublasDgemm() cont.

```
// initiate the CUBLAS context
cublasCreate(&handle);
// allocate device memory
cudaMalloc(&Adev, N * N * sizeof(double));
cudaMalloc(&Bdev, N * N * sizeof(double));
cudaMalloc(&Cdev, N * N * sizeof(double));
// copy contents of arrays A[], B[] into device memory
cublasSetMatrix(N, N, sizeof(double), &A[0], N, Adev, N);
cublasSetMatrix(N, N, sizeof(double), &B[0], N, Bdev, N);
// carry out CUBLAS operation
cublasDgemm(handle, Atype, Btype, N, N, N, &alpha, Adev, N, Bdev, N, &beta, Cdev, N);
// retrieve results from device memory
cublasGetMatrix(N, N, sizeof(double), Cdev, N, &C[0], N);
// do something with the result
... = C[]
// make clean
cudaFree(Cdev); cudaFree(Bdev); cudaFree(Adev);
cublasDestroy(handle);
// free host memory and return
free(C); free(B); free(A);
return(0);
}
```

→ <https://docs.nvidia.com/cuda/cublas/index.html>

Example cublasDgemm() cont.

```
// initiate the CUBLAS context
cublasCreate(&handle);
// allocate device memory
cudaMalloc(&Adev, N * N * sizeof(double));
cudaMalloc(&Bdev, N * N * sizeof(double));
cudaMalloc(&Cdev, N * N * sizeof(double));
// copy contents of arrays A[], B[] into device memory
cublasSetMatrix(N, N, sizeof(double), &A[0], N, Adev, N);
cublasSetMatrix(N, N, sizeof(double), &B[0], N, Bdev, N);
// carry out CUBLAS operation
cublasDgemm(handle, Atype, Btype, N, N, N, &alpha, Adev, N, Bdev, N, &beta, Cdev, N);
// retrieve results from device memory
cublasGetMatrix(N, N, sizeof(double), Cdev, N, &C[0], N);
// do something with the result
... = C[]
// make clean
cudaFree(Cdev); cudaFree(Bdev); cudaFree(Adev);
cublasDestroy(handle);
// free host memory and return
free(C); free(B); free(A);
return(0);
}
```

Step 1

→ <https://docs.nvidia.com/cuda/cublas/index.html>

Example cublasDgemm() cont.

```
// initiate the CUBLAS context
cublasCreate(&handle);
// allocate device memory
cudaMalloc(&Adev, N * N * sizeof(double));
cudaMalloc(&Bdev, N * N * sizeof(double));
cudaMalloc(&Cdev, N * N * sizeof(double));
// copy contents of arrays A[], B[] into device memory
cublasSetMatrix(N, N, sizeof(double), &A[0], N, Adev, N);
cublasSetMatrix(N, N, sizeof(double), &B[0], N, Bdev, N);
// carry out CUBLAS operation
cublasDgemm(handle, Atype, Btype, N, N, N, &alpha, Adev, N, Bdev, N, &beta, Cdev, N);
// retrieve results from device memory
cublasGetMatrix(N, N, sizeof(double), Cdev, N, &C[0], N);
// do something with the result
... = C[]
// make clean
cudaFree(Cdev); cudaFree(Bdev); cudaFree(Adev);
cublasDestroy(handle);
// free host memory and return
free(C); free(B); free(A);
return(0);
}
```

Step 1

Step 2

→ <https://docs.nvidia.com/cuda/cublas/index.html>

Example cublasDgemm() cont.

```
// initiate the CUBLAS context
cublasCreate(&handle);
// allocate device memory
cudaMalloc(&Adev, N * N * sizeof(double));
cudaMalloc(&Bdev, N * N * sizeof(double));
cudaMalloc(&Cdev, N * N * sizeof(double));
// copy contents of arrays A[], B[] into device memory
cublasSetMatrix(N, N, sizeof(double), &A[0], N, Adev, N);
cublasSetMatrix(N, N, sizeof(double), &B[0], N, Bdev, N);
// carry out CUBLAS operation
cublasDgemm(handle, Atype, Btype, N, N, N, &alpha, Adev, N, Bdev, N, &beta, Cdev, N);
// retrieve results from device memory
cublasGetMatrix(N, N, sizeof(double), Cdev, N, &C[0], N);
// do something with the result
... = C[]
// make clean
cudaFree(Cdev); cudaFree(Bdev); cudaFree(Adev);
cublasDestroy(handle);
// free host memory and return
free(C); free(B); free(A);
return(0);
}
```

Step 1

Step 2

Step 3

→ <https://docs.nvidia.com/cuda/cublas/index.html>

Example cublasDgemm() cont.

```
// initiate the CUBLAS context
cublasCreate(&handle);
// allocate device memory
cudaMalloc(&Adev, N * N * sizeof(double));
cudaMalloc(&Bdev, N * N * sizeof(double));
cudaMalloc(&Cdev, N * N * sizeof(double));
// copy contents of arrays A[], B[] into device memory
cublasSetMatrix(N, N, sizeof(double), &A[0], N, Adev, N);
cublasSetMatrix(N, N, sizeof(double), &B[0], N, Bdev, N);
// carry out CUBLAS operation
cublasDgemm(handle, Atype, Btype, N, N, N, &alpha, Adev, N, Bdev, N, &beta, Cdev, N);
// retrieve results from device memory
cublasGetMatrix(N, N, sizeof(double), Cdev, N, &C[0], N);
// do something with the result
... = C[]
// make clean
cudaFree(Cdev); cudaFree(Bdev); cudaFree(Adev);
cublasDestroy(handle);
// free host memory and return
free(C); free(B); free(A);
return(0);
}
```

Step 1

Step 2

Step 3

Step 4

→ <https://docs.nvidia.com/cuda/cublas/index.html>

Example cublasDgemm() cont.

```
// initiate the CUBLAS context
cublasCreate(&handle);
// allocate device memory
cudaMalloc(&Adev, N * N * sizeof(double));
cudaMalloc(&Bdev, N * N * sizeof(double));
cudaMalloc(&Cdev, N * N * sizeof(double));
// copy contents of arrays A[], B[] into device memory
cublasSetMatrix(N, N, sizeof(double), &A[0], N, Adev, N);
cublasSetMatrix(N, N, sizeof(double), &B[0], N, Bdev, N);
// carry out CUBLAS operation
cublasDgemm(handle, Atype, Btype, N, N, N, &alpha, Adev, N, Bdev, N, &beta, Cdev, N);
// retrieve results from device memory
cublasGetMatrix(N, N, sizeof(double), Cdev, N, &C[0], N);
// do something with the result
... = C[]
// make clean
cudaFree(Cdev); cudaFree(Bdev); cudaFree(Adev);
cublasDestroy(handle);
// free host memory and return
free(C); free(B); free(A);
return(0);
}
```

Step 1

Step 2

Step 3

Step 4

Step 5

→ <https://docs.nvidia.com/cuda/cublas/index.html>

Example cublasDgemm() cont.

```
// initiate the CUBLAS context
cublasCreate(&handle);
// allocate device memory
cudaMalloc(&Adev, N * N * sizeof(double));
cudaMalloc(&Bdev, N * N * sizeof(double));
cudaMalloc(&Cdev, N * N * sizeof(double));
// copy contents of arrays A[], B[] into device memory
cublasSetMatrix(N, N, sizeof(double), &A[0], N, Adev, N);
cublasSetMatrix(N, N, sizeof(double), &B[0], N, Bdev, N);
// carry out CUBLAS operation
cublasDgemm(handle, Atype, Btype, N, N, N, &alpha, Adev, N, Bdev, N, &beta, Cdev, N);
// retrieve results from device memory
cublasGetMatrix(N, N, sizeof(double), Cdev, N, &C[0], N);
// do something with the result
... = C[]
// make clean
cudaFree(Cdev); cudaFree(Bdev); cudaFree(Adev);
cublasDestroy(handle);
// free host memory and return
free(C); free(B); free(A);
return(0);
}
```

Step 1

Step 2

Step 3

Step 4

Step 5

Step 6

→ <https://docs.nvidia.com/cuda/cublas/index.html>

Example cublasDgemm() cont.

```

// initiate the CUBLAS context
cublasCreate(&handle);
// allocate device memory
cudaMalloc(&Adev, N * N * sizeof(double));
cudaMalloc(&Bdev, N * N * sizeof(double));
cudaMalloc(&Cdev, N * N * sizeof(double));
// copy contents of arrays A[], B[] into device memory
cublasSetMatrix(N, N, sizeof(double), &A[0], N, Adev, N);
cublasSetMatrix(N, N, sizeof(double), &B[0], N, Bdev, N);
// carry out CUBLAS operation
cublasDgemm(handle, Atype, Btype, N, N, N, &alpha, Adev, N, Bdev, N, &beta, Cdev, N);
// retrieve results from device memory
cublasGetMatrix(N, N, sizeof(double), Cdev, N, &C[0], N);
// do something with the result
... = C[]
// make clean
cudaFree(Cdev); cudaFree(Bdev); cudaFree(Adev);
cublasDestroy(handle);
// free host memory and return
free(C); free(B); free(A);
return(0);
}

```

Step 1

Step 2

Step 3

Step 4

Step 5

Step 6

Step 7

→ <https://docs.nvidia.com/cuda/cublas/index.html>

Words of caution:

- Always start with a small explicit example to check/confirm each individual step
- Formats of matrices are crucial ! 1D-representation, column-wise linearized;
- (double) on consumer grade cards is slower by $\approx 1/64$ than (float)
- Two new APIs, CUBLASXT (for CUDA ≥ 6.0) and cuBLASLt (for CUDA 10.1)
- Pity that we can't use `cudaMallocManaged()` at this point

→ <https://docs.nvidia.com/cuda/cublas/index.html>

→ <https://devtalk.nvidia.com/default/topic/1047981/b/t/post/5318441>

Next level of problems:

$$\underbrace{\begin{pmatrix} 1.96 & -6.49 & -0.47 & -7.20 & -0.65 \\ -6.49 & 3.80 & -6.39 & 1.50 & -6.34 \\ -0.47 & -6.39 & 4.17 & -1.51 & 2.67 \\ -7.20 & 1.50 & -1.51 & 5.70 & 1.80 \\ -0.65 & -6.34 & 2.67 & 1.80 & -7.10 \end{pmatrix}}_A$$

- For example, what are the eigenvalues and corresponding eigenvectors of A ?
- $Ax = \lambda x$
- A typical LAPACK problem — usually tightly coupled to BLAS
- Not on the GPU ! CUBLAS is BLAS only !
- However, there are cuSolver and MAGMA

→ <https://docs.nvidia.com/cuda/cusolver>

→ https://icl.cs.utk.edu/projectsfiles/magma/doxygen/group__cublas__const.html

Making use of CUSOLVER is schematically very similar to using CUBLAS:

1. Initiate the CUSOLVER context (`cusolverDnCreate()`)
2. Allocate device memory using `cudaMalloc()` !
3. Transfer content of host arrays to the device (`cudaMemcpy()`)
4. Semi-automatically set up working space required by CUSOLVER routine, e.g. `WORK`, `LWORK` in LAPACK jargon
5. Call a specific CUSOLVER routine, e.g. `cusolverDnDsyevd()`
6. Transfer back the result from device memory to host (`cudaMemcpy()`)
7. Free device memory
8. Destroy CUSOLVER context

→ https://docs.nvidia.com/cuda/cusolver/index.html#eig_examples

Example cusolverDnDsyevd()

```
#include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <cuda_runtime.h>
#include "cusolverDn.h"

int main(int argc, char **argv)
{
    int i, j, lwork, info_gpu, *devInfo;
    double *d_A, *d_W, *d_work;
    cusolverDnHandle_t cusolverH;
    cusolverStatus_t cusolver_status;
    cusolverEigMode_t jobz;
    cublasFillMode_t uplo;
    cudaError_t cudaStat;
    const int m = 5;
    const int lda = m;
    double W[m];
    double V[lda*m];
    double A[lda*m] = { 1.96, -6.49, -0.47, -7.20, -0.65,
                       -6.49, 3.80, -6.39, 1.50, -6.34,
                       -0.47, -6.39, 4.17, -1.51, 2.67,
                       -7.20, 1.50, -1.51, 5.70, 1.80,
                       -0.65, -6.34, 2.67, 1.80, -7.10};
```

CUSOLVER CONT.

CUDA SDK CONT.

Example cusolverDnDsyevd()

```
#include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <cuda_runtime.h>
#include "cusolverDn.h"

int main(int argc, char **argv)
{
    int i, j, lwork, info_gpu, *devInfo;
    double *d_A, *d_W, *d_work;
    cusolverDnHandle_t cusolverH;
    cusolverStatus_t cusolver_status;
    cusolverEigMode_t jobz;
    cublasFillMode_t uplo;
    cudaError_t cudaStat;
    const int m = 5;
    const int lda = m;
    double W[m];
    double V[lda*m];
    double A[lda*m] = { 1.96, -6.49, -0.47, -7.20, -0.65,
                       -6.49, 3.80, -6.39, 1.50, -6.34,
                       -0.47, -6.39, 4.17, -1.51, 2.67,
                       -7.20, 1.50, -1.51, 5.70, 1.80,
                       -0.65, -6.34, 2.67, 1.80, -7.10};
```

CUSOLVER header
for 'dense' subset

CUSOLVER CONT.

CUDA SDK CONT.

Example cusolverDnDsyevd()

```
#include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <cuda_runtime.h>
#include "cusolverDn.h"

int main(int argc, char **argv)
{
    int i, j, lwork, info_gpu, *devInfo;
    double *d_A, *d_W, *d_work;
    cusolverDnHandle_t cusolverH;
    cusolverStatus_t cusolver_status;
    cusolverEigMode_t jobz;
    cublasFillMode_t uplo;
    cudaError_t cudaStat;
    const int m = 5;
    const int lda = m;
    double W[m];
    double V[lda*m];
    double A[lda*m] = { 1.96, -6.49, -0.47, -7.20, -0.65,
                      -6.49, 3.80, -6.39, 1.50, -6.34,
                      -0.47, -6.39, 4.17, -1.51, 2.67,
                      -7.20, 1.50, -1.51, 5.70, 1.80,
                      -0.65, -6.34, 2.67, 1.80, -7.10};
```

CUSOLVER header
for 'dense' subset

4 CUSOLVER context

Example cusolverDnDsyevd()

```

#include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <cuda_runtime.h>
#include "cusolverDn.h"

int main(int argc, char **argv)
{
    int i, j, lwork, info_gpu, *devInfo;
    double *d_A, *d_W, *d_work;
    cusolverDnHandle_t cusolverH;
    cusolverStatus_t cusolver_status;
    cusolverEigMode_t jobz;
    cublasFillMode_t uplo;
    cudaError_t cudaStat;
    const int m = 5;
    const int lda = m;
    double W[m];
    double V[lda*m];
    double A[lda*m] = { 1.96, -6.49, -0.47, -7.20, -0.65,
                       -6.49, 3.80, -6.39, 1.50, -6.34,
                       -0.47, -6.39, 4.17, -1.51, 2.67,
                       -7.20, 1.50, -1.51, 5.70, 1.80,
                       -0.65, -6.34, 2.67, 1.80, -7.10};

```

CUSOLVER header
for 'dense' subset

4 CUSOLVER context

CUSOLVER flags 4 LAPACK's

Example cusolverDnDsyevd() cont.

```
jobz = CUSOLVER_EIG_MODE_VECTOR;
uplo = CUBLAS_FILL_MODE_LOWER;
cusolver_status = cusolverDnCreate(&cusolverH);
assert(CUSOLVER_STATUS_SUCCESS == cusolver_status);
cudaStat = cudaMalloc ((void**)&d_A, sizeof(double) * lda * m);
cudaStat = cudaMalloc ((void**)&d_W, sizeof(double) * m);
cudaStat = cudaMalloc ((void**)&devInfo, sizeof(int));
cudaStat = cudaMemcpy(d_A, A, sizeof(double) * lda * m, cudaMemcpyHostToDevice);

cusolver_status = cusolverDnDsyevd_bufferSize(cusolverH, jobz, uplo, m, d_A, lda, d_W, &lwork);
cudaStat = cudaMalloc((void**)&d_work, sizeof(double)*lwork);
cusolver_status = cusolverDnDsyevd(cusolverH, jobz, uplo, m, d_A, lda, d_W, d_work, lwork, devInfo);
cudaStat = cudaDeviceSynchronize();

cudaStat = cudaMemcpy(W, d_W, sizeof(double)*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(V, d_A, sizeof(double)*lda*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(&info_gpu, devInfo, sizeof(int), cudaMemcpyDeviceToHost);
... do something with W[] and V[] ...
cudaFree(d_A); cudaFree(d_W); cudaFree(devInfo); cudaFree(d_work);
cusolverDnDestroy(cusolverH);
cudaDeviceReset();

return(0);
}
```

Example cusolverDnDsyevd() cont.

establish CUSOLVER context

```
jobz = CUSOLVER_EIG_MODE_VECTOR;
uplo = CUBLAS_FILL_MODE_LOWER;
cusolver_status = cusolverDnCreate(&cusolverH); ←
assert(CUSOLVER_STATUS_SUCCESS == cusolver_status);
cudaStat = cudaMalloc ((void**)&d_A, sizeof(double) * lda * m);
cudaStat = cudaMalloc ((void**)&d_W, sizeof(double) * m);
cudaStat = cudaMalloc ((void**)&devInfo, sizeof(int));
cudaStat = cudaMemcpy(d_A, A, sizeof(double) * lda * m, cudaMemcpyHostToDevice);

cusolver_status = cusolverDnDsyevd_bufferSize(cusolverH, jobz, uplo, m, d_A, lda, d_W, &lwork);
cudaStat = cudaMalloc((void**)&d_work, sizeof(double)*lwork);
cusolver_status = cusolverDnDsyevd(cusolverH, jobz, uplo, m, d_A, lda, d_W, d_work, lwork, devInfo);
cudaStat = cudaDeviceSynchronize();

cudaStat = cudaMemcpy(W, d_W, sizeof(double)*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(V, d_A, sizeof(double)*lda*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(&info_gpu, devInfo, sizeof(int), cudaMemcpyDeviceToHost);
... do something with W[] and V[] ...
cudaFree(d_A); cudaFree(d_W); cudaFree(devInfo); cudaFree(d_work);
cusolverDnDestroy(cusolverH);
cudaDeviceReset();

return(0);
}
```


Example cusolverDnDsyevd() cont.

```
jobz = CUSOLVER_EIG_MODE_VECTOR;
uplo = CUBLAS_FILL_MODE_LOWER;
cusolver_status = cusolverDnCreate(&cusolverH);
assert(CUSOLVER_STATUS_SUCCESS == cusolver_status);
cudaStat = cudaMalloc ((void**)&d_A, sizeof(double) * lda * m);
cudaStat = cudaMalloc ((void**)&d_W, sizeof(double) * m);
cudaStat = cudaMalloc ((void**)&devInfo, sizeof(int));
cudaStat = cudaMemcpy(d_A, A, sizeof(double) * lda * m, cudaMemcpyHostToDevice);

cusolver_status = cusolverDnDsyevd_bufferSize(cusolverH, jobz, uplo, m, d_A, lda, d_W, &lwork);
cudaStat = cudaMalloc((void**)&d_work, sizeof(double)*lwork);
cusolver_status = cusolverDnDsyevd(cusolverH, jobz, uplo, m, d_A, lda, d_W, d_work, lwork, devInfo);
cudaStat = cudaDeviceSynchronize();

cudaStat = cudaMemcpy(W, d_W, sizeof(double)*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(V, d_A, sizeof(double)*lda*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(&info_gpu, devInfo, sizeof(int), cudaMemcpyDeviceToHost);
... do something with W[] and V[] ...
cudaFree(d_A); cudaFree(d_W); cudaFree(devInfo); cudaFree(d_work);
cusolverDnDestroy(cusolverH);
cudaDeviceReset();

return(0);
}
```

establish CUSOLVER context

error check

Example cusolverDnDsyevd() cont.

```

jobz = CUSOLVER_EIG_MODE_VECTOR;
uplo = CUBLAS_FILL_MODE_LOWER;
cusolver_status = cusolverDnCreate(&cusolverH);
assert(CUSOLVER_STATUS_SUCCESS == cusolver_status);
cudaStat = cudaMalloc ((void**)&d_A, sizeof(double) * lda * m);
cudaStat = cudaMalloc ((void**)&d_W, sizeof(double) * m);
cudaStat = cudaMalloc ((void**)&devInfo, sizeof(int));
cudaStat = cudaMemcpy(d_A, A, sizeof(double) * lda * m, cudaMemcpyHostToDevice);

cusolver_status = cusolverDnDsyevd_bufferSize(cusolverH, jobz, uplo, m, d_A, lda, d_W, &lwork);
cudaStat = cudaMalloc((void**)&d_work, sizeof(double)*lwork);
cusolver_status = cusolverDnDsyevd(cusolverH, jobz, uplo, m, d_A, lda, d_W, d_work, lwork, devInfo);
cudaStat = cudaDeviceSynchronize();

cudaStat = cudaMemcpy(W, d_W, sizeof(double)*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(V, d_A, sizeof(double)*lda*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(&info_gpu, devInfo, sizeof(int), cudaMemcpyDeviceToHost);
... do something with W[] and V[] ...
cudaFree(d_A); cudaFree(d_W); cudaFree(devInfo); cudaFree(d_work);
cusolverDnDestroy(cusolverH);
cudaDeviceReset();

return(0);
}

```

establish CUSOLVER context

error check

standard procedure

Example cusolverDnDsyevd() cont.

adjust
work
arrays

establish CUSOLVER context

error check

standard procedure

```

jobz = CUSOLVER_EIG_MODE_VECTOR;
uplo = CUBLAS_FILL_MODE_LOWER;
cusolver_status = cusolverDnCreate(&cusolverH);
assert(CUSOLVER_STATUS_SUCCESS == cusolver_status);
cudaStat = cudaMalloc ((void**)&d_A, sizeof(double) * lda * m);
cudaStat = cudaMalloc ((void**)&d_W, sizeof(double) * m);
cudaStat = cudaMalloc ((void**)&devInfo, sizeof(int));
cudaStat = cudaMemcpy(d_A, A, sizeof(double) * lda * m, cudaMemcpyHostToDevice);

cusolver_status = cusolverDnDsyevd_bufferSize(cusolverH, jobz, uplo, m, d_A, lda, d_W, &lwork);
cudaStat = cudaMalloc((void**)&d_work, sizeof(double)*lwork);
cusolver_status = cusolverDnDsyevd(cusolverH, jobz, uplo, m, d_A, lda, d_W, d_work, lwork, devInfo);
cudaStat = cudaDeviceSynchronize();

cudaStat = cudaMemcpy(W, d_W, sizeof(double)*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(V, d_A, sizeof(double)*lda*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(&info_gpu, devInfo, sizeof(int), cudaMemcpyDeviceToHost);
... do something with W[] and V[] ...
cudaFree(d_A); cudaFree(d_W); cudaFree(devInfo); cudaFree(d_work);
cusolverDnDestroy(cusolverH);
cudaDeviceReset();

return(0);
}
    
```

Example cusolverDnDsyevd() cont.

adjust
work
arrays

establish CUSOLVER context

error check

standard procedure

```

jobz = CUSOLVER_EIG_MODE_VECTOR;
uplo = CUBLAS_FILL_MODE_LOWER;
cusolver_status = cusolverDnCreate(&cusolverH);
assert(CUSOLVER_STATUS_SUCCESS == cusolver_status);
cudaStat = cudaMalloc ((void**)&d_A, sizeof(double) * lda * m);
cudaStat = cudaMalloc ((void**)&d_W, sizeof(double) * m);
cudaStat = cudaMalloc ((void**)&devInfo, sizeof(int));
cudaStat = cudaMemcpy(d_A, A, sizeof(double) * lda * m, cudaMemcpyHostToDevice);

cusolver_status = cusolverDnDsyevd_bufferSize(cusolverH, jobz, uplo, m, d_A, lda, d_W, &lwork);
cudaStat = cudaMalloc((void**)&d_work, sizeof(double)*lwork);
cusolver_status = cusolverDnDsyevd(cusolverH, jobz, uplo, m, d_A, lda, d_W, d_work, lwork, devInfo);
cudaStat = cudaDeviceSynchronize();

cudaStat = cudaMemcpy(W, d_W, sizeof(double)*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(V, d_A, sizeof(double)*lda*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(&info_gpu, devInfo, sizeof(int), cudaMemcpyDeviceToHost);
... do something with W[] and V[] ...
cudaFree(d_A); cudaFree(d_W); cudaFree(devInfo); cudaFree(d_work);
cusolverDnDestroy(cusolverH);
cudaDeviceReset();

return(0);
}
    
```

call
solver

CUSOLVER CONT.

CUDA SDK CONT.

Example cusolverDnDsyevd() cont.

adjust
work
arrays

```
jobz = CUSOLVER_EIG_MODE_VECTOR;
uplo = CUBLAS_FILL_MODE_LOWER;
cusolver_status = cusolverDnCreate(&cusolverH);
assert(CUSOLVER_STATUS_SUCCESS == cusolver_status);
cudaStat = cudaMalloc ((void**)&d_A, sizeof(double) * lda * m);
cudaStat = cudaMalloc ((void**)&d_W, sizeof(double) * m);
cudaStat = cudaMalloc ((void**)&devInfo, sizeof(int));
cudaStat = cudaMemcpy(d_A, A, sizeof(double) * lda * m, cudaMemcpyHostToDevice);

cusolver_status = cusolverDnDsyevd_bufferSize(cusolverH, jobz, uplo, m, d_A, lda, d_W, &lwork);
cudaStat = cudaMalloc((void**)&d_work, sizeof(double)*lwork);
cusolver_status = cusolverDnDsyevd(cusolverH, jobz, uplo, m, d_A, lda, d_W, d_work, lwork, devInfo);
cudaStat = cudaDeviceSynchronize();

cudaStat = cudaMemcpy(W, d_W, sizeof(double)*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(V, d_A, sizeof(double)*lda*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(&info_gpu, devInfo, sizeof(int), cudaMemcpyDeviceToHost);
... do something with W[] and V[] ...
cudaFree(d_A); cudaFree(d_W); cudaFree(devInfo); cudaFree(d_work);
cusolverDnDestroy(cusolverH);
cudaDeviceReset();

return(0);
}
```

establish CUSOLVER context

error check

standard procedure

call
solver

back copy results

Example cusolverDnDsyevd() cont.

adjust
work
arrays

establish CUSOLVER context

error check

standard procedure

```

jobz = CUSOLVER_EIG_MODE_VECTOR;
uplo = CUBLAS_FILL_MODE_LOWER;
cusolver_status = cusolverDnCreate(&cusolverH);
assert(CUSOLVER_STATUS_SUCCESS == cusolver_status);
cudaStat = cudaMalloc ((void**)&d_A, sizeof(double) * lda * m);
cudaStat = cudaMalloc ((void**)&d_W, sizeof(double) * m);
cudaStat = cudaMalloc ((void**)&devInfo, sizeof(int));
cudaStat = cudaMemcpy(d_A, A, sizeof(double) * lda * m, cudaMemcpyHostToDevice);

cusolver_status = cusolverDnDsyevd_bufferSize(cusolverH, jobz, uplo, m, d_A, lda, d_W, &lwork);
cudaStat = cudaMalloc((void**)&d_work, sizeof(double)*lwork);
cusolver_status = cusolverDnDsyevd(cusolverH, jobz, uplo, m, d_A, lda, d_W, d_work, lwork, devInfo);
cudaStat = cudaDeviceSynchronize();

cudaStat = cudaMemcpy(W, d_W, sizeof(double)*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(V, d_A, sizeof(double)*lda*m, cudaMemcpyDeviceToHost);
cudaStat = cudaMemcpy(&info_gpu, devInfo, sizeof(int), cudaMemcpyDeviceToHost);
... do something with W[] and V[] ...
cudaFree(d_A); cudaFree(d_W); cudaFree(devInfo); cudaFree(d_work);
cusolverDnDestroy(cusolverH);
cudaDeviceReset();

return(0);
}
    
```

call
solver

back copy results

finalize

- Very straightforward procedure offering significant speed-up for large scale problems at very minor porting effort
- Again, (float) \leftrightarrow (double) gap in terms of performance, $\approx 1/64$ on consumer cards !
- A related problem of considerable interest may be addressed via `cusolverEigType_t` — generalized symmetric-definite eigenvalue problem (section 2.2.1.4)
 $A x = \lambda B x$
- Sparse problems covered too — `cuSolverSP`

→ <https://docs.nvidia.com/cuda/cusolver/index.html>

MORE LIBRARIES

CUDA SDK CONT.

- CUFFT — CUDA Fast Fourier Transform
- CURAND — CUDA random number generation library
- CUDNN & TENSORRT — Deep Learning, training & inference
- CUSPARSE — Basic linear algebra for sparse matrices
- NVBLAS — Another replacement for BLAS calls with little effort of porting (re-linking or LD_PRELOAD)
- NVGRAPH — Big Data analytics via graph problems

→ <https://docs.nvidia.com/cuda/cufft/index.html>

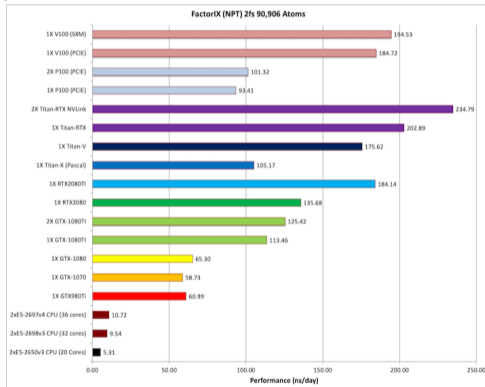
→ <https://docs.nvidia.com/cuda/curand/index.html>

→ <https://docs.nvidia.com/deeplearning/sdk/index.html>

NUMERICAL ACCURACY & PERFORMANCE

CUDA SDK CONT.

How come that scientific apps, e.g. AMBER's pmemd.cuda are doing so well on consumer grade GPUs ?



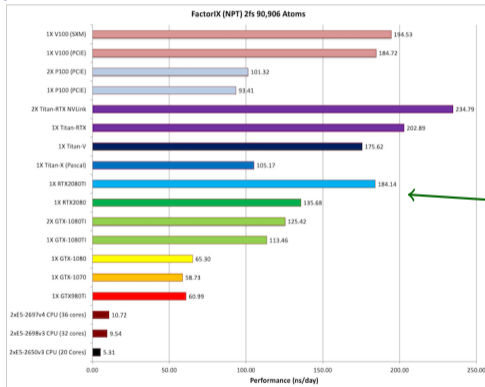
→ <http://ambermd.org/GPUPerformance.php>

→ <https://www.sciencedirect.com/science/article/pii/S0010465512003098>

NUMERICAL ACCURACY & PERFORMANCE

CUDA SDK CONT.

How come that scientific apps, e.g. AMBER's pmemd.cuda are doing so well on consumer grade GPUs ?



Combine several 32bit types (float, int) to approximate double

→ <http://ambermd.org/GPUPerformance.php>

→ <https://www.sciencedirect.com/science/article/pii/S0010465512003098>

Dekker (1971) and Knuth (1969)

Quasi-double precision method

- Combine two floating-point numbers to express one variable

$$\begin{array}{ll} x + y = z + zz & x = 10.2, y = 0.345 \\ \text{(where } |z / zz| < 2^{24}\text{)} & \text{(three significant digits)} \\ z = \text{fl}(x + y) & z = 10.5 \\ w = \text{fl}(z - x) & w = 0.3 \\ zz = \text{fl}(y - w) & zz = 0.045 \\ & (z + zz = 10.545, \text{ no cancellation}) \end{array}$$

x, y, z, zz : FP32 variable

$\text{fl}(\text{operation})$: FP32 arithmetic operation

- More than twice operations are needed for similar accuracy to double-precision

⇒ Only works for consumer GPUs

(slide courtesy of Prof Tetsu Narumi, UEC, Tokyo)

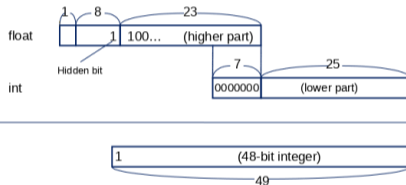
→ <https://link.springer.com/article/10.1007/BF01397083>

→ https://en.wikipedia.org/wiki/The_Art_of_Computer_Programming

Combine floating- and fixed-

point

- and INT32 for one variable
- Initialized with a large number beforehand
- Similar to 48-bit integer



Tetsu Narumi et al., International Journal of Computational Methods, vol. 8, No. 3, pp. 561-581 (2011)

(slide courtesy of Prof Tetsu Narumi, UEC, Tokyo)

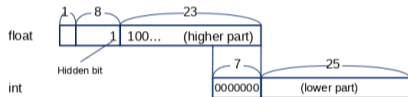
→ [doi:10.1142/S0219876211002708](https://doi.org/10.1142/S0219876211002708)

→ [doi:10.1016/j.cpc.2012.09.022](https://doi.org/10.1016/j.cpc.2012.09.022)

Combine floating- and fixed-

point

- and INT32 for one variable
- Initialized with a large number beforehand
- Similar to 48-bit integer



All elementary operations, +, -, *

Tetsu Narumi et al., International Journal of Computational Methods, vol. 8, No. 3, pp. 561-581 (2011)

(slide courtesy of Prof Tetsu Narumi, UEC, Tokyo)

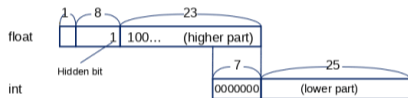
→ [doi:10.1142/S0219876211002708](https://doi.org/10.1142/S0219876211002708)

→ [doi:10.1016/j.cpc.2012.09.022](https://doi.org/10.1016/j.cpc.2012.09.022)

Combine floating- and fixed-

point

- and INT32 for one variable
- Initialized with a large number beforehand
- Similar to 48-bit integer



Similar strategy in AMBER's pmemd.cuda

All elementary operations, +, -, *

Tetsu Narumi et al., International Journal of Computational Methods, vol. 8, No. 3, pp. 561-581 (2011)

(slide courtesy of Prof Tetsu Narumi, UEC, Tokyo)

→ [doi:10.1142/S0219876211002708](https://doi.org/10.1142/S0219876211002708)

→ [doi:10.1016/j.cpc.2012.09.022](https://doi.org/10.1016/j.cpc.2012.09.022)

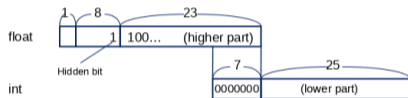
NUMERICAL ACCURACY & PERFORMANCE CONT.

CUDA SDK CONT.

Combine floating- and fixed-

point

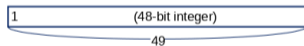
- and INT32 for one variable
- Initialized with a large number beforehand
- Similar to 48-bit integer



Also see Mandelbrot example in SDK, 5_Domain_Specific

All elementary operations, +, -, *

Similar strategy in AMBER's pmemd.cuda



Tetsu Narumi et al., International Journal of Computational Methods, vol. 8, No. 3, pp. 561-581 (2011)

(slide courtesy of Prof Tetsu Narumi, UEC, Tokyo)

→ [doi:10.1142/S0219876211002708](https://doi.org/10.1142/S0219876211002708)

→ [doi:10.1016/j.cpc.2012.09.022](https://doi.org/10.1016/j.cpc.2012.09.022)

5_DOMAIN_SPECIFIC/MARCHINGCUBES

CUDA SDK CONT.

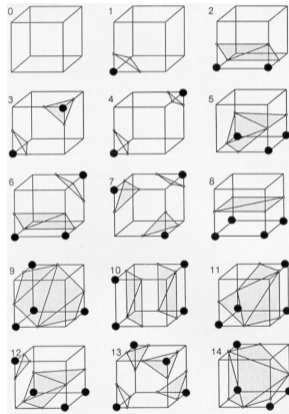
Practical example, of frequent use in biophysical chemistry

- Extracts an isosurface from a volumetric dataset using the 'marching cubes algorithm'
- OpenGL interoperation
- Chosen example is the electron density of C_{60}

→ `./run_marchingCubes.sh`

5_DOMAIN_SPECIFIC/MARCHINGCUBES CONT.

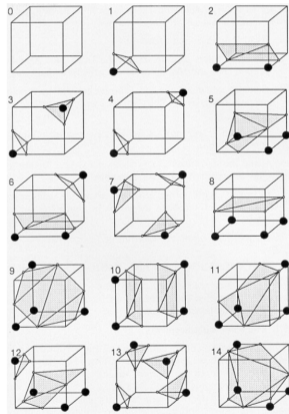
CUDA SDK CONT.



→ <https://dl.acm.org/citation.cfm?doid=37402.37422>

5_DOMAIN_SPECIFIC/MARCHINGCUBES CONT.

CUDA SDK CONT.



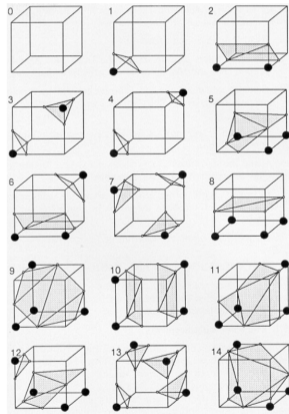
Assign vertices 0 if below iso-value, 1 if above

→ <https://dl.acm.org/citation.cfm?doid=37402.37422>

5_DOMAIN_SPECIFIC/MARCHINGCUBES CONT.

CUDA SDK CONT.

Neighbouring vertices of 0, 1 assignment determine iso-surface edges



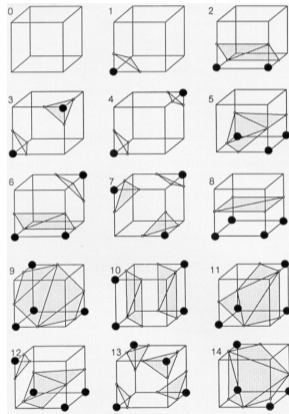
Assign vertices 0 if below iso-value, 1 if above

→ <https://dl.acm.org/citation.cfm?doid=37402.37422>

5_DOMAIN_SPECIFIC/MARCHINGCUBES CONT.

CUDA SDK CONT.

Neighbouring vertices of 0, 1 assignment determine iso-surface edges



14 elementary cases, the rest (256) from symmetry

Assign vertices 0 if below iso-value, 1 if above

→ <https://dl.acm.org/citation.cfm?doid=37402.37422>

5_DOMAIN_SPECIFIC/MARCHINGCUBES CONT.

CUDA SDK CONT.

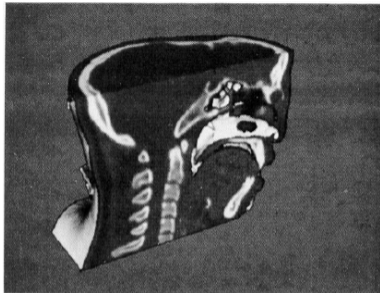


Figure 11. Sagittal Cut with Texture Mapping.

CT data, shows the slice data in relation to the constructed

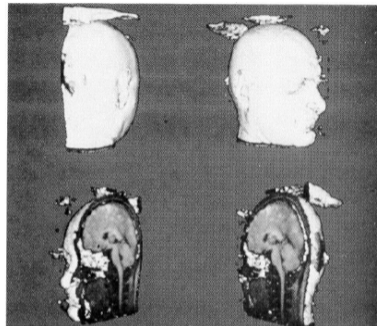


Figure 12. Rotated Sequence of Cut MR Brain.

Interesting for medical
images (CT,MR)

→ <https://dl.acm.org/citation.cfm?doid=37402.37422>

PROFILING CUDA CODE

CUDA SDK CONT.

The main profiling tool in CUDA nowadays is NVIDIA Nsight Compute

- nsys nvprof (quick and easy)
- nsys/nsys-ui (for system traces and timeline analysis – streams)
- ncu/ncu-ui (for detailed kernel optimization)

→ <https://docs.nvidia.com/nsight-compute/NsightComputeCli/index.html>

PROFILING CUDA CODE

CUDA SDK CONT.

The main profiling tool in CUDA nowadays is NVIDIA Nsight Compute

- nsys nvprof (quick and easy)
- nsys/nsys-ui (for system traces and timeline analysis – streams)
- ncu/ncu-ui (for detailed kernel optimization)

```
cuda-zen sh@n3073-009:~$ nvcc ./mmm_example_1.cu
cuda-zen sh@n3073-009:~$ ncu -f -set full -o profile ./a.out
cuda-zen sh@gui3068-010:~$ ncu-ui ./profile.ncu-rep
```

→ <https://docs.nvidia.com/nsight-compute/NsightComputeCli/index.html>

PROFILING CUDA CODE

CUDA SDK CONT.

The main profiling tool in CUDA nowadays is NVIDIA Nsight Compute

- nsys nvprof (quick and easy)
- nsys/nsys-ui (for system traces and timeline analysis – streams)
- ncu/ncu-ui (for detailed kernel optimization)

1. Compilation

```
cuda-zen sh@n3073-009:~$ nvcc ./mmm_example_1.cu  
cuda-zen sh@n3073-009:~$ ncu -f -set full -o profile ./a.out  
cuda-zen sh@gui3068-010:~$ ncu-ui ./profile.ncu-rep
```

→ <https://docs.nvidia.com/nsight-compute/NsightComputeCli/index.html>

PROFILING CUDA CODE

CUDA SDK CONT.

The main profiling tool in CUDA nowadays is NVIDIA Nsight Compute

- nsys nvprof (quick and easy)
- nsys/nsys-ui (for system traces and timeline analysis – streams)
- ncu/ncu-ui (for detailed kernel optimization)

1. Compilation

```
cuda-zen sh@n3073-009:~$ nvcc ./mmm_example_1.cu  
cuda-zen sh@n3073-009:~$ ncu -f -set full -o profile ./a.out  
cuda-zen sh@gui3068-010:~$ ncu-ui ./profile.ncu-rep
```

2. Write profile

→ <https://docs.nvidia.com/nsight-compute/NsightComputeCli/index.html>

PROFILING CUDA CODE

CUDA SDK CONT.

The main profiling tool in CUDA nowadays is NVIDIA Nsight Compute

- nsys nvprof (quick and easy)
- nsys/nsys-ui (for system traces and timeline analysis – streams)
- ncu/ncu-ui (for detailed kernel optimization)

1. Compilation

```
cuda-zen sh@n3073-009:~$ nvcc ./mmm_example_1.cu  
cuda-zen sh@n3073-009:~$ ncu -f -set full -o profile ./a.out  
cuda-zen sh@gui3068-010:~$ ncu-ui ./profile.ncu-rep
```

2. Write profile

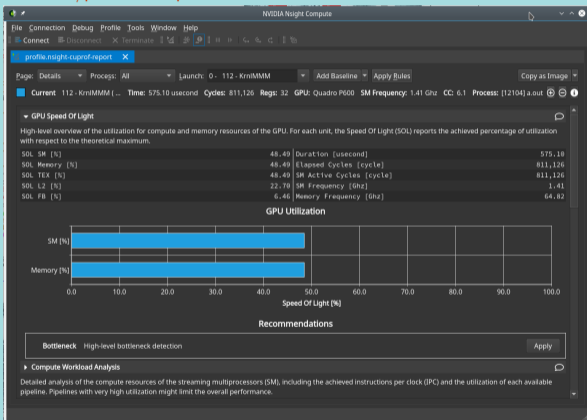
3. Visualize profile

→ <https://docs.nvidia.com/nsight-compute/NsightComputeCli/index.html>

PROFILING CUDA CODE CONT.

CUDA SDK CONT.

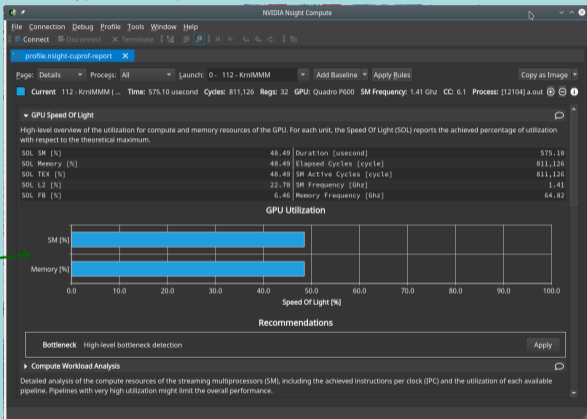
```
cuda-zen sh@gui3068-010:~$ ncu-ui ./profile.ncu-rep
```



PROFILING CUDA CODE CONT.

CUDA SDK CONT.

```
cuda-zen sh@gui3068-010:~$ ncu-ui ./profile.ncu-rep
```



GPU features evaluated w.r.t theoretical max (SOL)

PROFILING CUDA CODE CONT.

CUDA SDK CONT.

```
cuda-zen sh@gui3068-010:~$ ncu-ui ./profile.ncu-rep
```

The screenshot displays the NVIDIA Night Compute interface. At the top, the application title is "profile.ncu-rep". Below the title bar, there are menu options: File, Connection, Debug, Profile, Tools, Window, Help. A toolbar contains buttons for Connect, Disconnect, Terminate, and various navigation icons. The main content area shows the following information:

Page: Details | Process: All | Launch: 0 - 112 - KniMMM | Add Baseline | Apply Rules | Copy as Image

Current: 112 - KniMMM [...] Time: 575.10 usecond Cycles: 811,126 Regs: 32 GPU: Quadro P600 SM Frequency: 1.41 Ghz CC: 6.1 Process: [12104] a.out

Compute Workload Analysis

Detailed analysis of the compute resources of the streaming multiprocessors (SM), including the achieved instructions per clock (IPC) and the utilization of each available pipeline. Pipelines with very high utilization might limit the overall performance.

Executed Ipc Elapsed [1inst/cycle]	0.91	SM Busy [%]	48.49
Executed Ipc Active [1inst/cycle]	0.92	Issue Slots Busy [%]	15.37
Issued Ipc Active [1inst/cycle]	0.92	-	-

Memory Workload Analysis

Detailed analysis of the memory resources of the GPU. Memory can become a limiting factor for the overall kernel performance when fully utilizing the involved hardware units (Mem Busy), exhausting the available communication bandwidth between those units (Max Bandwidth), or by reaching the maximum throughput of issuing memory instructions (Mem Pipes Busy). Depending on the limiting factor, the memory chart and tables allow to identify the exact bottleneck in the memory system.

Memory Throughput [Gbyte/second]	66.96	Mem Busy [%]	48.49
L1 Hit Rate [%]	0	Max Bandwidth [%]	48.49
L2 Hit Rate [%]	19.87	Mem Pipes Busy [%]	24.57

Scheduler Statistics

Summary of the activity of the schedulers issuing instructions. Each scheduler maintains a pool of warps that it can issue instructions for. The upper bound of warps in the pool (Theoretical Warps) is limited by the launch configuration. On every cycle each scheduler checks the state of the allocated warps in the pool (Active Warps). Active warps that are not stalled (Eligible Warps) are ready to issue their next instruction. From the set of eligible warps the scheduler selects a single warp from which to issue one or more instructions (Issued Warp). On cycles with no eligible warps, the issue slot is skipped and no instruction is issued. Having many skipped issue slots indicates poor latency hiding.

Active Warps Per Scheduler [warp/cycle]	7.97	Instructions Per Active Issue Slot [1inst/issue]	1.38
Eligible Warps Per Scheduler [warp/cycle]	0.20	No Eligible [%]	83.60
Issued Warp Per Scheduler [1issue/cycle]	0.17	One or More Eligible [%]	16.86

Warp State Statistics

Analysis of the states in which all warps spent cycles during the kernel execution. The warp states describe a warp's readiness or inability to issue its next instruction. The warp cycles per instruction define the latency between two consecutive instructions. The higher the value, the more warp parallelism is required to hide this latency. For

PROFILING CUDA CODE CONT.

CUDA SDK CONT.

```
cuda-zen sh@gui3068-010:~$ ncu-ui ./profile.ncu-rep
```

profile.nsigth-cuprof-report

Page: Details Process: All Launch: 0 - 112 - KrmMMM Add Baseline Apply Rules Copy as Image

Current: 112 - KrmMMM [...] Time: 575.10 usecond Cycles: 811,126 Regs: 32 GPU: Quadro P600 SM Frequency: 1.41 Ghz CC: 6.1 Process: [12104] a.out

Compute Workload Analysis

Detailed analysis of the compute resources of the streaming multiprocessors (SM), including the achieved instructions per clock (IPC) and the utilization of each available pipeline. Pipelines with very high utilization might limit the overall performance.

Executed Ipc Elapsed [1inst/cycle]	0.91	SM Busy [%]	48.49
Executed Ipc Active [1inst/cycle]	0.92	Issue Slots Busy [%]	15.37
Issued Ipc Active [1inst/cycle]	0.92	-	-

Memory Workload Analysis

Detailed analysis of the memory resources of the GPU. Memory can become a limiting factor for the overall kernel performance when fully utilizing the involved hardware units (Mem Busy), exhausting the available communication bandwidth between those units (Max Bandwidth), or by reaching the maximum throughput of issuing memory instructions (Mem Pipes Busy). Depending on the limiting factor, the memory chart and tables allow to identify the exact bottleneck in the memory system.

Memory Throughput [Gbyte/second]	66.96	Mem Busy [%]	48.49
L1 Hit Rate [%]	0	Max Bandwidth [%]	48.49
L2 Hit Rate [%]	19.87	Mem Pipes Busy [%]	24.57

Scheduler Statistics

Summary of the activity of the schedulers issuing instructions. Each scheduler maintains a pool of warps that it can issue instructions for. The upper bound of warps in the pool (Theoretical Warps) is limited by the launch configuration. On every cycle each scheduler checks the state of the allocated warps in the pool (Active Warps). Active warps that are not stalled (Eligible Warps) are ready to issue their next instruction. From the set of eligible warps the scheduler selects a single warp from which to issue one or more instructions (Issued Warp). On cycles with no eligible warps, the issue slot is skipped and no instruction is issued. Having many skipped issue slots indicates poor latency hiding.

Active Warps Per Scheduler [warp/cycle]	7.97	Instructions Per Active Issue Slot [1inst/issue]	1.38
Eligible Warps Per Scheduler [warp/cycle]	0.20	No Eligible [%]	83.60
Issued Warp Per Scheduler [1issue/cycle]	0.17	One or More Eligible [%]	16.86

Warp State Statistics

Analysis of the states in which all warps spent cycles during the kernel execution. The warp states describe a warp's readiness or inability to issue its next instruction. The warp cycles per instruction define the latency between two consecutive instructions. The higher the value, the more warp parallelism is required to hide this latency. For

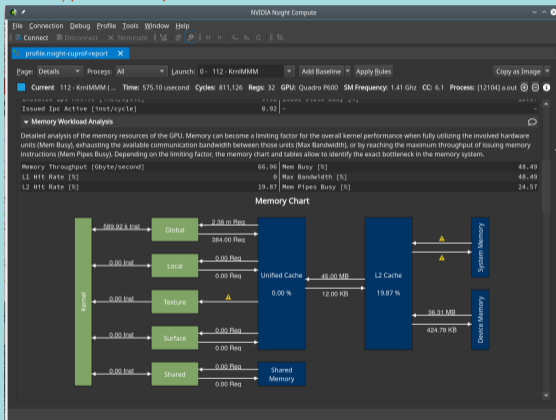
Different categories' evaluation

PROFILING CUDA CODE CONT.

CUDA SDK CONT.

```
cuda-zen sh@gui3068-010:~$
```

```
ncu-ui ./profile.ncu-rep
```

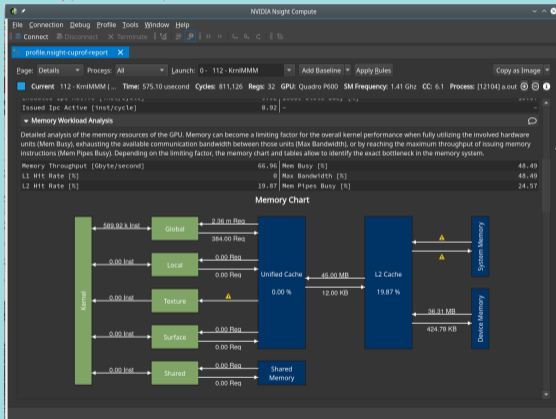


PROFILING CUDA CODE CONT.

CUDA SDK CONT.

```
cuda-zen sh@gui3068-010:~$
```

```
ncu-ui ./profile.ncu-rep
```



- The Heterogeneous Interface for Portability (HIP) is AMD's dedicated GPU programming environment (e.g. to program MI250X devices)
- Very similar to CUDA
- HIP code will also run on NVIDIA platforms
- HIP forms another ecosystem with tools, libraries, etc
- Every basic CUDA construct has a direct counterpart in HIP

→ <https://developer.amd.com/wp-content/resources/ROCM%20Learning%20Centre/chapter3/HIP-Coding-3.pdf>

Vector Addition with HIP

```
// HIP kernel. Each thread takes care of one element of c;
__global__ void vecAdd(double *a, double *b, double *c, int n)
{
    // Get our global thread ID
    int id = blockIdx.x * blockDim.x + threadIdx.x;
    // Make sure we do not go out of bounds
    if (id < n)
        c[id] = a[id] + b[id];
}

int main()
{
    ...
    // Execute the kernel with n threads
    hipLaunchKernelGGL(vecAdd, dim3(gridSize), dim3(blockSize), 0, 0, d_a, d_b, d_c, n);
    hipDeviceSynchronize();
}
```

→ <https://developer.amd.com/wp-content/resources/ROCm%20Learning%20Centre/chapter3/HIP-Coding-3.pdf>

HIP CONT.

HIP EXAMPLE

Vector Addition with HIP

```
// HIP kernel. Each thread takes care of one element of c;
__global__ void vecAdd(double *a, double *b, double *c, int n)
{
    // Get our global thread ID
    int id = blockIdx.x * blockDim.x + threadIdx.x;
    // Make sure we do not go out of bounds
    if (id < n)
        c[id] = a[id] + b[id];
}

int main()
{
    ...
    // Execute the kernel with n threads
    hipLaunchKernelGGL(vecAdd, dim3(gridSize), dim3(blockSize), 0, 0, d_a, d_b, d_c, n);
    hipDeviceSynchronize();
}
```

same kernel
declaration
specifiers

→ <https://developer.amd.com/wp-content/resources/ROCm%20Learning%20Centre/chapter3/HIP-Coding-3.pdf>

HIP CONT.

HIP EXAMPLE

Vector Addition with HIP

```
// HIP kernel. Each thread takes care of one element of c;
__global__ void vecAdd(double *a, double *b, double *c, int n)
{
    // Get our global thread ID
    int id = blockIdx.x * blockDim.x + threadIdx.x;
    // Make sure we do not go out of bounds
    if (id < n)
        c[id] = a[id] + b[id];
}

int main()
{
    ...
    // Execute the kernel with n threads
    hipLaunchKernelGGL(vecAdd, dim3(gridSize), dim3(blockSize), 0, 0, d_a, d_b, d_c, n);
    hipDeviceSynchronize();
}
```

same kernel
declaration
specifiers

same built-in
variables,
e.g. `threadIdx.x=0,1,2...`

→ <https://developer.amd.com/wp-content/resources/ROCm%20Learning%20Centre/chapter3/HIP-Coding-3.pdf>

HIP CONT.

HIP EXAMPLE

Vector Addition with HIP

```
// HIP kernel. Each thread takes care of one element of c;
__global__ void vecAdd(double *a, double *b, double *c, int n)
{
    // Get our global thread ID
    int id = blockIdx.x * blockDim.x + threadIdx.x;
    // Make sure we do not go out of bounds
    if (id < n)
        c[id] = a[id] + b[id];
}

int main()
{
    ...
    // Execute the kernel with n threads
    hipLaunchKernelGGL(vecAdd, dim3(gridSize), dim3(blockSize), 0, 0, d_a, d_b, d_c, n);
    hipDeviceSynchronize();
}
```

same kernel
declaration
specifiers

same built-in
variables,
e.g. `threadIdx.x=0,1,2...`

different ker-
nel execution
configuration

→ <https://developer.amd.com/wp-content/resources/ROCm%20Learning%20Centre/chapter3/HIP-Coding-3.pdf>

HIP CONT.

COMPILE AND RUN AND MONITOR

```
(zen3) [sh@some ~ mi250x ~]$ hipcc vadd_hip.cpp
(zen3) [sh@some ~ mi250x ~]$ ./a.out

0 100.000000
1 100.000000
2 100.000000
3 100.000000
4 100.000000
5 100.000000
6 100.000000
7 100.000000
8 100.000000
9 100.000000
10 100.000000
11 100.000000
...
99 100.000000
```

- <https://developer.amd.com/wp-content/resources/ROCm%20Learning%20Centre/chapter3/HIP-Coding-3.pdf>
- https://rocm-docs.amd.com/en/latest/Programming_Guides/HIP-porting-guide.html
- <https://forums.extremehw.net/topic/782-solvedimage-upload-broken-on-firefox>

HIP CONT.

COMPILE AND RUN AND MONITOR

```
(zen3) [sh@some ~ mi250x ~]$ hipcc vadd_hip.cpp  
(zen3) [sh@some ~ mi250x ~]$ ./a.out
```

```
0 100.000000  
1 100.000000  
2 100.000000  
3 100.000000  
4 100.000000  
5 100.000000  
6 100.000000  
7 100.000000  
8 100.000000  
9 100.000000  
10 100.000000  
11 100.000000  
...  
99 100.000000
```

```
tictoc@TickTockArch $ rocm-smi
```

```
=====ROCM System Management Interface=====
```

GPU	Temp	AvgPwr	SCLK	MCLK	Fan	Perf	PwrCap	VRAM%	GPU%
0	42.0c	294.0W	1925Mhz	1000Mhz	0.0%	high	450.0W	17%	100%
1	50.0c	219.0W	1772Mhz	1000Mhz	0.0%	high	220.0W	4%	100%
2	45.0c	273.0W	1925Mhz	1000Mhz	0.0%	high	300.0W	4%	100%
3	55.0c	163.0W	1700Mhz	1000Mhz	0.0%	high	170.0W	4%	100%

```
=====End of ROCm SMI Log =====
```

```
tictoc@TickTockArch $ _
```

- <https://developer.amd.com/wp-content/resources/ROCM%20Learning%20Centre/chapter3/HIP-Coding-3.pdf>
- https://rocm-docs.amd.com/en/latest/Programming_Guides/HIP-porting-guide.html
- <https://forums.extremehw.net/topic/782-solvedimage-upload-broken-on-firefox>

TAKE HOME MESSAGES

- CUDA can be easily extended to Fortran programs

TAKE HOME MESSAGES

- CUDA can be easily extended to Fortran programs
- BLAS calls may be substituted with calls to CUBLAS potentially offering significant acceleration

TAKE HOME MESSAGES

- CUDA can be easily extended to Fortran programs
- BLAS calls may be substituted with calls to CUBLAS potentially offering significant acceleration
- Basic LAPACK calls may be replaced with CUSOLVER routines

TAKE HOME MESSAGES

- CUDA can be easily extended to Fortran programs
- BLAS calls may be substituted with calls to CUBLAS potentially offering significant acceleration
- Basic LAPACK calls may be replaced with CUSOLVER routines
- Sparse problems are supported with specific libraries, other basic tasks as well, CUFFT, CURAND, CUDNN

TAKE HOME MESSAGES

- CUDA can be easily extended to Fortran programs
- BLAS calls may be substituted with calls to CUBLAS potentially offering significant acceleration
- Basic LAPACK calls may be replaced with CUSOLVER routines
- Sparse problems are supported with specific libraries, other basic tasks as well, CUFFT, CURAND, CUDNN
- Taking advantage of fast FP32 operations with quasi-FP64 accuracy on consumer grade cards requires a couple of algorithmic changes

TAKE HOME MESSAGES

- CUDA can be easily extended to Fortran programs
- BLAS calls may be substituted with calls to CUBLAS potentially offering significant acceleration
- Basic LAPACK calls may be replaced with CUSOLVER routines
- Sparse problems are supported with specific libraries, other basic tasks as well, CUFFT, CURAND, CUDNN
- Taking advantage of fast FP32 operations with quasi-FP64 accuracy on consumer grade cards requires a couple of algorithmic changes
- NVIDIA's NSIGHT compute is an advanced profiling tool for in-depth kernel optimization

TAKE HOME MESSAGES

- CUDA can be easily extended to Fortran programs
- BLAS calls may be substituted with calls to CUBLAS potentially offering significant acceleration
- Basic LAPACK calls may be replaced with CUSOLVER routines
- Sparse problems are supported with specific libraries, other basic tasks as well, CUFFT, CURAND, CUDNN
- Taking advantage of fast FP32 operations with quasi-FP64 accuracy on consumer grade cards requires a couple of algorithmic changes
- NVIDIA's NSIGHT compute is an advanced profiling tool for in-depth kernel optimization
- AMD's HIP is very similar to CUDA and kernel code can be directly interchanged